# Prior Information and Generalized Questions

**Jörg C. Lemm**

This publication can be retrieved by anonymous ftp to publications.ai.mit.edu.
The pathname for this publication is: ai-publications/1500-1999/AIM-1598.ps.Z

## Abstract

In learning problems available information is usually divided into two categories: examples of function values (or training data) and prior information (e.g. a smoothness constraint).

This paper 1.) studies aspects on which these two categories usually differ, like their relevance for generalization and their role in the loss function, 2.) presents a unifying formalism, where both types of information are identified with answers to generalized questions, 3.) shows what kind of generalized information is necessary to enable learning, 4.) aims to put usual training data and prior information on a more equal footing by discussing possibilities and variants of measurement and control for generalized questions, including the examples of smoothness and symmetries, 5.) reviews shortly the measurement of linguistic concepts based on fuzzy priors, and principles to combine preprocessors, 6.) uses a Bayesian decision theoretic framework, contrasting parallel and inverse decision problems, 7.) proposes, for problems with non–approximation aspects, a Bayesian two step approximation consisting of posterior maximization and a subsequent risk minimization, 8.) analyses empirical risk minimization under the aspect of non-local information 9.) compares the Bayesian two step approximation with empirical risk minimization, including their interpretations of Occam's razor, 10.) formulates examples of stationarity conditions for the maximum posterior approximation with nonlocal and nonconvex priors, leading to inhomogeneous nonlinear equations, similar for example to equations in scattering theory in physics.

In summary, this paper focuses on the dependencies between answers to different questions. Because not training examples alone but such dependencies enable generalization, it emphasizes the need of their empirical measurement and control and of a more explicit treatment in theory.

# 1 Introduction

To clarify the aim of the paper, our use of the term "prior information" and to define vocabulary and notation we analyze the fundamental problem of generalization for a simple toy example:

Let us assume that we are interested in the answers $y_1(x_1)$ and $y_2(x_2)$ to two questions $x_1$ and $x_2$. We will call these questions *relevant* to our *application*. Let us assume we have found for $y_1$ the value $y_1 = 2$ in a first measurement. We are interested in the future output when measuring $x_1$ and $x_2$. Clearly this is not possible without further assumptions. These assumptions are used in choosing a set $F^0$ of 'possible states of nature' $f^0$. If we assume that the problem is deterministic, then the problem is already solved for $x_1$. We call this assumption *local* as they refer to a single question. If the problem is probabilistic and we want to determine the probability of future answers $y_1$ to $x_1$ this requires the assumption of a set of stationary distributions $p(y_1|x_1, f^0)$ corresponding to possible states $f^0$ for $y_1$ so we can use the data to reweight the probabilities $p(f^0)$. The collection of $p(f^0)$ is what we call a 'state of knowledge' and can be used to predict future outcomes $y_1$.

But in most real world problems there are new questions $x_2$ for which no training examples are available. This is the *generalization* problem. It seems even harder than the learning problem for $x_1$ for which at least local data are available. To remedy the situation we clearly need answers which depend on $y_2$ at least indirectly. We may directly use a *nonlocal* assumption like $y_2 = y + y_1$. This is trivial if we know $y$. But this may not be the case in the beginning and this paper analyses how information about $y$ may be obtained. Therefore, we want to start with independent relevant outcome $y_1$ and $y_2$ (what we will call a factorial state) and try to relate all dependencies to additional information. That means we see $y$ as an answer to a question $q$ with $y(q) = y_2 - y_1$. This *nonlocal q* question represents a measurement device for differences. We call questions depending on more than one $y_i$ *generalized* questions. Examples are nonlocal questions like $y_2 - y_1$ or questions referring to repeated measurements of the same questions like $y_1 + y_1'$. We can now separate the nonlocal information in two parts: the structural information, that is the definition of the nonlocal question $q$ and an empirical or controllable part, that is the result of measuring $q$ or controlling its answer. The definition of $q$ alone does not yet determine the value or probability of $y_2$. Only an actual measurement of $q$ relates those two questions.

This is a conceptually clear starting point to analyze generalization. We start with an unrelated set of relevant questions $x_i$ in which we are interested and formulate all nonlocal information as answers to nonlocal questions. How can we know what $q$ measures? Clearly the definition of $q$, that is the structural information, will contain some assumptions, specifically some kind of stationarity. The definition of $q$ is usually also based on previous empirical information. For example, we could have tested a difference device many times before using it and found it working correctly (or at least indicating correctly its failure). Then the stationarity assumption that it will work correctly also for our task, seems reasonable. (In a Bayesian approach we can give it a high subjective probability.) Thus, structural information is based on *transfer* of knowledge between tasks. Transfer means generalization with respect to a task variable. In statistics a nonlocal question corresponds for example to the choice of a smoothness functional to be used as prior. Likewise, *rules* in an expert system with dependencies modeled by *logic, fuzzy logic* or *Bayesian belief networks* and *macroscopic variables* (e.g. energy) in physics are special cases of nonlocal questions.

Data are pairs of question and answer and we refer to the data related to relevant questions $x_1$, $x_2$ as *test data* and to data related to questions $x_1$, $q$ with answers available as *training data*. Then for independent test questions the training data necessarily have to contain answers to nonlocal questions to allow generalization. It is from this point of view that we want to analyze prior information: we start with independent relevant questions and try to give explicit explanations, namely measurable prior information, for their dependencies. We summarize three aspects we have been mentioning

1. Generalization aspect: not all relevant questions are available for training,

2. Application aspect: not all available data correspond directly to relevant questions,

3. Transfer aspect: Some aspects of the setting, like the nature of the measuring devices, have to be known.

Now let us switch to more realistic statistical problems and consider a given function $f : y = f(x)$ with $x \in X$ and $y \in Y$ for some set $X$ and some set $Y$. We allow as response to $x$ instead of a deterministic $f(x)$ a random variable which takes values $y$ with probability $p(y|x, f)$. Let us consider the following two kinds of information: A typical training example, that is a pair $(x, y)$, and some smoothness constraint. Let the latter be given by a bound on a smoothness functional, like, for example, the discretized version $\sum_i (\frac{E(f(x_i)) - E(f(x_i + \Delta x_i))}{\Delta x_i})^2 < \theta$, where $E(f(x))$ denotes the regression function or expectation of $y = f(x)$ at point $x$.

The first one can clearly be interpreted as a pair of question and answer. We see $x = q_x$ as a question for state $f$ about its value at $x$ and denote the answer by $y = q_x(f) = f(x)$ (i.e. $f(x)$ denotes a random variable), generated by a probability distribution $p(y|x, f)$. We write such a pair as $(x, y) = (q_x, y_{q_x})$. Information about smoothness can be seen as answer to a more generalized 'smoothness' question $q_{s,\theta} = \text{sign}\left(\sum_i (\frac{E(f(x_i)) - E(f(x_i + \Delta x_i))}{\Delta x_i})^2 - \theta\right)$. We write this pair as $(q_{s,\theta}, y_{s,\theta})$ with $q_{s,\theta}(f)$ equal to $\pm 1$ and remark that replacing the fixed $\theta$ by a random variable $\Theta$ with mean $\theta$ is one easy way to generalize to the case of a probabilistic answer. Now note that in contrast to a simple example data question the smoothness question depends on more than one $x$−value of the function $f$. In this sense the smoothness question may be seen as nonlocal. Nonlocality is also present in symmetry questions. e.g. like $q_{p,\theta} = \text{sign}\left(\sum_i (E(f(x_i)) - E(f(-x_i)))^2 - \theta\right)$ or more

general $q_{s,\theta} = \text{sign}\left(\sum_i (E(f(x_i)) - E(f(sx_i)))^2 - \theta\right)$ where $s$ represents any symmetry operator under which invariance of $f$ is tested. Also bounds for the nonlocal maximum function $\max_x E(f(x))$ are widely used. Those three types of bounds, on smoothness, symmetry, and maximum functionals, are the the most frequently used forms of prior information in practice.

There is one practical distinction between nonlocal prior information like smoothness and standard training examples: the latter are assumed to be empirically measurable while the 'measurability' of nonlocal data like smoothness is often quite unclear. The difficulty lies in the fact that these nonlocal questions depend on the whole function $f(x)$, and for continuous $x$ there is an infinite number of function values. But, if we want to treat nonlocal information in a similar way as standard training examples we also have to discuss the measurement of such data. This paper shows that, besides cases where highly parallel measurement devices are available, also restrictions of measurement devices and the definition of situations of interest and their control can be interpreted as measurements of such questions.

There is another group of nonlocal questions, for which data are more easily available in practice. These are nonlocal questions depending on only a finite (and in practice 'small enough') number of $f(x_i)$. Empirical answers to these questions can be found by measuring standard training examples $f(x)$ and applying some operations to the results which define the questions. Examples include differences $f(x_i) - f(x_j)$, weighted averages $\sum_i w_i f(x_i)$ like discrete wavelet components, symmetries $s$ for specific points $f(x_i) - f(sx_i)$, and measurements of $f(x)$ with input noise. When the expectation $E(f(x_i))$ can be approximated by $1/n \sum_i f(x_i)$, that is by using repeated measurements, we can, at least approximately, also measure examples depending on $E(f(x_i))$. Sometimes only nonlocal data are available, as in models where the variables of interest are hidden but related by structural assumptions to several observable variables, like for example in hidden Markov models. However, even when the variables of interest are observable there can be measurement devices which measure finite nonlocal questions with greater accuracy than by using the measuring devices for $f(x)$. Specifically, differences are often measured better directly than by first looking for $f(x_i)$ then for $f(x_j)$ and then taking the difference. This is always the case if part of the measurement error appears in a final step common to all measuring devices, like an output scale with fixed finite resolution or memory errors. We also remarked that the definition of a generalized question is called a rule in expert systems. Thus, available rules collected from experts may be incorporated as nonlocal question. Also, many examples of measuring nonlocal questions can be found in physics. The energy function of a macroscopic system is highly nonlocal with respect to the microscopic constituents. It seems that nonlocal data are often available, however in many cases their practical combination with priors of infinite nonlocality is probably hindered due to the difficulties of the Frequentist approach in statistics in dealing with nonlocal questions. That means that for data with infinite nonlocality the measurement problem represents an empirical difficulty, while including complex forms of nonlocal data, even with finite nonlocality, yields easily to mathematical difficulties related to solving nonlinear equations.

The paper formulates a theoretical approach in which answers to generalized questions are treated similar to standard training examples. In Section 2 the theoretical framework is presented, paying particular attention to the generalization problem. It is clarified that not the (standard training) data, but the nonlocal (prior) data enable learning.

Section 3 presents generalized questions in more detail. The importance of nonlocal information for learning makes it necessary to discuss the fundamental measurement problem for nonlocal questions. For example, restrictions of measurement devices correspond to nonlocal measurements.

Section 4 relates the classical approximate symmetry priors (e.g. smoothness) to nonideal measurement devices, the specific restriction being input noise.

Usually learning problems are defined and controlled by man. Consequently, nonlocal dependencies may be caused by internal human concepts. In addition experts may contribute knowledge in form of verbal statements. Section 5 discusses principal variants of including subjective priors (which is meant to include priors caused by subjects), like an interface with fuzzy priors. The Section does not aim to give an overview over this active and growing area of research, but presents some principles and aims in explaining the origin of nonlinearities which appear as technical difficulties in Section 9.

Section 6 extends the theoretical framework of Section 2 to decision problems using the language of generalized questions. Parallel and inverse settings are contrasted.

Section 7 contains a discussion of the Bayesian approach, concentrating an the saddle point approximation, leading to the two step MaP–MiR procedure.

Section 8 discusses the Frequentist approach of empirical risk minimization from the Bayesian decision theoretic point of view. Especially the relation to the maximum posterior approximation in approximation and non–approximation situations is discussed in detail.

Section 9 shows the possible use of generalized information by discussing maximum posterior variational (mean field) equations for several variants of nonlinear regularization procedures. In general those have the form of inhomogeneous integro–differential equations, like they appear for example similarly in a time independent formulation of quantum mechanical scattering theory.

Finally, Section 10 exemplifies the ideas of nonlinear regularization on a numerical example.

## 2 A constructive Bayesian approach

### 2.1 Model and vocabulary

In this section we choose a Bayesian approach (See for example Berger, 1985; Haussler, 1995; Bishop, 1995a, Wolpert, 1996a) and attempt to separate in a clear manner the local from the nonlocal parts in our model.

Therefore, we begin with the construction of the local components: We enumerate a set of basis questions $x$ necessary to determine the possible states we are interested in. The term *constructive* means that we do not start with a model given by some chosen parameterization (see Wolpert, 1994a), however, we use the relevant questions we are interested in to construct the related model. For every question we give a set of possible answers together with a set of question specific answer distributions $p(y|x, f_x^0)$. Each such answer distribution represents one possible pure local state $f_x^0$. Its index $x$ allows the local states to be defined independently and individually for every basis question $x$. In this paper a pure state is indicated by an superscript 0. We assume the nature to be in exactly one pure state. [1] Which one, is usually unknown to us. Thus, in contrast an $f_x$ without the superscript 0 denotes not a real pure state but a state of knowledge being equivalent to an assignment of a probability $p(f_x^0|f_x)$ to every pure state, where we often skip $f_x$ in the notation. In a learning situation we assume the actual, usually unknown, pure state (of nature) to be constant[2] for the time under study, while the state of knowledge, which reflects the learning process, changes with our information.

The independently defined local components are related by nonlocal states. A pure nonlocal state $f^0$ is defined through an assignment of one local state $f_x^0$ to every basis question $x$ and a nonlocal state of knowledge $f$ through an assignment of a probability $p(f^0)$ to every pure nonlocal state, with $p(f^0)$ implicitly understood as $p(f^0|f)$. The probabilities $p(f^0)$ contain all the relations between local components and therefore the nonlocal information.

We now define the ingredients of the theoretical approach more formally:

1. A *local basic model* consisting of

    a. a (finite[3]) set $X$ of *basis questions* $x$,

    b. a (measurable[4]) space $Y_x$ of possible answers $y_x$ for every question[5],

    c. a (finite[6]) set of probability distributions (local elementary functions) over answers $F_x^0 = \{p(y|x, f_x^0)\}$ for every question, the *pure local states*,

    d. a set $F_x$ of local states of knowledge $f_x$ defined by assigning a probability $p(f_x^0)$ to every local pure state $f_x^0$.

A local state $f_x^0$ can be considered indexing a vector $p(y|x, f_x^0)$ in the linear space $\mathcal{F}(Y, x)$ of functions defined on $Y = Y_x$, i.e. $f_x^0 \in \mathcal{F}(Y, x)$, with norm $|f_x^0| = \sum_y p(y|x, f_x^0) = \sum_y p(y|x, f^0) = 1$. A state of knowledge $f_x$ is in the convex hull of $F_x^0$, with $\sum_{f_x^0} p(f_x^0) = 1$. We denote by $\mathcal{C}(V)$ the convex hull of a set $V$ of vectors $v_i$ generated by linear combinations $\sum_i a_i v_i$ with coefficients fulfilling $\sum_i a_i = 1$ (e.g. $i = f_x^0$, $a_i = p(f_x^0|f_x)$), and by $\mathcal{L}$ the linear span, generated by unrestricted linear combinations. So we have

$$F_x^0 \subseteq F_x = \mathcal{C}(F_x^0) \subset \mathcal{L}(F_x^0) \subseteq \mathcal{F}(Y, x).$$

Now we construct the nonlocal parts.

2. *Pure nonlocal states* $f^0$, or, more shortly, pure states, are defined by a set of pairs $(p(y|x, f_x^0), x)$ containing exactly one local state for each question: $\{(p(y|x, f_x^0), x) : f_x^0 \in F_x^0, x \in X\}$ to give $p(y|x, f^0) = p(y|x, f_x^0)$. They form the set $F^0$ of pure states or elementary functions, which is a subset of the linear space of functions defined on $Y_x, X$, $f^0 \in \mathcal{F}(Y_x, X)$. We implicitly understand states (functions) as equivalence classes defined with respect to $X$ by $[f^0]_X = [f'^0]_X \Leftrightarrow \forall x \in X, \forall y \in Y_x : p(y|x, f^0) = p(y|x, f'^0)$.

3. A *nonlocal state of knowledge* $f$, or, more shortly, a state of knowledge is defined through assigning probabilities $p(f^0|f)$, or often shortly $p(f^0)$, to the pure states $f^0 \in F^0$ of the model. Thus, $f$ denote convex combinations of pure states $f^0$, i.e. $p(y|x, f) = \sum_{f^0 \in F^0} p(f^0) p(y|x, f^0)$, with $\sum_{f^0 \in F^0} p(f^0) = 1,.$ The set $F$ of possible states of knowledge (also called mixture states) $f$ of the model form the convex hull $F = \mathcal{C}(F^0)$ of $F^0 \subset \bigotimes_x \mathcal{L}(F_x^0)$. By construction we have $p(f^0|x) = p(f^0)$ stating independence between states $f^0$ and basic questions $x$. A state of knowledge is called *factorial* with respect to $X$ iff $p(f^0) = \prod_x p(f_x^0); \forall f^0 \in F^0$. It is fully determined by its local probabilities $p(f_x^0)$ assigned to each local state for question $x$.[7]

This paper discusses predefined dependencies between answers to different questions. The role of such *structural information* is discussed in section 2.3. Structural information is implemented by defining generalized questions.

---

[1] Thus, a pure state is maximally specified. We allow here a pure state to be probabilistic, so cases are possible where one pure state has the same probability distribution as a mixture of other pure states. In contrast to such an equivalent mixture a pure state is a fixed point under arbitrary learning.

[2] Constant is meant relative to a chosen reference system, which might also, for example, be varying in time.

[3] For infinite $X$ the corresponding functional integration represents a stochastic process in the language of mathematics or a field theory in the language of physics and must be consistently defined for every finite subset. Despite tehe existence of interacting field theories in physics, only functional integrals with Gaussian measures are mathematically well defined objects (Gardiner, 1990; van Kampen, 1992; Schervish, 1995). Non–Gaussian functional integrals have to be be defined for example by perturbation theory using the Feynman–Kac formula or by discretization, i.e. an ultraviolet cutoff (Glimm–Jaffe, 1987, Zinn–Justin, 1989, Bialek, Callan, Strong, 1996, Balasubramanian, 1996). For the use of Gaussian processes especially in Bayesian statistics see e.g. Williams, Rasmussen, 1996; Barber, Williams, 1997; Neal, 1997.

[4] In the mathematical sense.

[5] The set $Y_x$ can be assumed $x$–independent without loss of generality, so we will usually write $Y_x = Y$.

[6] The case of an infinite space $F_x^0$ requires the definition of a measure $df_x^0$.

[7] $\sum_{f_x^0 \in F_x^0} p(f_x^0) = 1, \forall x \in X$ ensures $\sum_{f^0 \in F^0} p(f^0) = 1$.

3

4. A *generalized question* is any probabilistic functional $q(f)$ of $f$. A probabilistic functional $q$ is a probabilistic mapping from the space $F$ of functions $f$ to a space $y^q$ of answers $y = q(f)$ given by probability distributions $p(y|q, f) = \sum_{f^0 \in F^0} p(f^0)p(y|q, f^0)$ defined by $p(y|q, f^0)$. We require the $f^0$–dependency of $p(y|q, f^0)$ to be expressible in terms of only the $p(y|x, f^0)$.[8][9][10] We define some set $Q$ of generalized questions to be *measurable* (or observable) by assuming, for every question in $Q$, the existence of a measuring device producing the corresponding answer. As consistency condition we require any question depending on a finite number of answers to measurable questions also to be measurable. We call the set of $x \in X$ on which the $q \in Q$ depend the *basis*[11] of $Q$ and denote it by $X^Q$. Skipping the unobservable elements $x \notin X^Q$ we write $X = X^Q$. Data $D$ are pairs of questions and corresponding answers $(q, y^q)$. In section 3.1 we will show how to write generalized questions explicitly, and section 3.2 discusses the measurement processes.

For applications we use a decision theoretic framework (see also Section 6) and separate the two subsets

5. $Q^D \subseteq Q$ of *available* or *training questions* for which answers are available.[12] The basis of $Q^D$ will be denoted by $X^{Q^D} = X^D \subseteq X$. $D_{Q^D}$ denotes the set of all possible data which only depend on questions $q \in Q^D$. We denote by $q^D$, and analogous for other sets of questions, vectors with components $q^D_i \in Q^D$. Here, $q^D_i = q^D_j \in Q^D, i \neq j$ is also allowed.[13] By $q^D \in Q^D$ we mean that $q^D_i \in Q^D$ for every component of the vector. If a vector appears in $Q^D$

we understand not the vector but its components to be elements. An data vector $D$ is a question–answer pair $(q^D, y^D)$. We may decompose $D$ into lower dimensional vectors $D_i = (q^D_i, y^D_i) \subseteq D_{Q^D}$, $i \in I$.

6. $Q^l \subseteq Q$ of *relevant* or *test questions* defining the possible application situations. The basis of $Q^l$ is denoted by $X^{Q^l} = X^l \subseteq X$, $D_{Q^l}$ the set of all possible data depending on $Q^l$, and we use $D^l$ to denote a test data vector[14] i.e. a set of pairs $(q^l, y_{q^l}) \in D_{Q^l}$.

The set $Q^l$ is called relevant, because

7. we assume a given *loss function* $l(q^l, y_{q^l}, z)$, defined for all $q^l \in Q^l$, i.e. application situations. The loss function depends on the question $q^l$, the answer $y_{q^l}$ and potentially also on additional variables $z \in Z$. We will call $z$ the *action variables* as we usually allow them to be controlled.

We include the possibility to control the action variables $z \in Z$ within the loss function by an active choice of an action state $\hat{f} \in \hat{F}$. Thus, we define

8. a family $\hat{f} \in \hat{F}$ of possible *action states* producing the probabilistic action $z \in Z$ according to $p(z|q^l, y, \hat{f})$ for $q^l \in Q^l$.

We write $\tilde{l}(q^l, y, \hat{f})$ for the effective loss if the $z$ variables in $l(q^l, y, z)$ can be integrated out.

The possible application situations, i.e. the relevant or test questions $q^l$, are generated by

9. a *test question* $q^l$–producing device $p(q|y_c, z_c)$ which, conditioned on a subset $c$ of 'past' values of $y_c \in Y$ and $z_c \in Z$, does not depend on $f^0$, $f$ or $\hat{f}$.

The probability distributions $p(y|q, f^0)$, $p(q|y_c, z_c)$, and $p(z|q^l, y, \hat{f})$ define a $\hat{f}$–dependent loss distribution $p(l|f, \hat{f})$:[15]

$$ p(l|f, \hat{f}) = \int dq \int dy \int dz \, p(q|y_c, z_c) $$

$$ \times p(y|q, f)p(z|q, y, \hat{f})\delta(l(q, y, z) - l). $$

The normative component is represented by the requirement to minimize

10. a *risk functional* $r[p(l|f, \hat{f})]$, which is a mapping from the loss distribution $p(l|f, \hat{f})$ into a subset of the real numbers, bounded from below. A common risk functional is the expectation of $l$ or *expected risk* $r(f, \hat{f}) = \int dl \, p(l|f, \hat{f})l$. [16]

New available data $D$ require updating of an initial state of knowledge $f^I$ to obtain a new $f$.

---

[8]Aside from the possibility that a functional is not defined for a specific $f^0$, there exists the possibility of a functional having only a certain probability of being not defined if applied to a function $f^0$, like $q(f^0) = 1/y$ if $p(y = 0|x, f^0) \neq 0$. Formally one can add the value 'undefined' to the space $Y$, as it is common in programming.

[9]One could allow a dependence of the definition of $q$, i.e. the $p(y|q, f, f^0)$, from the state of knowledge $f$ (which is known, as the name indicates), i.e. from the $p(f^0|f)$. However, this is only important when studying the dependence from $f$. A dependence $p(f^0|q)$ would mean that selecting $q$ already changes the hidden variable $f^0$. By definition of the model we have $p(y|x, f, f^0) = p(y|x, f^0)$ and $p(f^0) = p(f^0|x)$ for the basic questions.

[10]Compare for example with similar concepts in Ratsaby & Maiorov, 1996 and Smola & Schölkopf, 1997.

[11]This is not a linear basis of a vector space (Compare Jeffrey, 1968).

[12]This is related to learning (see 11., below). The set $Q^D$ depends on what part of information is considered part of the prior A prior $f^I$ can be expressed depending on some data $D^0$ and another prior $f^{I'}$ according to $p(f^0|D) = p(f^0|f(D, f^I))$ $= p(f^0|f(D, f^I(D^0, f^{I'}))) = p(f^0|D, D^0)$. Note that here the sloppy notation $p(f^0|D) = p(f^0|D, D^0)$ refers to different priors on the right and left hand side which are not explicitly indicated.

[13]In as far as no linear structure is needed $q^D$ can also be seen as a set, with dummy indices for repeated elements.

[14]Not to be confused with the validation set of empirical test data like it is used in cross–validation (see Section 8.4).

[15]See Section 3 and Section 6 for details and justification of the chosen components and their probability distributions.

[16]Minimizing an expected loss was already proposed by (Laplace, 1810ab) and later revived by (Wald, 1938) (See the historical remarks in (Le Cam, Yang, 1990)). A general formalization can be found in (Le Cam, 1986).

11. A *learning* model is a mapping $f = f(D, f^I)$ from $F$ to $F$ parameterized by $D$. The initial state will be called *prior* state. The Bayesian learning model is defined by

$$p(f^0 | f(D, f^I)) = \frac{p(y^D | q^D, f^0) p(f^0 | f^I)}{p(y^D | q^D)}.$$

We will from now on skip $f^I$ in the notation and write $p(f^0 | f(D, f^I)) = p(f^0 | D)$ and $p(f^0 | f^I) = p(f^0)$. We can write this in a form

$$p(f^0 | D) = \frac{\sum_{f'^0} T^D(f^0, f'^0) p(f'^0)}{\sum_{f''^0} \sum_{f'^0} T^D(f''^0, f'^0) p(f'^0)},$$

or shortly, in matrix notation

$$p^D = \frac{Tp}{\mathrm{Tr}\, TP},$$

with $P(f'^0, f''^0) = p(f'^0)$ and $\mathrm{Tr}$ denoting the trace. The so defined matrix

$$T^D(f'^0, f^0) = \delta_{f'^0, f^0} p(y^D | q^D, f^0),$$

is diagonal in $f^0$–representation. This shows that the pure states $f^0$ represent the possible fixed points of learning.[17]

Section 2.2 discusses learning in factorial states.

The local basic model including the local probabilities $p(f_x^0)$ represent the *local part of prior information*, the definition of generalized questions the *structural part of prior information*. Thus, a factorial state has only local and structural prior information and in any state of knowledge the nonlocal, nonstructural information should be information resulting from measurement or control added to a factorial starting state. After clarifying the importance of nonlocal information we will discuss their possible measurement. Table 1 summarizes some of the notations.

**Examples of basis questions**

The basis questions $p(y | x, f^0)$ define the answer probabilities for the available data $Q^D$ and relevant questions $Q^l$. Their definition is therefore not independent of the available measurement devices for $Q^D$ and $Q^l$, and they cannot be chosen arbitrarily. Choosing a set of $p(y | x, f^0)$ is part of the local prior knowledge. Mainly interested in generalization we do here not concentrate on the local part. We assume a formulation of the problem with a local prior $p(f_x^0)$ where the local part is learnable, i.e. where the local state of knowledge $p(f_x^0 | f_x(D^x))$ asymptotically converges to one $f_x^0$ under local questions $x$. according to some convergence criteria of our choice. This is without loss of generality, because we could always consider questions about the possible densities $p(y | x, f^0)$, i.e. study the density approximation problem. Practically, this does not change the situation, but splits formally the local part into several parts, so the

---

**Notations**

| | |
|---|---|
| $x \in X$ | basis questions |
| $q \in Q$ | generalized questions |
| $y \in Y$ | (probabilistic) answers |
| $\hat{y}, \hat{q}$ | (probabilistic) actions |
| $z$ | internal variables |
| $D = (D^q, D^y)$ | data |
| $f^0 \in F^0$ | pure states |
| $f \in F$ | states of knowledge |
| $\hat{f} \in \hat{F}$ | action states |
| $p(y | q, f^0)$ | probability (or density) for answer $y$ under question $q$ in state $f^0$ |
| $p(z | x, y, q)$ | probability (or density) for internal $z$ given data $(x, y)$ in question $q$ |
| $p(\hat{y} | q, \hat{f})$ | probability (or density) for action $\hat{y}$ under question $q$ in action state $\hat{f}$ |
| $p(\hat{q} | y, \hat{f})$ | same in inverse model for action $\hat{q}$ |
| $p(f^0 | f)$ | probability of $f^0$ under $f$ |
| $\bar{y}_x(f^0)$ | regression function at $x$ in state $f^0$ |
| $L = \ln p$ | log-probability |
| $l(q, y, z)$ | loss function |
| $\tilde{l}(q, y, \hat{f})$ | loss function integrated over probabilistic action or for deterministic action |
| $r$ | risk functional |

Table 1: Some notations frequently used in this paper

---

[17]Notice, that if $T^D$ has degenerate eigenvalues, also non-pure states may be unchanged. See Theorem in Section 2.2.

problem appears as a nonlocal one. [18] Also, in a possible loss function for density

Also, $x$ must not necessarily be a single minimal component, but we can combine many $x$ (with e.g. previously learned dependencies) to a larger $x$ vector. Technically, one $x$ just denotes one independently parameterized of subset of $F^0$, and the question of generalization is the question of generalization between such sets.

Basis questions can be Gaussian

$$p(y|x,f^0) \propto e^{-\frac{(y-\bar{y}_x(f^0))^2}{2\sigma_x^2}},$$

so that states $f^0$ are parameterized by their regression function $\bar{y}_x(f^0)$. In general, the parameterization of the $p(y|x,f^0)$ (and therefore of the states) can be arbitrary, e.g. also the variance $\sigma_x^2$ or higher order moments can be $f^0$–dependent.

Consider as a more complex example image $y$ producing states (generative models), e.g. with $x$ having the values face and non–face. States $f^0$ are defined by their generation probabilities for images of faces $p(y|\text{face}, f^0)$ and of non–faces $p(y|\text{non–face}, f^0)$. Generation of faces in a state $f^0$ could be defined

$$p(y|\text{face}, f^0) = \int dv\, p(v|\text{face}, f^0) p(y|v, \text{face}, f^0),$$

with $v$ being an index for the different variants of a face. Possible variants include for example interpersonal differences, varying view points or changing illumination conditions. Using some interpolation scheme, like optical flow and correspondence of some reference points, a continuous $v$ could be constructed out of a discrete set of examples. Also, human prior knowledge may be that faces have constituents $j$ like two eyes, mouth and nose appearing and being combined in different variants. In the easiest version with independent constituents one could choose

$$p(y|v, \text{face}, f^0) \propto e^{-\sum_{i,j} \frac{(d_i^{v,j})^2}{2\sigma^2(i,j,v)}},$$

using some distance $(d_i^{v,j}(y_i))^2 = ||y_i - y_{v,j}^i||^2$ with $i$ being the pixel index and $y_i^{v,j} = y_i^{v,j}(f^0)$ being a template for variant $v$ for constituent $j$. Then the $p(v|\text{face}, f^0)$ parameterize the face states.

## 2.2 Analysis of generalization

### 2.2.1 Minimal models and sufficient data

The ability to generalize is essential for any real learning. It is easy to see that generalization requires *nonlocal dependencies* contained in the $p(f^0|f)$. Let us denote by $D_{X\setminus x}$ data which do not depend on $x$ and combine all $f_{x'}^0$ for $x' \neq x$ into $f_{X\setminus x}^0$. Then using $\sum_{f_{X\setminus x}^0} p(f_{X\setminus x}^0|D_{X\setminus x}) = 1 = \sum_{f_{X\setminus x}^0} p(f_{X\setminus x}^0)$ we see that in a factorial state where $p(f_x^0|f_{X\setminus x}^0) = p(f_x^0)$

$$p(y|x, f'(D_{X\setminus x}, f)) = \sum_{f^0} p(y|x, f^0) p(f^0|D_{X\setminus x})$$

[18]In a usual density estimation problem there is only one $x$ (with the meaning 'get the next $y$'), and, besides a local positivity restriction, a natural (nonlocal) normalization condition.

$$= \sum_{f_x^0} p(y|x, f_x^0) \sum_{f_{X\setminus x}^0} p(f_{X\setminus x}^0|D_{X\setminus x}) p(f_x^0|f_{X\setminus x}^0) = p(y|x, f).$$

This means that data not depending on $x$ can never change the answer probabilities to $x$. We can say that a factorial state allows no inference and represents a 'tabula rasa' situation with respect to generalization. Thus, starting from a factorial state we necessarily need nonlocal information to enable any nonlocal learning. We can always relate $f = f(D, f_{fact})$ to a factorial state $f_{fact}$ by

$$p(f^0|f) = p(f^0|D) = p(D|f^0) p_{fact}(f^0)/p(D).$$

Now we formulate this observation in a bit more general way and show under which conditions the conclusion can be reversed. We begin with a simple example to outline the general idea. We consider an example, where $F^0$ does not only consist of extremal points of the convex set $F$. Let a space $F^0$ with three possible pure states, be defined by $p(y = 1|x, f_1^0) = 1$, $p(y = -1|x, f_2^0) = 1$, $p(y = \pm 1|x, f_3^0) = 0.5$. This may be the probability that a certain gender is marked on an application form for girl schools, boy schools, and coeducated schools. For example a state of knowledge $f$ with $p(y = \pm 1|x, f) = 0.5$ can be expressed as

$$p(y|x, f) = \frac{1-a}{2}(p(y|x, f_1^0) + p(y|x, f_2^0)) + a p(y|x, f_3^0)$$

$$= p(f_1^0) p(y|x, f_1^0) + p(f_2^0) p(y|x, f_2^0) + p(f_3^0) p(y|x, f_3^0)$$

for any $0 \leq a = p(f_3^0) = 1 - \sum_{i=1}^2 p(f_i^0) \leq 1$. Now assume, that some data not related to $x$ change the probability of $p(f_3^0) = a$. Obviously, this has no influence on $p(y|x, f)$. We will call such a space $F^0$ non–minimal with respect to $x$, and the set of data $(x, y)$ not sufficient with respect to $F^0$. Now consider a space $Q_x$ of local questions $q_x$ which includes also repeated measurements of $x$. The probability for a repeated measurement $y(x), y'(x)$ is the product $p(y, y'|x, f) = p(y|x, f) p(y'|x, (y, x))$ which changes with $p(f_3^0)$. The probability for a measurement $(y = 1, y' = -1)$ would be zero if $p(f_3^0) = -1$, but $1/4$ if $p(f_3^0) = 1$. Then, the coefficients $p(f_i^0)$ which define a state $f$ are unique, and if they change, they change the probability distribution of some question $q_x$. We will say the space $F^0$ is minimal with respect to the set of data $D_x = (q_x, y_{q_x})$, and data $D_x$ are sufficient for $F^0$.

We define equivalence classes $f_{D_Q}^0 = [f^0]_{D_Q}$ of $D_Q$–*equivalent* states for the set of data $D_Q$ of the form $(q, y)$ with $q \in Q$ by

$$f_{D_Q}^0 = f_{D_Q}^0{}' \Leftrightarrow \forall (q, y) \in D_Q : p(y|q, f^0) = p(y|q, f'^0),$$

forming the set $F_{D_Q}^0$. In the case in which the data contain all $y \in Y_q$ for every question $q \in Q$ we speak of $Q$–equivalent states and write $f_Q^0$. We defined $F^0 = F_X^0$. The same constructions can be applied to states of knowledge yielding $f_Q$ and $F_Q$.

We define data $D_Q$ to be *sufficient* with respect to $F^0$ or equivalently the set $F^0$ to be *minimal* with respect to $D_Q$ iff all states of knowledge $f \in F(F_{D_Q}^0)$ are uniquely

decomposable into the $f^0_{D_Q}$, that is iff there is no solution $p(f^0_{D_Q})$ of the following system of homogeneous linear equations

$$0 = \sum_{f^0_{D_Q}} p(y|q, f^0_{D_Q}) p(f^0_{D_Q}).$$

This means that no pure state can be expressed by others

$$p(y|q, f^0_{D_Q}) = \sum_{f'^0_{D_Q} \neq f^0_{D_Q}} p(y|q, f'^0_{D_Q}) p(f'^0_{D_Q}).$$

Then the corresponding system of inhomogeneous linear equations for $p(y|q, f^0_{D_Q})$ is overdetermined and therefore there exists at most one solution for the state $p(f'^0_{D_Q})$. (At least one solution exists by construction.) This means that sufficient data determine the state of knowledge $f$ uniquely. In other words, the convex hull $F_{D_Q}$ of $F^0_{D_Q}$ does not contain equivalent states, i.e. $[F(F^0_{D_Q})]_{D_Q} = F(F^0_{D_Q})$.

To shorten the notation we can write for

$$p(y|q, f) = \sum_{f^0_{D_Q}} p(y|q, f^0{}_{D_Q}) p(f^0_{D_Q}),$$

introducing indices $i = (q, y)$, $j = f^0_{D_Q}$

$$p_i = \sum_j A_{ij} p^f_j,$$

i.e.

$$p = A p^f.$$

The matrix (or integral operator with kernel $A_{ij}$) $A = A(D, F^0)$ describes the model $F^0$ with components $f^0$ on a data vector $D$ with components $(q, y)$ and $p^f$ with components $p(f^0|f)$ is the state of knowledge. Minimality requires the number of independent pure states $f^0_{D_Q}$ to be not larger than the number of data, i.e. question–answer pairs $(q, y)$, which means that the rank of the matrix $A_{((y,q), f^0_{D_Q})} = p(y|q, f^0_{D_Q})$ is equal to the number of $f^0_{D_Q}$. Summarizing, minimality of the model with respect to data, or sufficiency of the data for a model, is defined as a situation where the 'model-data matrix' $A$ has rank equal to the dimension of $p^f$, i.e the number of $f^0_{D_Q}$.

Minimality/sufficiency can be achieved by deleting some $f^0$ or by including more $y$ or new $q$ in the minimality condition. This can be done by including new independent questions, using a finer scale or add new dimensions to $y$. Examples include repeated measurements of the same $x$, $Q_x^{(n)}$ consisting of $n$-tuples $(x_1 = x, x_2 = x, \cdots, x_n = x) = \times_{i=1}^n x$, $Q_x = \bigcup_{n=1}^\infty Q_x^{(n)}$, multiple measurements within $X$, $Q_X^{(n)}$ consisting of $n$-tuples $\times_{i=1}^n x_i$, with $F_{Q_X^{(n)}} \subset \bigotimes_i \mathcal{L}(F_{x_i})$ or of varying length $Q_X = \bigcup_{n=1}^\infty Q_X^{(n)}$ with $F_{Q_X} \subset \bigoplus_{n=1}^\infty \bigotimes_i^n \mathcal{L}(F_{x_{i,n}})$[19].

[19]Which is similar to the construction of the Fock space for a many particle system, except that here the space is restricted to the convex hull (constant $L_1$ norm) and not to a region with constant $L_2$ norm.

For two–components data vectors from $Q_X^{(2)}$ the probability $p(y_1, y_2|x_1, x_2, f) = p_{ij}^{(2)}$ has to fulfill,

$$p_{ij}^{(2)} = \sum_k A_{ik} A_{jk} p_k^f = \sum_k A_{ik} A_{k'j}^T \delta_{k,k'} p_{k'}^f,$$

and therefore

$$p^{(2)} = (A \otimes A) P^{f,(2)},$$

with diagonal $P_{k,k'}^{f,(2)} = \delta_{k,k'} p_k^f$. For general data vectors $Q_X^{(n)}$

$$p^{(n)} = (\bigotimes_n A) P^{f,(n)},$$

with $P_{k,k_1,k_2,\cdots,k_{n-1}}^{f,(n)} = \prod_{i=1}^{n-1} \delta_{k,k_i} p_k^f$. A model minimal for single measurements is minimal for multiple measurements. For a minimal model the $f^0$ are linear independent, and in case the number of data (equations, conditions) is larger than the number of $f^0$ there must exist a reduced $n \times n$ system with nonzero determinant, so the solution is unique. Thus, we choose a decomposition $A = \binom{A'}{A''}$ with $A'$ a square $n \times n$ matrix, so its determinant is defined, and for a minimal model there exists an $A'$ with $\det A' = 0$. The relation $\det(A' \otimes A') = (\det A')^{2n}$ for $n \times n$ matrices shows that the determinant for multiple measurements for the reduced system $A'$ is nonzero if it is nonzero for single measurements. If a given solution is already unique for a reduced system, it is also unique in an extended system, where just more conditions are added, consistent with the solution. Thus, models minimal for single measurements $Q_X^{(1)}$ are also minimal for multiple measurements, i.e. vectors $Q_X^{(n)}$.

Non–minimal local spaces $F_x^0$ are not commonly used. Minimal for example is a local space consisting of Gaussians at different locations. Then a convex linear combination is not Gaussian but a Gaussian mixture state.

### 2.2.2 Factorial priors

Let us now consider two sets of data, $Q^l$ (e.g. relevant questions) and $Q^D$ (e.g. training questions) with corresponding (test) data $D^l$ and (possible sets of training) data $D = \{D_i, i \in I\}$. The $D_i = (y_i, q_i)$ are allowed to be data vectors and may represent one possible collection of data which can be obtained during training.

In a minimal model the following theorem states that the prior probabilities reflect already the possibility or impossibility of generalization:

**Theorem**: For a set of (test) data $D^l \subseteq D_{Q^l}$ sufficient for $F_l^0 = F_{Q^l}^0$ (or equivalently for $F_l^0$ minimal with respect to $Q^l$) and another set of (potential training) data $D \subseteq D_{Q^D}$ sufficient for $F_D^0 = F_{Q^D}^0$ (or equivalently $F_D^0$ minimal with respect to $D$) the following proposition holds

$$\forall (q, y) \in D^l, \forall D_i \in D : p(y|q, f) = p(y|q, f'(D_i, f))$$
$$\Leftrightarrow \forall f^0_{l,D} : p(f^0_{l,D}) = p(f_l^0, f_D^0) = p(f_l^0) p(f_D^0),$$

where $f^0_{l,D} = f^0_{l \cup D} = [f^0]_{D_{Q^l} \cup D}$. The backward direction does not require the two sufficiency conditions.

The theorem gives conditions under which conditional independence of all relevant data of all training data, is equivalent to independence between all $f_l^0$ and $f_D^0$. The stated conditional independence also means that the actual state of knowledge is an eigenstate for all matrices $T^{D_i}(f'^0, f^0) = \delta_{f'^0, f^0}\, p(y^{D_i}|q^{D_i}, f^0)$, or that all *mutual informations*

$$\ln \frac{p(y|q^l, f, D_i)}{p(y|q^l, f)} = \ln \frac{p(y, q^l, D_i|f)}{p(y, q^l|f)p(D_i|f)},$$

are zero, and therefore also all averages of them.

*Proof:* We show that factorial priors do not allow generalization, and that sufficiency of $D^l$ and $D$ excludes the stated no–generalization property for other priors.

For $q \in Q^l$, abbreviating $p(y|q, f'(D_i, f))$ by $p(y|q, D_i)$, we write for the no–generalization condition

$$p(y|q, f) = p(y|q, D_i) = \sum_{f^0} p(y|q, f^0)p(f^0|D_i)$$

$$= \sum_{f_l^0} \sum_{f_D^0} p(y|q, f_l^0)p(f_l^0|D_i)p(f_D^0|f_l^0, D_i)$$

$$= \sum_{f_l^0} p(y|q, f_l^0)p(f_l^0|D_i),$$

because for $q \in Q^l$ the probability $p(y|q, f^0)$ only depends on $f_l^0$ and $\sum_{f_D^0} p(f_D^0|f_l^0, D_i) = 1$. Another summation over finer classes up to $D_X$ is not necessary because $\sum_{f_{l \cup D \cup D_X}^0} p(f_{l \cup D \cup D_X}^0|f_{l \cup D}^0) = 1$. Setting $p(f_l^0|D_i) = p(f_l^0)$ yields $p(y|q, f) = p(y|q, D_i)$ giving one solution of the nonlearnability condition. For a $D^l$–sufficient model the state $p(f_l^0|D_i)$ is uniquely determined by the probabilities for the relevant data $p(y|q, f^0)$, $q \in Q^l$ and $y \in y^q$ so $p(f_l^0|D_i) = p(f_l^0)$ is also the only solution. Thus, sufficiency of $D^l$ excludes the possibility that the influence of data only consists in switching between $D^l$–equivalent states.

Now we show that independence of $f_l^0$ of the data $D_i$ together with the minimality of $F_D^0$ with respect to $D$ only allows a factorized prior probability. We insert a summation over $f_D^0$ into $p(f_l^0|D_i)$ and write for the condition $p(f_l^0|D_i) = p(f_l^0)$

$$p(f_l^0) = p(f_l^0|D_i) = \sum_{f_D^0} p(f_D^0|D_i)p(f_l^0|f_D^0).$$

One sees that the condition $p(f_l^0|D_i) = p(f_l^0)$ is fulfilled for $p(f_l^0|f_D^0) = p(f_l^0)$. This already solves the backward direction without the need of any minimality or sufficiency condition. For a model not minimal on $F_D^0$ there still might be different states on $D$ leading to the same posterior $p(f_l^0|D_i)$ on $F_l^0$. We now use minimality of the data $D$, to exclude the possibility, that there are dependencies which cannot be explored by $D$. The probability $p(f_D^0|D_i)$ is related to the definition of the training questions, i.e. to the $p(y_i^D|q_i^D, f_D^0)$, by

$$p(f_D^0|D_i) = \frac{p(y_i^D|q_i^D, f_D^0)p(f_D^0)}{\sum_{f_D^0} p(y_i^D|q_i^D, f_D^0)p(f_D^0)}.$$

Inserting this equation into the above equation for $p(f_l^0|D_i)$ shows that according to the assumption of sufficient data $D$ the coefficient vector $a(f_D^0)$ multiplying the matrix $A_{i,f_D^0} = p(y_i^D|q_i^D, f_D^0)$ must be unique:

$$\frac{p(f_l^0)p(f_D^0)}{\sum_{f_D^0} p(y_i^D|q_i^D, f_D^0)p(f_D^0)} = \frac{p'(f_l^0|f_D^0)p'(f_D^0)}{\sum_{f_D^0} p(y_i^D|q_i^D, f_D^0)p'(f_D^0)},$$

where $p'(f_l^0|f_D^0)p'(f_D^0) = p'(f_l^0, f_D^0)$ denotes another solution of the joint probability. Summation over $f_l^0$ and $f_D^0$ gives equality for the $f^0$–independent denominators on both sides so that $p(f_l^0, f_D^0) = p(f_l^0)p(f_D^0)$ is the only solution. q.e.d.

Without restriction to sufficient data there might exist spurious dependencies between equivalent states of knowledge, which are not observable within the given set of relevant questions.

The formal structure of the Theorem and of its proof becomes clearer if written in a more abstract matrix formulation. With $i = (q, y)$, $j = (q^D, y^D)$, $k = f_l^0$, $l = f_D^0$, $A_{ik} = p(y|q, f_l^0)$, $B_{jl} = p(y^D|q^D, f_D^0)/p(y^D|q^D)$ (where for observed data the denominator is unequal to zero), we can write for the posterior $p(y|q, D)$, in components

$$p_{ij} = \sum_{k,l} M_{ij,kl} p_{kl}^f$$

$$= \sum_{k,l} A_{ik} B_{jl} p_{kl}^f,$$

which reads for matrices

$$p = (A \otimes B)p^f.$$

**Theorem (matrix formulation):**
For a (sub)system of equations with invertible $n \times n$ matrix[20] $A$ and $m \times m$ matrix $B$, i.e. with $\det A \neq 0 \neq \det B$ (Minimality/Sufficiency), the following holds

$$p_l \otimes 1 = p_l \otimes p_D = (A \otimes B)p^f \Leftrightarrow p^f = p_l^f \otimes p_D^f.$$

Indeed, according to

$$\det(A \otimes B) = (\det A)^n (\det B)^m$$

with $A$ and $B$ also $A \otimes B$ is invertible and using

$$(A \otimes B)^{-1} = (A^{-1} \otimes B^{-1})$$

and

$$(A \otimes B)(p_l \otimes p_D) = (Ap_l \otimes Bp_D)$$

we find

$$p^f = (A^{-1}p_l \otimes B^{-1}p_D) = p_l^f \otimes p_D^f,$$

thus $p^f$ factorizes. Formulated for probabilities conditioned on $D$, we used $p_D = 1$ for the theorem. Formulating the theorem for joint probabilities $p(y, q, D|f)$, or similarly $p(y, y^D|q, q^D, f)$, symmetrizes the formulation and gives $A_{ik} = p(y, q|f_l^0)$, $B_{jl} = p(y^D, q^D|f_D^0)$,

---

[20]For simplicity we skip the prime for $A$ and $B$ which we used previously for reduced matrices.

$(p_l \otimes p_D)_{ij} = p(y, q|f_l^0)p(y^D, q^D|f_D^0)$, with still $M = A \otimes B$ according to the definition of the model.

Formulated in a basis $X$ the theorem gives the

**Lemma** : For $(q_x, y) \in D^x \subseteq D_x$, $D^{X\backslash x} \subseteq D_{X\backslash x}$ and requiring for the forward direction sufficiency of the set $D^x$ for $F_x^0$ for all $x$ and sufficiency of the set $D^{X\backslash x}$ for $F_{X\backslash x}^0$ for all $x$ the following holds

$$\forall x \in X, \forall (q_x, y) \in D^x, \forall D_i^{X\backslash\{x\}} \in D^{X\backslash\{x\}}:$$

$$p(y|q_x) = p(y|q_x, D_i^{X\backslash\{x\}})$$

$$\Leftrightarrow \forall f^0 \in F^0 : p(f^0) = \prod_{x \in X} p(f_x^0).$$

Remarks:

1. The theorem states that a factorial state remains factorial after local learning and that there is no generalization possible across a 'factorization border' for *any* learning algorithm. In factorial states any information concerning answers to question $x$ can only come from questions depending explicitly on the point $x$ of interest, but not necessarily depending only on $x$. Answers to questions depending only on $X \backslash X^D$ cannot be learnt under a factorial prior, and training data not depending on $X^l$ are uninformative with respect to relevant questions if not combined with other information. They can thus only have indirect influence. For analysis of prior information in terms of relevant questions we may choose $X = Q^l$. Then in order to enable learning for every $q^l = x$ for every $x$ we must have data $D$ depending (maybe not only) on that $x$. For example, a smoothness constraint in statistics depends on all $x$ and can allow generalization. Especially easy to analyze is the situation when $Q^D \subseteq X$. Then all relevant questions are defined with respect to the data. In this case a factorial prior on $X$ refers to the training questions itself and the relevant questions are directly defined in their dependency on the training questions. We discuss the relations in detail in the next subsection.

2. The forward direction, factorial state $\Rightarrow$ no generalization, is related to so called 'No Free Lunch'–theorems (Wolpert, 1996a, 1996b), generalized from uniform priors (or uniform meta priors) to factorial states, without explicitly referring to a specific form of loss function or algorithm. Indeed, *uniform priors* of the form $p(f_x^0) = p(f_{x'}^0), \forall x'$ — possibly also with uniform probabilities $p(f_x^0) = p(f_x'^{,0}), \forall f_x'^0, \forall x$ within single $x$–components so that $p(f^0) = p(f'^{,0}), \forall f'^0$ — are special factorial priors. With respect to nonlearnability of specific sets of questions the theorem is sharper than results from the theory of uniform convergence (Vapnik, 1982) or regularization theory (Tikhonov, 1963), and independent of a specific loss function. However, it does not give quantitative results. Also, a change of $p(y|q, f)$ does not necessarily need to change the

decision $\hat{f}$ (see 4.). A quantitative measure of dependency can be obtained by averaging the mutual information over a distribution for data $D_i$ and $D^l$, or by integrating out $D^l$ and calculate an average change of the risk under a data distribution $p(D_i)$.

3. The backward direction gives sufficient conditions for a no–generalization state to be factorial. The main point is that if the model has only one way to express the same state then it cannot switch between equivalent descriptions and there is no possibility to create formal dependencies without visible effects.

4. A change of $p(y|q, f)$ does not have to a change the final decision $\hat{f}$. It may be that

   a. the loss function is not sensitive to change in $p(y|x, f)$,

   b. the risk, i.e a certain property of $p(l|f, \hat{f})$ like its expectation, is not sensitive to a change in the loss,

   c. the decision is not sensitive to a change in the risk.

   Of principal interest are the optimality equivalence classes

   $$[f^0]_{r^*} = [f'^0]_{r^*} \Leftrightarrow$$
   $$\hat{f}^* = \mathrm{argmin}_{\hat{f} \in \hat{F}} r(f^0, \hat{f}) = \mathrm{argmin}_{\hat{f} \in \hat{F}} r(f'^0, \hat{f}),$$

   identify $f^0$ leading to the same decision $\hat{f}^*$. While this requires already calculation of the optimal decision it may be easier to calculate

   $$[f^0]_r = [f'^0]_r \Leftrightarrow \forall \hat{f} : r(f^0, \hat{f}) = r(f'^0, \hat{f}).$$

   However for both variants, available data may have differing probabilities for different $f^0$ within such classes, so that $p(f^0|D)$ can vary within a class. Thus, analog to $f_{l,D}$ in the theorem, one has to consider finer equivalence classes

   $$[f^0]_{r^*, D} \quad \mathrm{or} \quad [f^0]_{r, D},$$

   where only data $D^f$ which factorize, like $p([f^0]_{r,D}) = p([f^0]_{r,D^{nf}})p([f^0]_{D^f})$, can be skipped. All those equivalence classes can be finer than $[f]_{r^*}$, so this does not guarantee that data changing $p(f^0)$, i.e. the state of knowledge $f$, change the decision $\hat{f}^*$.

5. Pure states are special factorial states. In this sense factorial states are the possible starting and ending points of a learning process. For $x$–variables which are visible and known, $f^0$ can be seen as collection of the remaining (stationary distributed) random variables. Making part of these random variables $f^0$ visible, that is moving them into $x$, breaks the old pure states in new pure components. The learning process would ideally go from the old pure state, now assumed to be incompletely specified and therefore only a factorial state of knowledge, to one of the possible new pure states being factorial and maximally specified. We do not restrict pure states to be deterministic. Also probabilistic states can be not further decomposable in the situation under investigation.

### 2.2.3 Factor dimension

We now look for conditions under which learning can occur, i.e. when $p(y|q) \neq p(y|q, D_i)$ so not all mutual informations between available and relevant data are zero. We are mainly interested in the effects of generalization where $q$ itself is not part of $D$, Indeed, if for continuous $q^l$ the $q^l$ which are directly in the training set have usually measure zero, so learning has to go over generalization.

Data not changing $p(f^0)$ have equal probabilities under all pure states $f^0$ or, equivalently, equal diagonal elements (and therefore eigenvalues) of the matrix $T(f'^0, f^0)$, projected to the subspace $F_x^0$, for all $f_x^0$. A unique maximal $f^{0,*}$, i.e. one with maximum probability for $D_i$ and not excluded from $f$, would be sufficient to ensure learning for all $f_x$ which are not yet in the state $f_x^{0,*}$. A change in $p(f^0)$ may however correspond to *non-relevant learning*, changing not the relevant distributions $p(y|q, f)$ but only higher order interactions like $p(y, y'|q, q', f)$.

If the potential data $D$ are sufficient and the model is minimal for the relevant questions, then according to the negated version of the theorem for non–factorial priors there are training data $D_i$ so that relevant data $(q, y) \in D^l$ change, Consider a locally minimal model with $x = q^l \in X$ local data $D^x = (x^D, y^D)$, $x^D \in X$ and a given prior $p(f^0|f)$. Local data $q^D$ change the probability of $y$ under $x$ if $p(f^0) \neq p(f_x^0|D^x)$ and therefore $p(f^0) \neq p(f_x^0|f_D^0)$. Then within $p(f^0)$ states restricted to relevant and available data cannot factorize. Thus, the state of knowledge must fulfill the generalization condition for $x$ under $x^D$

$$p(f^0|f) \neq$$
$$p(f_x^0|f_{X\setminus(x,x^D)}^0, f)p(f_D^0|f_{X\setminus(x,x^D)}^0, f)p(f_{X\setminus(x,x^D)}^0, f).$$

We can characterize a given state of knowledge by its *factor dimension* $\dim_F(f)$ with respect to a set $X$, defined as the maximal number $n$ of $x_i \in \{x_1, \cdots, x_n\} = X^n \subseteq X$ so that still one $f_x^0$, with $x \in X^n$ the same for all $f^0$, can be factorized, conditioned on $f_{X\setminus X^n}^0$. Then learning has to occur under $f$ if the number of local data $n$ is at least the factor dimension of $f$, i.e. $n \geq \dim_F(f)$, as then no additional $x$ can factorize.

We give three examples:

Consider for a space $F^0$ parameterized by local means $\bar{y}_x$ for all $x$, defining a regression function $\bar{y}$ according to $\bar{y}(x) = \bar{y}_x$:

1. A nonlocal prior[21] $p(f^0 = \bar{y}) \propto e^{-c\sum_x (\bar{y}_x - c_x)^2}$. Even though this state explicitly depends on every single $x$ it is factorial in $X$.

2. A symmetry (e.g. smoothness) prior $p(f^0 = \bar{y}) \propto e^{-c\sum_x (\bar{y}_x - \bar{y}_{sx})^2}$, $s$ indicating some bijective symmetry transformation (e.g. translation) $x' = sx$. It has a degenerated maximum $\bar{y}_x = \bar{y}_{sx}$ which can be made unique by adding local data for every orbit of $s$. The orbit of $x$ under $s$ consists of all elements $s^i x$ which can be generated out of $x$ by applying $s$ any number of times, $0 \leq i \leq \infty$. The factor

dimension for each orbit is equal to one. In case of several orbits the total factor dimension is only $n - m + 1$, with $n$ the total number $x$, and $m$ the size of the smallest orbit.

3. Data with $p(f^0 = \bar{y}) \propto e^{c(\sum_i^n \bar{y}_{x_i})^2}$ (parity for zero/one variables), which require additional data about $n - 1$ different $x_i$ to give a unique maximum. The sum can be seen as a local question in a space $F^0$ with a linearly transformed basis of $f^0$. The factor dimension is $n - 1$.

We make the side remark that the factor dimension has similarities to the concept of VC dimension. In both cases generalization is impossible if the number of data is smaller than the corresponding dimension. The latter is usually applied to a family of loss functions indexed by $\hat{f} \in \hat{F}$ in the context of empirical risk minimization. A small $\dim_{VC}$ indicates that the loss function cannot vary too much, so the difference between minimal empirical and minimal expected risk can be bounded. These bounds are independent of the family $F^0$, as long as certain minimal conditions, e.g. bounded $p(y|q, f^0)$, are fulfilled. In contrast, the factor dimension, as used here, is independent of the loss function and therefore of $\hat{F}$. The factor dimension has a similar interpretation for $F^0$ instead for $\hat{F}$: Nonconstant, independent data require a family $F^0$ at least 'large enough' to contain all the different combinations, while a smaller $F^0$ which excludes certain combinations implies the possibility of generalization. We already mentioned that not only natural choices of initial states but also the asymptotical final, i.e. pure, states of a learning process are factorial states, having maximal factor dimension, Hence, learning can as well decrease as also increase the factor dimension. Usually one would expect a U-shape like dependence of the factor dimension on the learning process, analogous to the mutual information, which provides a quantitative measure of dependency between relevant and available questions. Commonly prior information reduces the factor dimension and, at least deterministic, local data increase it.

Let $\dim_{VC}^{\alpha_{x,y}}(F^0)$ denote the ($\alpha_{x,y}$–dependent) VC dimension of a family of functions $p(y|x, f^0), f^0 \in F^0$, which is the maximal number of different pairs $(x, y)$ so that for all of them exist $f^0, f'^0 \in F^0$ with $p(y|x, f^0) > \alpha_{x,y}$ and $p(y|x, f'^0) \leq \alpha_{x,y}$, $0 \leq \alpha_{x,y} < 1$, i.e. the maximum number of points $(x, y)$ which can be shattered by $F^0$. Assume now that there exists an $\alpha_{x,y}$ with $p(y|x, D_>^x) > \alpha_{x,y} \geq p(y|x, D_<^x)$ for at least one pair of local data $D_>^x = (x, y_>)$ and $D_<^x = (x, y_<)$, Thus, $\alpha_{x,y}$ separates the posterior probabilities $p(y|x, D_i)$ of at least two possible data $D_i$, and the actual state cannot be a pure state. Then with $f^0 \in F^0$ the components of $f$ with $p(f^0|f) \neq 0$

$$\dim_F(f) \leq \dim_{VC}^{\alpha_{x,y}}(F^0).$$

Indeed, if we assume $\dim_{VC}^{\alpha_{x,y}}(F^0) = n$ then according to the definition of the VC dimension within a set of $n + 1$ questions $x$ there is at least one for which either $p(y|x, f^0) > \alpha_{x,y}$ or $p(y|x, f^0) \leq \alpha_{x,y}$ is impossible

---

[21] We will discuss the measurement, i.e. the corresponding nonlocal questions in Section 3.2.

for all $f^0 \in F^0$. As convex combination the probability $p(y|x, f)$ for state $f$ can only lie between extremal points, i.e. pure states $f^0$. Therefore, the assumption of existence of both $D^x_<$ and $D^x_>$, used to construct $\alpha_{x,y}$, implies that there exist $f^0$ with $p(y|x, f^0) > \alpha_{x,y}$ as well as $f'^0$ with $p(y|x, f'^0) \leq \alpha_{x,y}$ with nonzero probability $p(f^0|f) \neq 0 \neq p(f'^0|f)$. Thus, the probability $p(f^0_x, f^0_{x_1}, \cdots, f^0_{x_n}|f)$ cannot factorize, which gives $\dim_F(f) \leq \dim_{VC}^{\alpha_{x,y}}(F^0)$.

If the number of data is smaller than the factor dimension of $f$ learning can still be relevant if $Q^l$ contains multiple measurements (with respect to what has been considered an element in calculating the factor dimension). Multiple measurements of single components $q^l$ means considering their interactions also to be relevant. Using a measurement vector with one answer for every $x \in X$, i.e. $q = X$, makes for locally minimal models every learning for local data relevant. For example, knowing $y_1 - y_2 = 0$ does not change $p(y_1)$ or $p(y_2)$ if they have equal prior. However, it drastically changes the probability for a joint measurement of $y_1$ and $y_2$. This is an example where learning only occurs for non–relevant higher order dependencies. If higher order dependencies are missing, so that for a state with factor dimension $n$ already $m$ relevant questions factorize, then only $n - m$ data are necessary to allow generalization to a single component $q$.

An example of taking into account dependencies between components $q^l$ by multiple measurements, is the special case of the risk of an algorithm which uses past test data $(y_i, q^l_i)$, $i < n$ to improve the selection of the next action for $q^l_n$. Here the next action and therefore also the loss function does depend on all previous test data $(y_i, q^l_i)$. Then the risk does not consist of a sum of terms depending on disjunct sets of $q^l_i$, and correlations between components $q^l_i$ up to order $n \geq i$ can be important. Thus, the effective $q^l$ is the vector of components $q^l_i$, and the expected risk (on–line risk of an algorithm) is an average over this vector $q^l$. Correlations between answers to different $q^l$ (vectors, not components $q^l_i$) are not measured by this risk. However, the dependencies between vectors $q^l$ should be smaller than between components $q^l_i$ if $p(f^0|(y, q^l))$ is nearer to a pure state then $p(f^0|y_i, q^l_i)$.

Loosely speaking, we only need to know what we can see, we only can know what we have seen, and it is structural information which allows us to infer indirectly.

### 2.2.4 Generalization–related sets of questions

We now give the definition of some sets of questions, related to the previous analysis of generalization. We choose $Q = Q^l \cup Q^D$, $X = X^l \cup X^D$ and define the following sets of questions: Training questions $q$ in $Q^D$ but not in $Q^l$ are called non–test questions $q^{\neg l} \in Q^D \setminus Q^l = Q^{\neg l}$, test questions corresponding to questions $q$ in $Q^l$ but not in $Q^D$ are non–training questions $q^{\neg D} \in Q^l \setminus Q^D = Q^{\neg D}$. Analogously, we write $x^{\neg D} \in X^l \setminus X^D = X^{\neg D}$ for the non–training basis with $x \in X^l$ but $x \notin X^D$, as well as $x^{\neg l} \in X^D \setminus X^l = X^{\neg l}$ for the non test basis with $x \in X^D$ but $x \notin X^l$. $Q^{l,D} = Q^l \cap$

$Q^D$ denotes the set of common questions and $X^{l,D} = X^l \cap X^D$ the common basis. Clearly, we have $Q^{l,D} \cap Q^{\neg D} = \emptyset$, $Q^{l,D} \cap Q^{\neg l} = \emptyset$, $Q^l = Q^{l,D} \cup Q^{\neg D}$, $Q^D = Q^{l,D} \cup Q^{\neg l}$ and the corresponding relations for $X$. Notice that the common (or non test, non–training) basis is not necessarily the basis of the common (or non–test, non–training) questions.

Let us further introduce for sets $Q' \subseteq Q$, $X' \subseteq X$ the notation $Q'_{X'} \subseteq Q' \subseteq Q$ for the set of all questions within $Q'$ depending only on $X'$, i.e. with a basis completely within $X' \subseteq X$ so that $X^{Q'_{X'}} \subseteq X' \subseteq X$. Similarly, $Q'_{(X')} \subseteq Q$ denotes the set of all questions depending (not necessarily only) on $X'$, i.e. with a basis having nonzero intersection with $X' \subseteq X$ so that $X^{Q_{(X')}} \cap X' \neq \emptyset$. Obviously, $Q'_{X'} \subseteq Q'_{(X')} \subseteq Q'$, especially $Q = Q_{X^Q} = Q_{(X^Q)}$, and $Q'_{X'} \subseteq Q'_X$, $Q'_{(X')} \subseteq Q'_{(X)}$ for $X' \subseteq X$. We can partition a set $Q' \subseteq Q$ into two disjunct subsets with respect to $X' \subseteq X$ according to $Q' = Q'_{X'} \cup Q'_{(X \setminus X')}$. In particular, $Q^l = Q_{X^{\neg D}} \cup Q^l_{(X^D)}$ and $Q^D = Q_{X^{\neg l}} \cup Q^D_{(X^l)}$ with

$$Q^l \supseteq Q^l_{(X^D)} \supseteq Q^l_{X^D} \supseteq Q^{l,D} \subseteq Q^D_{X^l} \subseteq Q^D_{(X^l)} \subseteq Q^D,$$

and accordingly

$$X^l = X^{Q^l} \supseteq X^{Q^l_{(X^D)}} \supseteq X^{l,D} \supseteq X^{Q^l_{X^D}} \supseteq X^{Q^{l,D}}$$

$$X^{Q^{l,D}} \subseteq X^{Q^D_{X^l}} \subseteq X^{l,D} \subseteq X^{Q^D_{(X^l)}} \subseteq X^{Q^D} = X^D.$$

(See Fig.1 and for more details Fig.2.)

We have shown in the previous Subsection that in a factorial state the non–training basis $X^{\neg D}$ is not learnable and the non–test basis $X^{\neg l}$ does not influence relevant $q^l$ directly. However, $X^{\neg l}$ might well have indirect influence and act like noise sources within other questions $Q^D$ depending on $X^{l,D}$. Thus, information about states corresponding to $X^{\neg l}$ enables learning about the noise structure within $Q^D$. Thus, with reference to factorial states, we will call the set of questions $Q^l_{X^{\neg D}} = Q_{X^{\neg D}} \subseteq Q^l$ depending only on $X^{\neg D}$ unlearnable questions, $Q^l_{(X^D)} = Q^l \setminus Q^l_{X^{\neg D}}$ (potentially) learnable questions, the set $Q^D_{X^{\neg l}} = Q_{X^{\neg l}} \subseteq Q^D$ indirect questions depending not on $X^l$, and $Q^D_{(X^l)} = Q^D \setminus Q^D_{X^{\neg l}}$ direct questions. Indirect questions alone cannot contribute to knowledge about $Q^l$, their influence is indirect by contributing information about the unknown noise sources for questions in $Q^D_{(X^l)}$. Note that $Q^{\neg D} \supseteq Q^l_{X^{\neg D}}$ and $Q^{\neg l} \supseteq Q^D_{X^{\neg l}}$ which means that non–training questions can (and hopefully do) depend also on $X^D$ (so indirect information about them may be available) and non–test questions also on $X^l$ (so they can contribute directly). The non–training basis $X^{\neg D}$ could be eliminated by integrating over a corresponding (factorial) prior. Then all relevant questions depend only on $X^{l,D}$, i.e. we have $Q^l_{(X^{l,D})} = Q^l_{X^{l,D}}$. Therefore, we will call $Q^l_{X^{l,D}} = Q^l_{X^D}$ the set of effective questions.

In the following we give some simple examples (compare Section 2.2) that for $Q^l_{(X^D)}$ learning is possible, but not necessarily for individual $q^l \in Q^l_{(X^D)}$ if the dependencies between different $q^l$ are not considered rele-
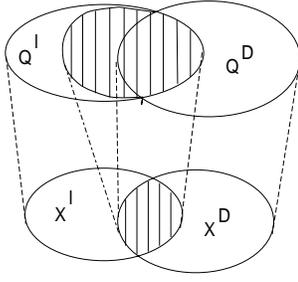
Figure 1: The figure shows the set of questions $Q$ and their corresponding basis $X$. The shaded area within $X$ represents the common basis $X^{l,D}$, the shaded area within $Q$ the learnable questions $Q^l_{(X^D)}$. Learning related to questions within the shaded area but not within the data $Q^D$ is called generalization. See Fig.2 for more details and notation.
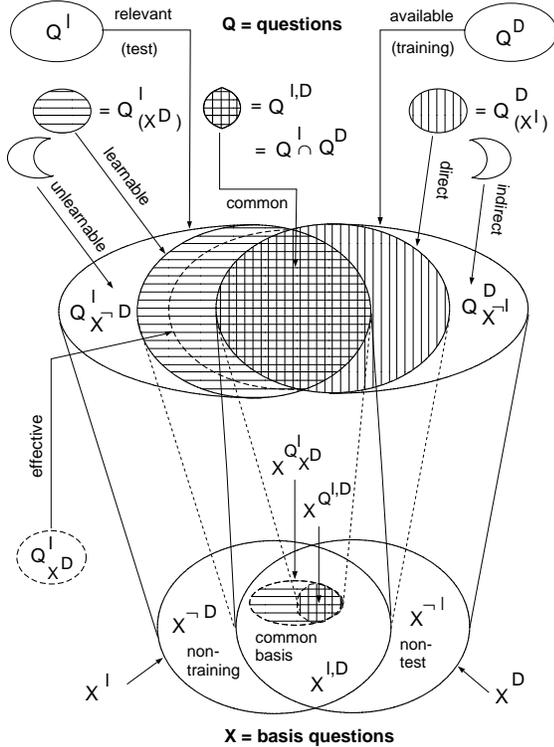


Figure 2: This figure shows the relations in more detail than Fig.1. The basis of relevant questions $Q^l$ is $X^l$, and for questions $Q^D$ corresponding to available data the basis is $X^D$. While $Q^D$ is assumed to be always finite $Q^l$, $X^l$ and $X^D$ can be infinite. A model describing learning can be restricted to the common basis $X^{l,D} = X^l \cap X^D$ which is the double shaded area below. Only for questions depending on that intersection learning can occur. This set $Q^l_{X^D}$ is symbolized by the upper horizontally shaded area. Learning the set $Q^l_{X^D} \setminus Q^{l,D}$ is called generalization. Questions in $Q^l_{X^{\neg D}}$, being in $Q^l$ but not in its shaded area, are unlearnable. Questions $Q^D_{(X^l)}$ can contribute information about $Q^l$ directly. Questions in $Q^D_{X^{\neg l}}$, being in $Q^D$ but not in its vertically shaded area $Q^D_{(X^l)}$, are not directly related to relevant questions and can only contribute indirectly, i.e. in combination with $Q^D_{(X^l)}$.

vant. Consider deterministic $x_1$ and $x_2$ with possible values $y_i = \pm 1$ with independent prior probabilities equal $1/2$. Then knowing only $y_1 y_2 = \pm 1$ without knowing something about $y_2$ says nothing about $y_1$. Thus, the prior factorizes with respect to the equivalence classes $[f^0]_{x_1 \cup (x_1 x_2)}$ or $[f^0]_{x_2 \cup (x_1 x_2)}$, if we choose $x_1 x_2$ and $x_1$ (resp. $x_2$) as basis questions. If we choose all three as basis then $[f^0]_{x_1 \cup x_2 \cup (x_1 x_2)}$ cannot be written in a factorial form. Compare with the question $y_1 + y_2$ for which also no information about $y_1$ results from $y_1 + y_2 = 0$, however $y_1 + y_2 = 2$ determines $y_1$ and $y_2$. As soon as either $y_1$ or $y_2$ is known both nonlocal questions give the value of the missing $y_i$. If we do not choose $Q^l = \{x_1, x_2\}$ but instead e.g. $Q^l = \{|y_1 - y_2|\}$ the prior can be written in a factorial form with respect to $y(x) = |y_1 - y_2|$ and $y'(x') = y_1 + y_2$: $p(f^0_x)p(f^0_{x'})$ as those are in the given situation independent events.

## 2.3 Why structural information?

Structural information is our knowledge about the definition of questions corresponding to predetermined dependencies between answers, or in other words, our knowledge of what we are measuring. One might wonder whether such predetermined dependencies between answers are necessary and corresponding data can be available and useful.

First, we remark, as discussed in the previous section, that without nonlocal information generalization is not possible. Any nonlocal information has two parts, a structural one represented by the definition and the answer to be determined by empirical measurement or control. Thus, in other words, in any case there must be predefined dependencies between answers to some questions in order to be able to infer to unseen questions. So, as we will discuss in more detail below, requiring a bound on a certain smoothness or symmetry property one must for example know: 1) the definition of smoothness or symmetry in mind corresponding to the structural part of information or the predefined dependency of the smoothness question from the other questions, 2) a bound on the allowed function values obtained by empirical measurement, active enforcing or pure assumption.

Secondly, we refer to the observation that it is very common to assume dependencies between answers. In logic these dependencies are called rules. With zero–one function values the logical operations can be written using multiplication, addition and the step function. Fuzzy set theory deals with unsharp rules and Bayesian belief networks are the general probabilistic formulation. Compared with logic based artificial intelligence expert systems which construct many rules with often low nonlocality (humans can better deal with lower order correlations) usual statistical models use only one, but highly nonlocal rule, e.g. some bound on a smoothness functional. However this difference is more of practical than of principal theoretical nature in as far as any number of rules, with maybe individually low nonlocality, can be combined into one ('vector') rule which can have high nonlocality. The definition a macroscopic observable in its dependence on microscopic variables is an example for structural information. One may say the usual macro-

12

scopic variables in form of integrals over all $x$ have maximal nonlocality, depending on all basis questions. The definition of an 'energy question' could have the following form in a real, scalar, Euclidean field theory

$$p(y|q_E, f^0) = \delta \left( \int dx \left( (\frac{d}{dx} f(x))^2 + m(x) f^2(x) \right) \right.$$

$$\left. + \int dx \left( h(x) f(x) + E_0(x) \right) - y \right),$$

where the first term has the form of a 'smoothness prior'. In terms of this analogy to physics, a statistical approximation problem with mean square error $\sum_i a_{x_i} (f(x_i) - y_i^D)^2$, data $y_i^D$ and a smoothness prior might be seen as minimizing a 'free' massless kinetic term, a $x_i$–dependent 'mass' $m(x) = a_x \sum_i \delta(x - x_i)$ and 'external field' $h(x) = -2a_x \sum_i \delta(x - x_i) y_i^D$ coupling to $f(x)$ at the data points $x_i$ and a constant $E_0 = a_x \sum_i \delta(x - x_i)(y_i^D)^2$. Higher order (local or nonlocal) interactions correspond to other forms of nonlocal questions. We will discuss in more detail below how answers to questions with infinite nonlocality can be enforced by measuring devices or by controlling the situations under investigation.

Finally, we point out that one always has to know what one measures, at least in a probabilistic sense, and therefore structural information is necessary in any case. When repeating measurement of basis questions $x$ in a fixed state $f^0$ we assume that we use the same device with a stationary answer distribution. This means, for example, that if we measure in the same state hundred times $y = 5$ followed by a hundred $y = 42$ we might assume that our question $x$ is incompletely specified and another (hidden) variable should have been included which probably has changed its value after hundred measurements. Such not directly observable or hidden variables define the state $f^0$ and we can say that changing the value of this additional variable represents a change in the unknown state $f^0$. However, we have to define the (not directly observable) state $f^0$ to be constant to allow us to infer something about future repeated measurements[22]. All controllable aspects of the model are attributed to $x$, so controllable actions like replacing, moving, transforming an object to be measured are part of $x$. On the other hand variables which are themself stationary distributed just increase the noise and do not necessarily have to be included in $x$. Thus,

_____

[22] If there is a known dependency on time (external time or internal time corresponding to the measurement history like the number of repetitions of a certain measurement) then the time variable is part of $x$. A measurement at a certain time can be repeated if the time variable can be reset (e.g. to zero). If time cannot be reset repeating the same measurement is not possible. One can think of other restrictions, so that a measurement of the same question cannot be repeated. For example, if continuous questions $x$ are generated by a random process, the probability of repeating a measurement usually has measure zero. Those cases show the importance of information about relations between answers to different questions and not only between repeated identical questions, i.e. of nonlocal dependencies.

stationarity of $p(y|x, f^0)$ is a form of structural information.

In the next sections we will discuss possibilities of the inclusion of a larger number and variety of generalized data in statistical decision making processes (including classification and approximation) which are based on structural information. Those data correspond to prior information in the statistical language, rules in the logical language, interactions in the physical language. From logical (fuzzy set, Bayesian belief network) expert systems one can learn how to deal with a lot of different, heterogeneous rules and from physics the treatment of highly nonlocal (macroscopic) variables different from smoothness.

## 3 Generalized questions

### 3.1 How to write them

A generalized questions is defined by giving a set of required basis questions and defining a function $y^q = q(\vec{y}(\vec{x}))$ to be applied to their results. To simplify notation we skip the vector arrow and use the same letter $q$ for the defining function as well as for the functional itself, that is $q(y) = q_q(\vec{y})$. Generalized questions are fully defined by their answer distributions in the pure states $f^0 \in F^0$. As we always assume the model (of nature) to be in a possibly unknown but pure state the $p(y^q|q, f^0)$ represents the real processes and we write all formulas within this section for the $f^0$. This formulation without reference to some state of knowledge $f$ fits also into the Frequentist interpretation of statistics meaning that the definitions of generalized questions do not require to specify an explicit model of $f^0$ and are therefore not only meaningful for Bayesians. This holds as far as we make the definition of a question not explicitly $f$–dependent. Thus, in the following we form products for pure states $f^0$ and not for states of knowledge $f$,

$$p(y, y'|q, q', f) = \int df^0 p(f^0|f) p(y, y'|q, q', f^0)$$

$$= \int df^0 p(f^0|f) p(y|q, f^0) p(y'|q', f^0)$$

$$\neq p(y|q, f) p(y'|q', f),$$

including the case $q = x$ and $q' = x'$. The fact that $p(y, y'|q, q', f)$ has not necessarily to factorize enables learning. But notice, that we always assume the different parts of a question to be measured in the same pure state $f^0$. For $f^0$ we postulated independence, i.e. factorization

$$p(y, y'|q, q', f^0) = p(y|q, f^0) p(y'|q', f^0),$$

meaning that different questions $q, q'$ use different realizations of answers to $x$ and that learning being in a pure state, i.e. $f = f^0$, is not possible. In this formulation results depending on the same realization are seen as components of one question. On the other hand the following is easily adapted for a notation which allows such dependency between non basis questions.

A simple example is a question asking for the sum of two independent $x_i$–measurements with answer probability

$$p(y|q, f^0) = p(\{y_1 + y_2 = y\} | x_1, x_2, f^0)$$

$$= \int dy_1 \int dy_2\, \delta(y_1 + y_2 - y) p(y_1|x_1, f^0) p(y_2|x_2, f^0),$$

where $\delta(y_1 + y_2 - y^q)$ is the indicator function of the event $y_1 + y_2 = y^q$ and as always in this paper the $\delta$–functional has to be understood as Kronecker–$\delta$ function $\delta_{y_1,y_2}$ in the discrete case.

### Deterministic questions

The last example can easily be generalized by replacing $y_1 + y_2$ by any other *defining function* $q(y)$ depending on a data vector $y$ with components $y_i$ [23]

$$p(y^q|q, f^0) = p(\{q(y) = y^q\}|\, x^q, f^0)$$

$$= \int \left( \prod_i dy_i p(y_i|x_i^q, f^0) \right) \delta(q(y_i) - y^q).$$

To simplify notation, we now use the symbol $x^q$ not for a single basis question but for a whole vector of them. Its components $x_i^q$, correspond to the $y_i$ necessary to evaluate $q$ and can include repeated measurements of the same basis question. Specifically, averages of actual measurements correspond to choosing $q(y) = 1/N \sum_i y_i$. The answer $y^q$ to question $q$ can be a vector with components $y_j^q$ meaning that

$$\delta(q(y) - y^q) = \prod_j \delta(q_j(y) - y_j^q).$$

For noisy $x$ this gives not the product

$$p(y^q|q, f) \neq \prod_j p(y_j^q|q, f),$$

because this would refer to a situation where the $y_j^q$ do not use the same realization but sample their own $y_i$ resulting in additional $y$–integrals.

While in the following no special attention is paid to the technical difficulties of the case of continuous $y$ we shortly discuss the definition of the $\delta$ for this case where $\delta$ stands not for the Kronecker function but for the $\delta$–functional. This is defined for general $q(y)$ for $\left.\frac{dq(y)}{dy}\right|_{y_0} \neq 0$ according to

$$\int dy p(y|x, f)\delta(q(y) - y^q) = \sum_{y_0} p(y_0|x, f)\left|\left.\frac{dq(y)}{dy}\right|_{y_0}\right|^{-1},$$

where $y_0$ are the solutions of $q(y) = y^q$, i.e. the zeros of the argument $q(y) - y^q$. The nonzero first derivative guarantees that $q(y)$ is locally invertible so we can write $y_0 = q^{-1}(y^q)$ with $q^{-1}$ defined at least in a neighborhood of $y^q$. In case $q(y) = y^q$ is fulfilled on a whole interval the $\delta$–functional has to be replaced by the characteristic function of that event which can be expressed by step functions $\Theta$.

---

[23] The components $y_i$ of the data vector as well as the answers $y^q$ can be vectors itself.

### Output noise

A question can also be constructed including additional internal random variables $z$ by using a *probabilistic defining function* $q(y, z)$ defined by $f^0$–independent output distributions $p(z|y, q) = p(z|y, q, f^0)$. The answer can be seen as created in a two step process

$$p(y^q|q, f^0) = \int dy\, p(y^q|y, q)p(y|x^q, f^0),$$

with

$$p(y^q|y, q) = \int dz\, p(z|y, q)\, \delta(q(y, z) - y^q),$$

and

$$p(y|x_q, f^0) = \prod_i p(y_i|x_i, f^0)$$

The vector of internal random variables $z$ increase the noise in answers while additional (repeated) measurements of the same basis questions reduce the noise. Note that the noise variables of different questions are independent, that is we assume $p(z_q, z_{q'}) = p(z_q)p(z_{q'})$ for $q \neq q'$, and questions always measure their own $y$–values according to $p(y|q, f^0)$ and do not refer to measurements of other questions. One specific realization of $y$ can be used multiple times only within the same question.

### Input noise

The vector $x$ of basis questions contained in $q$ and its dimension can be an $f^0$–independent random variable. We call this a situation with *input noise*. Then $q(x, y)$ depends on $x$ generated according to $p(x|q)$. The index $q$ indicates that this distribution is part of the definition of the functional $q$. Thus, we have

$$p(y|q, f^0) = \int dx\, p(x|q) \prod_i^{\dim(x)} p(y_i|x_i, f^0),$$

with $\int dx = \sum_j \int \prod_i^j dx_i$ and $p(x) = p(\dim = j)\, p(x|j)$. We understand all dependencies within the vector $x$ to be included in this notation. Products of $p(y_i|x_i, f^0)$ correspond to a logical AND applied to the results $y_i$, summations over $x_i$ and $z$ correspond to a logical OR. Input noise can be combined with a probabilistic $q$–function $q(x, y, z)$

$$p(y^q|q, f^0) = \int dy\, p(y|q, f^0)p(y^q|y, q)$$

$$= \int dx\, p(x|q) \int \left( \prod_i^{\dim(x)} dy_i p(y_i|x_i, f^0) \right) p(y^q|x, y),$$

$$(1)$$

where $p(y^q|x, y) = \int dz\, p(z|x, y, q)\, \delta(q(x, y, z) - y^q)$ and $\int dy = \sum_j \int \prod_i^j dy_i$. Finally, we remark that input noise has both more active and more passive interpretations. Changing the variable $x$ can be interpreted as changing the device or changing the situation. For example, measuring another location $x$ of an object can be done by moving the measuring device to location $x$ (passive interpretation) or by moving the object (active interpretation) or by moving both (mixed interpretation). The passive interpretation may be the usual one, but as we assume $f^0$ to be constant over the whole time of interest all active controllable aspects must be part of $x$.

## Causal chains

We can also allow *output dependent input noise* by choosing the $x_i$ depending also from results of previous measurements $y_j$ included in $q$. For example, active learning algorithms which select the next question according to the data in the past belong to this class of questions. We order the factors according to the causal dependencies of the generating process like

$$p(x, y|q, f^0) = p(x_1|q)p(y_1|x_1, f^0)$$
$$\times p(x_2|x_1, y_1, q)p(y_2|x_2, f^0)$$
$$\times p(x_3|x_2, y_2, x_1, y_1, q)p(y_3|x_3, f^0) \cdots$$
$$= p(x|y_c, q) \prod_i p(y_i|x_i, f^0).$$

where the subscript $c$ indicates a 'causal' ordering, meaning $p(x|y_c, q) = p(x_1|q) \prod_{i=2} p(x_i|\{x_j, y_j\}_{1 \leq j \leq i-1}, q)$, with (the 'last') one of the components of vector $y$ missing. Note that we understand the causal structure for components of $x$ to be implicit in the notation $p(x|y_c, q)$. Chains with variables depending only on those of the previous step $p(x|y_c, q) = p(x_1|q) \prod_{i=2} p(x_i|x_{i-1}, y_{i-1}, q)$, are sometimes called Markov chains.[24] For finite sequences such a representation can always be achieved by combining $x$, $y$ of different steps into one vector variable, or in general by including the relevant memory variables. In the extreme case this would lead back to the starting point $p(x|q)p(y|x, q)$. Including internal noise variables $z$ their probability $p(z|x, y)$ may also be written in a causal realization modeling the real causal processes

$$p(x, y, z|q, f^0) = p(x_1|q)p(y_1|x_1, f^0)p(z_1|x_1, y_1, q)$$
$$\times p(x_2|x_1, y_1, z_1, q)p(y_2|x_2, f^0)p(z_2|x_1, y_1, z_1, x_2, y_2, q) \cdots$$
$$= p(x|y_c, z_c, q)p(z|x, y, q) \prod_i p(y_i|x_i, f^0).$$

It is always possible to write a joint probability in this form and we always understand the indices of the variables $x$, $y$, $z$ to refer to the same ordering. This gives

$$p(y^q|q, f^0) = \int dx \int dy \int dz \left( \prod_i p(y_i|x_i, f^0) \right)$$
$$\times p(x|y_c, z_c, q)p(z|x, y, q)\,\delta(q(x, y, z) - y^q)$$
$$= \int dx \int dy \int dz\, p(y|x, f^0)$$
$$\times p(x, |y_c, z_c, q)p(z|x, y, q)p(y^q|x, y, z, q), \qquad (2)$$

showing the separation into $q$-dependent factors written in index form from and $f^0$-dependent factors. Special realizations of causal dependencies, including variables $z$, defining the answer producing process can be modeled by directed acyclic graphs, i.e. graphical models or belief networks (Pearl, 1988; Lauritzen, 1996, Jensen, 1996, Ripley, 1996). A pair of variables $R_q = (x, z)$ can be called realization of $q$ because it represents that part of the definition of a question determining a specific answer $y^q$ and $D_q = (x, y)$ the corresponding data. We denote the set of all questions of the form of Eq.(2) including only finite dimensional vectors $x$, $y$, $z$ by $\bar{Q}_X^{fin}$.

---

[24]See for example Golden, 1986. However the term Markov chain is used in many variations. See for example van Kampen, 1992, p.77 and its footnote on p.89 referring to a footnote on p.340 in Feller, 1957.

## Combinations, decompositions

Having defined a set $Q'$ of questions $q'$ we obviously can generate new questions $q$ according to Eq.(2) by replacing $p(y|x, f^0)$ with $p(y|q', f^0)$ and choosing any input functions $p(q'|y_c, z_c, q)$, internal noise functions $p(z|x, y, q)$, and defining functions $q(x, y, z)$. To avoid circular definitions the $q$ must of course not contain itself. If the new $q$ is expressed in terms of basis questions, i.e. the $p(y|x, f^0)$, it still has the form of Eq.(2), meaning it belongs to $\bar{Q}_X^{fin}$. Consequently, for arbitrary $q' \in \bar{Q}_X^{fin}$ a $q$ of the following *general q-form*

$$p(y|q, f^0) = \int dq' \int dy \int dz\, p(q'|y_c, z_c, q)$$
$$\times p(y|q', f^0)p(z|x, y, q)\delta(q'(x, y, z) - y) \qquad (3)$$

is also in $\bar{Q}_X^{fin}$. In this sense $\bar{Q}_X^{fin}$ is closed. Also for every $q \in \bar{Q}_X^{fin}$ there are $q' \in \bar{Q}_X^{fin}$ so that $q$ can linearly decomposed in that form. Specifically, the $q'$ can always be chosen as $x$.

In general a decomposition of $q$ in arbitrary components $q'$ can be written as

$$p(y|q, f^0) = \int dq'\, p(q'|q, f^0)p(y|q', q, f^0),$$

with $f^0$-dependent $p(q'|q, f^0)$ if $q'$ is sampled depending on $y$. We can get a $y$- and therefore $f^0$-independent decomposition of $q$ into lower noise components $q^0$ in analogy to $x$-independent decomposition of states $f$ into pure states $f^0$. To get this, we separate $R_q = (x, z)$ into a $y$-independent part $q^0$ and a $y$-dependent part. Note that at least one $x_i$ is in $q^0$ and arbitrary dependencies are allowed within $q^0$. Then we have an $f^0$-independent decomposition of $q$

$$p(y|q, f^0) = \int dq^0\, p(q^0|q)\,p(y|q^0, q, f^0), \qquad (4)$$

in analogy to a $x$-independent decomposition[25] of $f$

$$p(y|x, f) = \int df^0\, p(f^0|f)\,p(y|x, f^0),$$

but with, in general, differing 'local states' $p(y^q|q^0, q, f^0)$ for different $q$. Despite that analogy we do not assume the availability of data which can change the $p(q^0|q)$. This is without loss of generality as all data dependencies of question definitions can be incorporated by enlarging the space $F^0$ and adding questions with answers depending on $p(q^0|q)$. In fact, in practice an example generating distribution $p(x)$ is often unknown and estimated using the sampling distribution. Formally we can combine $x$ and $y$ into a new $y' = (x, y)$ and define functions $f'^0 = (f_x^0, f_y^0)$ by $p(y'|x', f'^0) = p(y|x, x', f_y^0)p(x|x', f_x^0)$. The new $x'$-variable can be skipped from the notation having only one value with the minimal meaning of requesting another $y' = (x, y)$ value. A state of knowledge $f'$ contains now also information about $p(x)$ and if

---

[25]For simplicity we use from now on the integral notation also for the $f^0$-variable assuming a well defined (possibly discrete) measure.

$p(f'^0) = p(f_x^0)p(f_y^0)$ we can estimate $p(x)$ independent of $p(y|x)$ by observing pairs $(x, y)$. But note that for continuous $X$ if not parameterized by a finite set of parameters the $f_x^0$–integration is a functional integration and not necessarily well defined.

If studying $f$–dependency we may also allow $f$–dependent definitions of $q$ as the state of knowledge is assumed to be known and therefore this knowledge could be incorporated into the answer producing process. The preceding equations include this case if $q$ is interpreted as double index $q \to (q, f) = q_f$. Notice that as answers change the state of knowledge $f$ and they also change the definition of $q$.

**General functionals**

We allowed the answer probability $p(y^q|q, f^0)$ at $y^q$ of a generalized question to be any functional of the $y$–distributions $p(y|x, f^0)$ (for all $x$ but fixed $f^0$). A parameterization of, or in the extreme case the $p(y|x, f^0)$ for all $x$ and $y$ itself, describe a state $f^0$ completely. Functionals depending on the $p(y|x, f^0)$ include for example

$$p(y^q|q, f^0) = \delta(y^q - \int dy \, y p(y|x, f^0)),$$

giving the local expectation, or in a noisy version

$$p(y^q|q, f^0) \propto e^{-\frac{1}{2}((y^q - \int dy \, y p(y|x, f^0))/\sigma)^2},$$

as well as

$$p(y^q|q, f^0) = \delta(y^q - p(y|x, f^0)),$$

giving a specific probability (density) or general

$$p(y^q|q, f^0) = \int dy p(y|f^0)\delta(y^q - q(z)).$$

We recognize that the questions we have constructed so far in the previous sections contain products of $p(y|x, f^0)$ (together with a sum implicit in the notation of the last sections) and are therefore functionals with a finite power series

$$p(y|q, f) = \sum_{n=0}^{n} \left( \prod_{i=1}^{n} \int dx_i dy_i p(y_i|x_i, f) \right) a_n(x, y).$$

Those generalized questions can be measured using a finite number of answers to basis questions, while questions with an infinite but converging power expansion might be measured approximately. Expectations and probabilities of events can sometimes be approximated using empirical sums but in general not a probability density $p(y|x, f^0)$ for continuous $y$. We will use the word measurement for the process of getting an answer to a question, but we will show in the next section that for questions with no finite or converging power series this has more the character of enforcing an answer or active control of $F^0$.

## 3.2   How to measure them

An answer to a question could be obtained using only measurements of basis questions $x$ when it has a defining function $q(x, y, z)$ depending only on a finite number of outcomes $y_i(x_i)$.

### 3.2.1   Finite case

We denote the set of questions with finite dimensional $y, z$ by $\bar{Q}_X^{fin}$. With $Q$ denoting the measurable questions we have therefore $X \subseteq Q \Rightarrow \bar{Q}_X^{fin} \subseteq Q$. These questions could (but do not have to) be measured by a finite number of possibly repeated measurements of basis questions $x \in X$ according the following steps:

i. Choose $x$ according to $p(x|q)$,

ii. Get $y$ for the basis questions in $x$,

iii. Get $z$ according to $p(z|x, y)$, (Repeat i., ii. and iii. in case of $p(x|y_c, z_c)$ dependencies)

iv. Insert result into $q(x, y, z)$.

While the measurements can be performed using only devices measuring basis questions, these questions do not necessarily need the full information contained in the answers to $x$. That means measuring devices designed specifically for them could be more effective. For example, using interference of waves differences can be measured in physics sometimes much more precisely than absolute values.

We discuss how in the finite case a direct measuring of $q$ might be preferable, a situation which often occurs in practice. Consider output noise which is added in a final step for all measurements during training, like observation and memory errors, or transformation to a lower resolution final output scale. Then measuring devices which are able to directly access the underlying lower noise function and perform the necessary operations before adding the output noise have higher accuracy. As simple example, take a basis set $X$ with output distributions

$$p(y|x, f^0) = N(\mu_x(f^0), \sigma)$$

and

$$\mu_x(f^0) = \int dy \, y \, p(y|x, f^0),$$

where $N(\mu, \sigma)$ stands for a Gaussian centered at $\mu$ with variance $\sigma^2$. Then, if measurable, the generalized question

$$p(y|q, f^0) = N(\mu_q(f^0), \sigma), \ \mu_q = \mu_{x_1}(f^0) + \mu_{x_2}(f^0),$$

obtains the sum with greater accuracy than using the sum $y_1 + y_2$ of two basis questions $x_1$ and $x_2$ where the independency of the noise gives $2\sigma^2$ for the variance. In this special case the question would be equivalent to four basis questions including one repeated measurement for every $x_i$ to get $(y_1 + y_1' + y_2 + y_2')/2$. Replacing the sum $\mu_{x_1} + \mu_{x_2}$ by an integral $\int \mu_x dx$ (infinite nonlocality) or setting $\sigma$ in $q$ equal to zero with retaining a finite $\sigma$ for the $x_i$ (infinite accuracy) the question depends on more than a finite number of answers to basis questions.

### 3.2.2   Infinite cases
**Asymptotic procedures**

In some cases a defining function can be found so that a well defined limit with the number of arguments going to infinity gives the desired $p(y^q|q, f^0)$. For example, there are measurements depending on a formally infinite number of basis questions using some inherent, infinite

16

parallelism of natural processes or a process having a well defined limit for the number of involved basis questions going to infinity. Scattering of waves on structures with specific translation and rotation invariances (crystals) create filters for specific Fourier components which depend on a (conceptually) infinite number of coordinate values. Measurement of macroscopic variables in physics, like energy or magnetization depends on many (normally in the order of $10^{23}$) microscopic variables. If the answers to those questions converge in the limit system size $n \rightarrow \infty$, they can be considered as measurement of an infinite system. In general, a meaningful statement about an infinite property, has ro rely on conditions which can be checked in finite times. In the following we relate such conditions to the preparation of an ensemble and possibilities of control.

## Measurement and preparation

We begin with the observation, that questions available during the training situation might be different from those during the preparation of the ensemble $f$.

Consider objects $f^0$ which, giving a parameter $\bar{y}$, produce as random output an output distribution with mean $\bar{y}$. We can form a population $F^0$ of such objects, notate, i.e. measure, the input parameter $\bar{y}$ for each member, and select, with a given (prior) probability $p(f^0) = p(\hat{y})$, one object $f^0$ with unknown $\bar{y}$. Then, we cannot measure the mean exactly using only training example $y_i$ of the output of $f^0$. In cases the mean exists, the sampling sum will, according to the theorem of Glivenko–Cantelli, usually asymptotically converge to it, (i.e. the probability $\delta$ for the deviation to be larger than some $\epsilon$ can be $n$–dependently bounded). Thus, what have been measurable during the preparation phase by a single measurement (i.e. reading the number $\bar{y}$) is, if at all, only asymptotically measurable during training.

In a general situation we consider a family of parameterized answer distributions $F^0$ with known parameters like expectation and variance. If we can measure the values of these parameters we can prepare a prior distribution (state $f$) depending on those parameters. For example, if a certain symmetry property is measurable we can choose an $f$ depending on this property, for example with possible states $f^0$ possessing (or concentrated around) a certain value of this symmetry. In a specific training situation, however, the symmetry may not be directly measurable, and an additional source of noise be present.

We will denote the set of questions with available answers during preparation by $X^P$ and will call it the *preparation set*. The preparation set $X^P$ is in general not equal to the set of training question $X$ considered actually available. Questions not finite or asymptotically measurable with respect to $X$ can be finite or asymptotically measurable with respect to $X^P$.

As far as generalization requires nonlocal questions not measurable with respect to $X$, like a symmetry or smoothness constraint for infinite $X$, there has to be another set $X^P$ to allow the necessary measurements to prepare the prior. Then the prior can be expressed by a question $q^P$ depending on a finite number of questions

$x^P \in X^P$ and the corresponding answer $y^q$

$$p(f^0|y^P, q^P) \propto p(f^0)p(y^P|q^P, f^0).$$

However it could not be written using a question depending on a finite dimensional vector $x \in X$ of training questions. The training data change the prior according to

$$p(f^0|y^D, q^D, y^P, q^P) \propto p(f^0|y^P q^P)p(y^D|q^D, f^0).$$

Preparation questions $x^P$ with deterministic answers, like an exact symmetry, correspond to restrictions of $F^0$. Probabilistic $q^P$ give probabilistic priors, like a preference for smooth functions or otherwise symmetric functions (approximate symmetries).

The distinction between a (momentarily considered) training set $X$ and a preparation set $X^P$ (on may say implicit training set) is more a practical than a formal aspect, For analysis of the generalization ability questions from $X^P$ have to be treated formally equal to those from $X$. Thus, effectively one deals with training questions depending on $X \cup X^P$.

Specifically, we will explain below how symmetry and smoothness can be generated by input noise or averaging. To recognize or prepare such a situation a measurement device without (or less) input noise ($X^P$) must be available.

We summarize, that what appears to be not measurable with a finite amount of data for $X$ may well be measurable within the larger set $X \cup X^P$. This is the case, when the set of actually considered training questions $X$ all share a common noise source, which is absent for $X^P$.

## Measurement and control

To enable learning we have to assume stationarity of $p(y|q, f^0)$. Thus, all factors changing the answer distributions have to be included in $q$ and $f^0$. But the model does not specify how stationarity is achieved. In practice, stationarity can result from an active control or just by not disturbing a constant part of nature.

In general, a measurement $y$ of $q$ in state $f^0$ is the result of an interaction of the 'active' part of posing a question $q$ and the 'passive' reaction of nature in state $f^0$. Usually, measuring a quantity $q$ emphasizes the passive picture where the value of $y$ reflects a permanent property of nature $f^0$, however, only seen when measuring $q$. For preparation questions the complementary active interpretation of measuring as control or selection may also be helpful. Then, the answer is seen as a property of nature enforced by the question $q$. In this interpretation a question is better called a control action and a measurement device a control device. Thus, $y$ can be seen as reaction of nature to $q$. In the general case of stochastic control different states $f^0$ react different to the control action $q$, described by $p(y|q, f^0)$. This point of view is especially suitable if the variability of $y$ between different $f^0$ is small under $q$. Both interpretations describe the same formalism, the difference being the larger emphasis of either the passive or the active part.

Interpretations may also refer to a more complicated model of interaction between $q$ and $f^0$. Such more complex models correspond formally to the introduction of hidden variables $z$,

$$p(y|q, f^0) = \sum_z p(y|q, f^0, z) p(z|q, f^0).$$

That means we think of measurement devices as generalized questions with respect to some underlying set of measurable questions. Then, the $x$ can be called effective questions $x^{eff}$ with respect to another, underlying $X$. If we do not want just to assume such a structure, we need other measurement devices, we may call them in analogy to the last Subsection $q^P \in \bar{Q}_{X^P}$, which allow to measure or control a given structure. We will see that this can be reasonable to assume and those additional measurement devices will have to be active only a finite number of times. If we can guarantee the structure for all $q$, this can be equivalent to dependencies between $q$ and therefore interpreted as nonlocal measurement.

For example one can imagine the value $y$ to be the result of an (ideal) measuring process with a following control device, e.g. a filter describing restrictions of the (real) measurement device. For this we separate the index $q$ into two components $q = (q^{ideal}, q^{filter})$ (or alternatively $f^0 = (f^{0,data}, f^{0,filter})$) and have

$$p(y|q, f^0) = \sum_{y^{ideal}} p(y|q^{filter}, y^{ideal}) p(y^{ideal}|q^{ideal}, f^0).$$

We will call this a model of *posterior control*.

Analogously, we can think of a scenario where first the functions are filtered and then measured, i.e. we split $q$ into $q = (q^{data}, q^{filter})$ (or alternatively $f^0 = (f^{0,data}, f^{0,filter})$) and have

$$p(y|q, f^0)$$
$$= \sum_{f^0_{filtered}} p(y|q^{data}, f^0_{filtered}) p(f^0_{filtered}|q^{filter}, f^0).$$

If we $f^0$ from the notation, write again $f^0$ for $f^0_{filtered}$, the filter $q^{filter} = (y^P, q^P)$ creates a prior ensemble. Here the $q$ independence allows us to have it prepared before the training starts. That is what we will call a model of *prior control*. The situation we discussed in the last Subsection can therefore be realized as prior control. Notice, that also under prior control stationarity of the probability distributions must be controlled, and in this sense there is always a posterior control component present. Adjusting the measurement device and repeat the measurement if some 'failure' is indicated is an example of such an posterior control, which is also present, if the prior ensemble is prepared in the past.

Now we turn again to the problem of questions depending on an infinite number of basis questions. In the last Subsection we used a decomposition $q' = (q, q^P)$ with $q \in \bar{Q}_X^{fin}$, $q^P \in \bar{Q}_{X^P}^{fin}$ and showed that this allows $q^P \notin \bar{Q}_X^{fin}$. Now we examine how also in the example of posterior control measurement of questions depending on an infinite number of basis questions is possible.

The key observation is that the stationarity condition for $p(y|q, f^0)$ only affects an always finite number of measurements. Hence. also (e.g. prior or posterior) control can be restricted to this finite number of measurements. This allows to understand or implement stationarity as part of the measuring process, e.g. as a filter. For example, nothing prevents us to assume that asking $q$ the first time causes stationarity for the following times. Therefore, posing a question $q$ can have the interpretation of an active control leading to some restrictions or dependencies for subsequent measurements, stationary for all $f^0$. Differentiation between different $f^0$ is only possible if the controller, but not necessarily the learner, has questions, $q \in \bar{Q}_{X^P}^{fin}$, available to distinguish between them.

As example, take a family of $q$ measurement devices, all capable only of producing answers smaller than $\theta$, independent of $f^0$. This can be ensured by one $q$–independent cutoff device as posterior control. This cutoff device depends itself only on finite dimensional incoming $y^{ideal}$ and has to be active only during the finite number of measurements. Nevertheless, the posterior control model is equivalent to a nonlocal measurement for *all* $q$, which is possibly an infinite number, with answer $q < \theta$. Here, the control is implemented by $f^0$–independent properties of the measurement device. The same effect appears if all objects under study underlay the same $f^0$–independent selection. So instead of using restricted measurement devices as posterior control, we could also study restricted objects or situations, i.e. an implementation as prior control. Then the controlling filter acts using $q^P \in \bar{Q}_{X^P}$ not on measurement values but on the objects or situations $f^0$ itself. In analogy to the output bound $\theta$ above, the length of object classes under study might be restricted by $\theta$ because they all arrive in the same box, must fit into a certain environment or are only produced in that way. Those are prior control devices parameterized by $\theta$. Also, if the number of objects is infinite, the prior control device (e.g. the box into which the object must fit) must only be active during the finite number of measurements.[26]

Thus, because control of stationarity, be it related to

---

[26]There can be practical differences between such devices, for example between a 'box' device, and a cutoff in the output scale. Indeed, a filter can be implemented as active or as passive filter. A passive filter does not always return an output. For values larger $\theta$ for example it may answer 'overflow'. For example, using the box device it may take considerable time to find an object which fits into it. There might even be no objects available which fit in the box, and we could need a formally infinite time to produce such an impossible state. An active filter always returns an output. For example we may assume for a cutoff device that it always returns $\theta$ in case the output is larger $\theta$. However, this has nothing to do with the difference between prior and posterior control. Both variants are also possible for prior control. Consider, for example, an 'active box' cutting everything to a fitting size Those aspects of differing complexity of single measurements or control actions are not included in the formalism. Yet, nothing prevents us from modeling the micro-structure of single measurements if necessary; we could use the same type of theory, just on another level.

training or preparation questions, only has to be active at the always finite number of times of actual measurements, there is no practical impossibility in measurements based upon control depending on an infinite number of basis questions. Nothing is actually done an infinite number of times.[27] It just could be done arbitrarily often because we defined the situation to be so. A filter bounding outcome values only has to be active a finite number of times for every actual measurement even if we defined the situation that it could do so arbitrarily often. This includes also the stationarity conditions for local questions.

We have seen that what exactly choosing question $q$ means in practice and how stationarity is achieved, be it by leaving nature alone or by active control does not enter the formalism. Dependency on an infinite number of questions is no practical impossibility in as far as control over stationarity conditions only has to be active at the always finite number of times of actual measurements.

Another aspect of control will be discussed in more detail in the next Sections: Often the presence of control is easy to recognize, however difficult to formalize. For example, creating a training set for a object recognition task by choosing images as training examples of faces and non–faces or chairs and non–chairs, depends on one's implicit definition of the concept of a face or chair, respectively. Drawings may be accepted as valid examples of the object class or not, while very unregular, random–like objects are not selected to represent a face or a chair. Consequently, the definition of such object classes is related to linguistic or implicit concepts representing the objects. Such defining concepts, involved for example in a prior or posterior control process, can correspond to a measurement of an infinite number of objects. Thus, it can be expected to be helpful having a method to formalize those concepts and include them as prior information into the statistical inference process. Accordingly, useful restrictions might be found from an analysis of the application situations or of the object or situation generating process: a pedestrian detection task might be restricted to pedestrians walking on or near the street in a certain distance and a car (ship, airplane, $\cdots$) detection task may take into account that only certain types of them have been produced up to now.

As man are always more or less involved in the definition of situations or problems of interest this shows the clear need of a *human interface* which enables human knowledge, like linguistic concepts, used to define and control the application situations to be incorporated into the learning process.

Summarizing, we list two variants in which measurement by control depending on a possibly infinite number of questions can appear:

1. prior control: controlling the process which generates the prior distribution, e.g. control of the situ-

ations or objects of interest,

2. posterior control: controlling the measurement values, e.g. restrictions of the measurement devices.

Stationarity is essential for all questions and the interpretation of measurement as control or active enforcing of stationarity can be applied to all questions, not only to those which cannot be measured using a finite number of basis questions or some asymptotic procedure. One can say that the ability to generalize is based upon the ability to control that the situation of interest and the measuring devices are kept within their restrictions. We have discussed, that the probability concept does in general not allow to exactly verify (even local) nondeterministic conditions using only the training data and no preparation questions. Approaches to test the conditions, i.e. the models $f$, using training data have to refer to meta models or use the classical testing of null hypotheses, i.e. calculate the probability of the data given the condition $f$. Applied to model testing this has also been called evidence approach (MacKay, 1992c).

We conclude, that the ability to generalize is based on the ability to measure or, emphasizing more the active connotation and formally equivalent, to control dependencies between basis questions. We have shown that this is also practically possible in case of an infinite $X$. From this point of view only measurement can generate new information. Actual learning, however, is not the discovery of something new, but the reformulation of actually imposed and measured conditions to give answers to relevant questions. Analogously, assumed learning is the reformulation of assumed conditions.

## 4 Priors

We discuss two examples of generalized questions often used as priors and give then in the next Section a general method to construct priors.

### 4.1 Bounds

In practice measurement values have a bounded range and even models using distributions with unbounded range for the variables use normally a bounded range for moments, like the mean. To discuss such bounds we present some variations of questions with $q$–functions calculating the maxima of different sets:

1. empirical maximum from a finite sample
$$m_1 = \max_i y_i(x_i, f^0),$$

2. maximum of the local expectation (regression function)
$$m_2 = \max_x \int dy\, y\, p(y|x, f^0) = \max_x \bar{y}_x(f^0),$$
allowing observed values $y_i > m_2$ generated by noise,

3. outcome (answer) being maximal with probability one
$$m_3 = \max_x y^*(x),$$

[27]This is related to a constructive view of infinity, not attributing 'existence' to an abstract infinite object itself, but to its constructing procedure. This is, for example, the position of Jaynes, 1996, who 'sails under the banner of Gauss, Kronecker, and Poincaré rather than Cantor, Hilbert, and Bourbaki.'

defining as local maximum $y^*(x)$ at $x$ the minimal $y$ with

$$\forall y' \in Y_x \ : \ p(y' > y|x, f^0) = 0,$$

or

4. maximal potential outcome

$$m_4 = \max_x \max_{y \in Y_x} y$$

being independent of $f^0$ if the $Y_x$ are defined independent of $f^0$.

Questions related to $m_1$ belong to $\bar{Q}_X^{fin}$ and can be empirically measured if $X$ can. In approximation problems where one is interested in modeling the regression function $\bar{y}_x(f^0)$ and not $f^0$ itself one normally refers to $m_2$ and allows for example Gaussian noise still to generate arbitrary large outcomes even for finite $m_2$. The answer to $m_2$ can be bounded only with access to the regression function, which in this case is interpreted as the true underlying function. This 'bounding device' has therefore to be applied before the measurement noise. The bound $m_3$ is the most interesting one in worst case considerations and is itself bounded by $m_4$. Fixing (the answer to question) $m_4$ is done by using cutoff devices. A question $q_A^{cut}$ corresponding to a cut-off device applied to (real) answers to question $q$ can be written

$$p(y^{cut}|q_A^{cut}, f^0) = \int dy \, [p(y > A|q, f^0)\delta(A - y^{cut})$$

$$+ p(y \le A|q, f^0)\delta(y - y^{cut})]$$

$$= \int dy \, p(y|q, f^0)[\delta(\Theta(y - A) - 1)\delta(A - y^{cut})$$

$$+ \delta(\Theta(A - y) - 1)\delta(y - y^{cut})].$$

The step function $\Theta(x)$ is defined to give 1 for $x \ge 0$ and 0 for $x < 0$. If only those cutoff questions are available we can restrict to effective states $f_{eff}^0$ defined by

$$p(y^{cut}|q_A^{cut}, f^0) = p(y|q_A, f_{eff}^0),$$

. We need this device for the next example.

## 4.2 Approximate invariances: Smoothness and symmetry

Smoothness and symmetries are probably the most important and most often used nonlocal priors. Here we want to show how restrictions of the measurement device (or equivalently restrictions in the situation (object) generating device) can lead to bounds on smoothness and symmetries.

Let us assume a group $S$ represented by operations $s_x$ acting on $x$. For simplicity we skip the index $x$ if not necessary and write simply $s$. Then for scalar states $f^0$ the action of operation $s$ is defined by requiring invariance $p(y|x, f^0) = p(y|sx, s^{-1}f^0)$ that is

$$p(y|x, sf^0) = p(y|s^{-1}x, f^0).$$

Generalization to states with components not invariant under $s$, like for example vector states, is straight forward

$$p(y|x, sf^0) = p(s_y y|s^{-1}x, f^0),$$

where $s_y$ is the action of $s$ in the $y$–representation of the group, which may be different from the $x$–representation $s_x = s$.

Measuring invariance of the regression function under a set of operations $S$ can for example be done by calculating a mean square error (weighted with $w(x, s)$ if necessary) writing $E(y(f^0, x)) = \bar{y}_x$

$$d_S^2(f^0) = \int dx \int ds \, w(x, s)(\bar{y}_x - \bar{y}_{sx})^2,$$

where $s$ denotes also the index of operator $s \in S$ and the integral notation valid for continuous (Lie) groups has to be replaced by a sum for finite groups. (See for example Ferraro, 1992, for Lie groups in pattern recognition). A distance $d_S$ can be used to construct priors, for example like $p(f^0) \propto e^{-cd_S^2(f^0)}$, if normalisable in $F^0$, or $p(f^0) = \Theta(d_{max}^2 - d_S^2(f^0))$. A relative difference is obtained using $w(x, s) = (1/(x - sx))^2$. One could also measure higher order differences like $((\bar{y}_x - \bar{y}_{sx}) - (\bar{y}_{sx} - \bar{y}_{s^2x}))^2$ corresponding in the infinitesimal case to second derivatives. Measuring smoothness is the special case of measuring infinitesimal translational invariance $s_0$. For example, written for the one–dimensional case with $w(x, s) = \delta(s - s_0)\frac{w(x)}{(x - s_0x)^2}$

$$d_S^2(f^0) = \int dx \, w(x)(\frac{d\bar{y}_x}{dx})^2.$$

The $x$–integration $p(x)dx$ can be written in the form $\sum_{i \in orbits} p(s'x_0^i)ds'$, with the orbit $i$ of $x_0^i$ defined as $Sx_0^i$ and the sum is restricted to one $x_0^i$ per orbit. If there is only one orbit, all $x$ can be generated as $x = s_x x_0$ out of one $x_0$ by repeated applications of the infinitesimal $s_0$ with $s_x = e^{xs^o}$ denoting the corresponding finite transformations.

For discrete symmetries integrals become sums. Examples of discrete symmetries include permutation of components in case $x$ is a vector. Function spaces of functions depending on vector arguments $x$ can be constructed as tensor product of function spaces depending on the single components of $x$. If every component of $x$ corresponds to a measurement of another object, exact permutation symmetry means indistinguishable objects.[28]

Now, let us assume that there is *input noise* or *averaging* associated to operations $s \in S$ of some group $S$. In spatial systems this is often also called *coarse graining* (see for example Balian, 1991, Goldenfeld, 1992). This is a very natural assumption for smoothness or infinitesimal translational invariance, as no real measurement device has infinite resolution. If only questions with input noise with respect to $S$ are available we can define an effective state $f_{eff}^0$. Including the identity into the set $S$ it is defined by

$$p(y|x_{eff}, f^0) = p(y|x, f_{eff}^0) = \int ds \, p(s|x)p(y|sx, f^0),$$

with the input noise characterized by $p(s|x)$, that is the probability of posing question $sx$ instead of question $x$.

---

[28] For example, in physics identical particles like bosons are related to an exact permutation symmetry.

20

The following is an example of averaging with respect to $S$ by the same weight factor $p(s|x)$

$$p(y|x_{eff}, f^0) = p(y|x, f^0_{eff})$$

$$= \int dy_s \prod_s p(y_s|sx, f^0)\delta(\sum_s p(s|x)y_s - y).$$

.

Now, we are interested in bounds on the approximate symmetry of the effective regression functions

$$\bar{y}^{eff}_x = E(y(f^0_{eff}, x)) = \int dy\, y\, p(y|x, f^0_{eff})$$

$$= \int dy \int ds\, y\, p(s|x)p(y|sx, f^0)$$

$$= \int ds\, p(s|x) E(y(f^0, sx)) = \int ds\, p(s|x)\bar{y}_{sx},$$

for the input noise version. The average versions give the same expectation but smaller variance. The result is a convex combination of the local averages $\bar{y}_{sx}$. Now we consider a measurement device with finite range (for its real valued components) of output and define according to the last section a new effective state which includes the cutoff with respect to an upper bound $A_x$ and lower bound $B_x$. Then the effective regression function can only take values between the extremal points $A_x = \min_s E(y(f^0, sx))$ and $B_x = \max_s E(y(f^0, sx))$ and by changing $p(s|x)$ we can obtain any value in between.

Analogously, differences are bounded. If we take for simplicity $p(s|x) = p(s|x') = p(s)$, we have

$$|\bar{y}^{eff}_x - \bar{y}^{eff}_{sx}| = |E(y(f^0_{eff}, x)) - E(y(f^0_{eff}, sx))|$$

$$= |\int ds'(p(s') - p(s's^{-1}))\bar{y}_{s'x}|,$$

for a parameterization $s'$ with $\frac{ds}{ds'} = 1$. Therefore, changing $p(s)$ allows to obtain any bound $d_{eff}$ for the norm of the difference between $d^{max}_{eff} = \max_x \max((A_x - B_x), (B_x - A_x))$ and $d^{min}_{eff} = 0$ (for $p(s) = $ const.) and we can achieve

$$d^2_S(f^0_{eff}) = d^2_S(f^0_{eff})\Theta(cd^2_{eff} - d^2_S(f^0_{eff})),$$

with $\int dx \int ds\, w(x, s) = c$. This bounds the smoothness or deviation from perfect symmetry by $0 \leq d_S \leq cd_{eff}$ for any $0 = d^{min}_{eff} \geq d_{eff} \geq d^{max}_{eff}$.

Thus, input noise or averaging in connection with a cutoff device can lead to dependencies between answers of an effective function by bounding their maximal differences, which is a symmetry or smoothness property.[29]

---

[29] For example, the support vector machine (Vapnik, 1995) applied to classification problems can be seen as such an approach. Here the input space is embedded in a (often much) higher dimensional feature space. The classification in feature space only requires the calculation of scalar products which are defined through a positive definite kernel, chosen so that it is easy to calculate, In feature space a linear separating hyperplane is constructed maximizing its distance to the nearest data points, which are also called support vectors.

For sampling of symmetry priors by virtual examples and references see Section 8.4.

In real measurement devices these combination of cutoff with input noise or averaging (coarse graining) can occur on many different levels and it is an interesting question whether the omnipresent smoothness phenomena in nature could be partly explained in that way.

# 5 Subjective priors

## 5.1 General priors: How probabilistic models are obtained

The preparation of an ensemble $F^0$ with a certain prior distribution requires measurement or control of certain properties of $f^0$. States $f^0$ are defined in terms of a parameterization of $p(y|x, f^0)$. Any prior $p(f^0|f^P)$ describing a state of knowledge (i.e. state of preparation) is therefore a deterministic functional of those $p(y|x, f^0)$, and also of their parameterization. Thus the set of $p(y|q, f^0)$ itself are answers to the maximal set $X^P_{max}$. The $x^p_{max}$ are the deterministic functions giving $p(y|x, f^0)$. This means the number $p(y|x, f^0)$ for given $x$, $y$, $f^0$, not the probabilistic questions $x$ giving answer $y$. The preparation questions $X^P_{max}$ allow to construct a set of well defined $p(y|x, f^0)$ producing devices. Usually then one of these devices is selected according to some $p(f^0)$, and one has to find out using new data $D$ which one is actually chosen. With respect to $X^P_{max}$ every prior, even if defined directly by a deterministic functional of the $p(y|x, f^0)$ can be reinterpreted (but not in a unique way) as resulting from a uniform prior $p(f^0)$ with additional given data $(y^P, q^P)$. We take the point of view that every nonuniform prior is caused by such data $D$, sometimes also denoted by $D^0$ if we want to distinguish them from other training data $D$.

Given data $D$ the corresponding state $p(f^0|D)$ can be calculated if $p(D|f^0)$ is known. The $p(D|f^0)$ are part of the structural knowledge. They must be determined independent of the actual task under consideration. Thus, their knowledge is always a transfer of knowledge from another task and assumes constancy of this distributions

---

Choosing the separating hyperplane with maximal distance to the nearest sample points is equivalent to maximizing the input noise around the sample points without changing the classification. The cutoff consists in the restriction to data within a certain radius. One may interpret the class boundaries resulting from equal class priors and a class membership probability of the form of a mixture of Gaussians centered at the data points, radially symmetric in feature space with respect to the distance induced by the selected kernel, and with equal variance $\sigma$. From this point of view one can say that the support vector machine obtains a solution with a maximal 'smoothness' of the class membership probability with respect to kernel induced distance in feature space, i.e. with maximal $\sigma$. Thus, the support vector machine implements a smoothness prior relative to the feature space. One may remark here, that the related VC dimension of the support vector machine can (up to now) not be calculated exactly, because the related function space $\hat{F}$ of optimal hyperplanes is not defined a priori but dependent on the $x$−values of the training data. (See Shawe-Taylor, Bartlett, Williamson, Anthony, 1996ab and their concept of a luckiness function.)

under transfer. Up to now we assumed the $p(D|f^0)$ to be given. Then the preparation of an ensemble $p(f^0)$ can be related to measurement devices with already known (empirically measured) answer probabilities under the various states. Here we discuss the measurement of $p(D|f^0)$. The main problem is, that we have to measure $p(D|f^0)$ in another situation, and therefore have to ensure its constancy (or make its approximate constancy plausible) to allow transfer.

Measurement or control 'devices' include humans e.g. an ensemble is prepared under human control or described by verbal statements. So, images of a training set may be labeled 'chair' or 'non–chair' by some 'expert in chairs'. To allow any meaningful generalization we must at least approximately (maybe implicitly) transform its concept of a chair into a chair approximator, applicable to all possible images, which then can be improved by training examples. Of course, one may rely on some smoothness condition or other implicit model restrictions (corresponding to a special chair approximator for all possible images) and the training examples alone. But obviously, besides the training examples, any verbal description of a chair also adds to the available information and one may have good reasons to believe that there is a better 'universal first guess chair approximator' than just using the implemented smoothness and other implicit model restrictions. Even if many training examples are available, it can help to use verbal descriptions to create also 'virtual' examples (which must not be an infinite set, assuming that the implicitly or explicitly implemented smoothness conditions interpolate in the neighborhood), because the virtual examples can include data which are not available as training data. They may teach the concept 'chair' far more effectively, than images of real chairs. Thus, we want to find a reliable relation of answers distributions of experts to possible states $f^0$, i.e. approximate the expert answer probability $p(y^E|q^E, f^0)$. This is obviously a very complicated task and as well a subject of psychological research as of statistics.

More general, one may even take the point of view that human experts are always involved: They have to describe 1. the single data probabilities $p(D|f^0)$ in an, already more precise or still less precise, verbal form, and 2. the dependency of the prior on the various data. E.g. single data have to be combined by AND, OR, or more complicated operations. These operations depend on the dependency structure of the single data. One must find a procedure to translate the related verbal information into numbers. This may be rather trivial, if the verbal statement is a reference to a (maybe empirically obtained) numbers. However, often this is not quite as easy. Assume we have data $(q, T_{y_q})$ available from some approximator $T$ of $p(y|q, f^0)$. We might know that it has been trained on some examples. But what would be $p(T|f^0)$ and $p(y_q|T)$? This clearly depends on the learning algorithm, the approximator was using, as well as on $f^0$. We may not know the details of the algorithm, and determining $p(T|f^0)$ consistently to our model of basis questions may well be a much harder task than solving our actual problem. Thus, we have to perform

an approximate 'internal integration over (hyper-)priors' by considering our experience with this and similar approximators in similar situations. The result will be a verbal statement describing a subjective concept related to $p(T|f^0)$, e.g.: 'We may trust the results of this simple approximator, not too much but a little in all situations which are not too similar to the examples $A$ or $B$'. This has to be changed into a numerical representation, which includes e.g. translation of 'not too much', '$A$', '$B$', 'similar to $A$, $B$', IF ... THEN, OR, AND. Still one can expect this to be in general a better solution than just ignoring unprecise information: Generalization requires implementation of nonlocal dependencies, and trying to match those nonlocal dependencies to unprecise verbal descriptions seems better than using nonlocal dependencies which are implicitly implemented without any further reasoning. Of course, there may be individual situations where the unknown implicit assumptions are better suited to the problem. But if this occurs regularly, the specific method used to include unprecise information has and can be adapted.

In principle, every prior could be seen as resulting from one combined prior question, and therefore be constructed like any generalized question. However, especially for priors, but possibly for every question, the necessary ingredients have to be constructed from more or less unprecise information in verbal form. We therefore shortly comment the standard operations AND and OR for probabilities from this point of view. Consider a prior, prepared by applying preparation question $q^P$ and enforcing or measuring an empirical answer $y^P$. One obtains

$$p(f^0|f^P) \propto p(f^0)p(y^P|q^P, f^0).$$

According to the construction of generalized questions we can translate our verbal statement that the data ($q_i^P$, $y_i^P$) should be combined by an AND for conditional independent events by forming the product

$$p(f^0|f^P) \propto p(f^0)\prod_i p(y_i^P|q_i^P, f^0).$$

But we might not be to sure about their independence, without knowing any explicit dependency structure. Again we may delegate the task to an expert and translate its verbal output either in 1. conditional probabilities determining the dependency between $y^p$ and apply the correct probability theoretical formula, or 2. directly into an adapted combination of the $p(y_i^P|q_i^P, f^0)$ which might well look different from the above formula for independent events.

Similarly, a partial sum (or integration for densities) correspond to a (weighted) OR would translate statements, referring to measurement noise for disjunct events:

$$p(f^0|f^P) \propto p(f^0)\sum_i p_i^0 p(y_i^P|q^P, f^0),$$

or

$$p(f^0|f^P) \propto p(f^0)\sum_{i,j} p_i^0 p(q_j^P)p(y_i^P|q_j^P, f^0).$$

The following statement 'The result was either 1 or 7, I am not sure.' would fit quite well into that category, provided we can relate the $p_i^0$ to subjective beliefs. Again,

instead of guessing $p_i^0$ and the dependency structure and applying the correct probability theoretical formula, one may translate a verbal statement directly into $p(f^0|f^P)$ and the resulting formula could look differently.

In general, instead of constructing the ingredients from verbal statements and then inserting them in probability theoretical formulas, one may also directly construct the result from unprecise knowledge. The results will differ if the subjective concepts do not represent probabilities. This is the case if the dependency structure is easier formulated in a verbal form than in terms of conditional probabilities. On the other hand, one may try to improve verbal statements, by enforcing experts to use probabilistic models. In all cases we start with unprecise verbal information and require a probabilistic model at the end. The transition can be done on various levels.

## 5.2 Fuzzy priors

### 5.2.1 Human control and subjective probabilities

Let us therefore consider in more detail the case where preparation and problem definition is under human control and expert knowledge about the problem domain is available, for example in a verbal form instead in terms of a probabilistic model. Those situations rely more on a subjective instead of the empirical interpretation of probability. 'Subjective probability'. may simply mean the answer of an expert, which has been asked to give a probability. But those guesses do not need to be very accurate. In cases empirical probabilities are available one can compare subjective probabilities obtained by different methods with empirical probabilities. Typical tendencies of deviations of subjective estimates from empirical probabilities have been studied (Tversky, 1972; Tversky & Kahneman, 1981; Kahneman & Tversky, 1979, 1982abc; Kahneman, Slovic, Tversky, 1983. See Lemm, 1984, for an example related to information costs in decisions). They include, besides many others, overestimation of small probabilities or of probabilities for easily retrievable, salient events, neglecting the base rate of events, the sample size or correlations, tendency towards a chosen reference point like underestimation for the sum of probabilities and overestimation for joint probabilities. Subjective estimates of probabilities do usually not obey the rules of probability theory, for example, independent obtained estimates for probabilities for events $A$ and NOT $A$ need not necessarily sum up to one. In addition to the difficulties caused by the deviations between subjective estimates and empirical probabilities it is often not obvious how to describe an event $A$ to which a subjective estimate is related. A subjective estimates for probability $p(A)$ refers to an internal representation of $A$. To relate the estimate to a probability for some external $A$ the internal representation of $A$ (e.g. concept of a mouth) must be related to the external event (e.g. images of mouths). That means $A$ has to be identified with a set of $f^0$, described by an explicit parameterization.

The process in obtaining subjective probabilities may be outlined as follows: We aim in producing a guess for probability $p(f^0)$. Let us call any deterministic question $C(f^0)$ (i.e. function or functional, respectively) of the parameters of $f^0$ a *property* of $f^0$ and $\tilde{C}$, its subjective representation, a *concept*. We want to construct a final property $C(f^0)$ which is related to $p(f^0)$ by some function $g$, i.e. $p(f^0) = g(C(f^0))$. For example $p(f^0) \propto C(f^0)$, or $p(f^0) \propto e^{C(f^0)}$. The function $g$ could also be used to compensate for known estimation biases. The subjective representation $\tilde{C}$ used to produce subjective probabilities might be difficult to relate directly to a specific function $C(f^0)$ of the parameters of $f^0$. It can be easier for other, simpler concepts $\tilde{C}_i$. For example, the concept $\tilde{C}$ of having something similar to a nose, a mouth and two eyes with certain possible spatial relations, might be difficult to relate to pixel values directly. But for a simple enough concept $\tilde{C}_i$ a property $C_i$ might be more accurately related.

A property $C_i$ measuring distance can be built out of a property $C_i'$ and a related *template* $T_i$, using a monotonic function of a 'meaningful' distance of the expectation, e.g.

$$C_i = ||C_i'(f^0) - T_i||^2.$$

$C_i'$ could be an arbitrary question. For $C_i(f^0) = \bar{y}_x(f^0)$ this measures the square distance of $\bar{y}_x(f^0)$ at point $x$ from some reference template $\bar{y}_x^T$, i.e. $||\bar{y}_x(f^0) - T_{\bar{y}_x}||^2$ and represents a usual mean square error term. Templates could also be defined relative to another question

$$||C_i(f^0) - C_{T(i)}(f^0)||,$$

like in the case of symmetries, but also in cases where not invariance but any arbitrary dependence between $C_i$ and $C_{T(i)}$ is measured.

We will write $C_i = G(\tilde{C}_i)$ for the relation between concepts and properties. Properties $C_i$ related to (linguistic) concepts $\tilde{C}_i$ have been called linguistic variables and have been used in the theory of *fuzzy sets* (see for example the collections Zadeh, 1987, or the more recent one, Zadeh, 1996). Subconcepts are modified and combined by concept functions $\tilde{F}$ to form the final concept $C$. For example, we may require: 'great similarity to $T_1$ and $T_2$ if not already very similar to $T_3$'. In analogy to concepts, also concept functions $\tilde{F}_i$ must be mapped to functions $F_i = H(\tilde{F}_i)$ acting on $C_i$. We will call functions $F$ related to concept functions $\tilde{F}$ linguistic functions. The mappings $G$ for concepts and $H$ for concept functions represent the subsymbolic level. The communication of the (symbolic) structure of $\tilde{C}$ in terms of concept variables and concept functions, i.e. $\tilde{C} = \tilde{F}(\{\tilde{C}_i\})$ is used for the approximation

$$C = F(\{C_i\}) = H(\tilde{F})(\{G(\tilde{C}_i)\}),$$

where $H(\tilde{F})$ stands for the function $F$ with all included subfunctions replaced according to $F_j = H(\tilde{F}_j)$. (See Fig.3). We try to achieve

$$G(\tilde{C}) = G(\tilde{F}(\{\tilde{C}_i\})) \approx F(\{C_i\}) = H(\tilde{F})(\{G(\tilde{C}_i)\}).$$

However, this is difficult to check in general if $G(\{\tilde{C}\})$ cannot be obtained consistently in a direct way. Indeed,
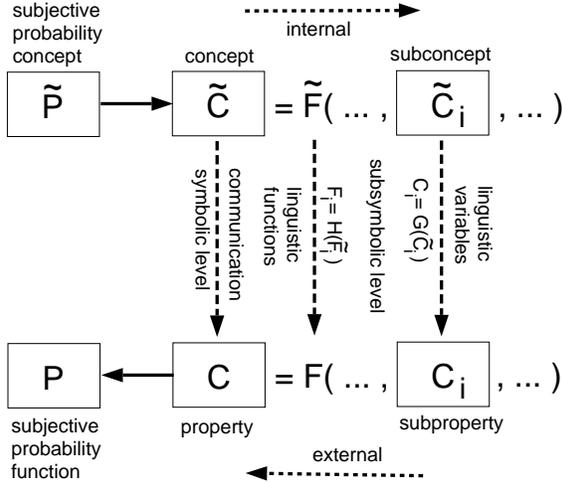
Figure 3: The relation between (internal) concepts and (external) properties. For complex concepts $\tilde{C}$ the relation to property $C$ is often easier to obtain and more invariant under transfer to new situations, if the mapping from concept to properties (subsymbolic mapping for linguistic variables) is done for simpler subconcepts $\tilde{C}_i$. This, however, requires also a (subsymbolic) mapping of concept functions $\tilde{F}$, acting on concepts $\tilde{C}_i$, to property (external, linguistic) functions, acting on properties $C$. Then, the property $C$ for a communicated symbolic structure $\tilde{C} = \tilde{F}(\{\tilde{C}_i\})$ can be approximated by $C = F(\{C_i\})$.

this difficulty is the reason for the decomposition into subconcepts. One can instead define $C$ by the right hand side and check variability of dependencies from those properties under transfer to new situations. Thus, this can be seen as a heuristic method to achieve dependencies with smaller variance under transfer.

We will call priors obtained by this method *fuzzy priors*. In the following we do not intend to give an introduction into fuzzy logic. We mainly want to stress that related techniques can be adapted to construct priors and the more technical point that this can lead to non-linearities in the equations determining the maximum probability $p(f^0)$ even if single $C_i$ correspond to Gaussian probabilities.

The method can be applied in two variants. The first way of defining $p(f^0|f^P)$ is directly to describe the functional in its dependence from the parameters without explicit reference to data. A second possibility is to define in a first step $p(y^P|q^P, f^0)$ in its complete dependence from all three variables and then choose in a second step $y^P$ to fix the prior. We will mainly study the first possibility and we will study linguistic functions $C_C = F(C_A, C_B)$ depending on two one–dimensional properties. Their combination gives functions depending on more than two variables.

### 5.2.2 Real valued extensions of logic

Linguistic functions can be constructed by fixing there values at certain, typical points and using some (e.g. smooth, symmetric) interpolation scheme. For example, practically important linguistic functions are such related to logical functions like AND, OR or NOT used in *fuzzy logic* (see e.g. Kandel, 1982, Klir & Yuan, 1995). From such functions we assume that they coincide with the corresponding binary logical functions at the four corners, i.e. where both arguments have a minimal or maximal value. We will mainly concentrate on this functions but also model combinations which do not correspond to binary logical functions at the four corner points.

For example, we may assume high probability for a function $f^0$ if it is smooth at $x_1$ AND $x_2$ AND $\cdots$. In the second variant of the method we begin with the construction of answer distributions $p(y|q, f^0)$ and a prior results by choosing data $(y_i^P, q_i^P)$. There, for example in face recognition, $q^P$ can be defined as question looking for components like eyes in the images $y$ produced by $p(y| x = \text{face}, f^0)$ and the $y^P$ can be examples of how eyes can look like. Requiring that $f^0$ produces (all of) several variants $V_i$ of eyes, i.e. its output is $V_1$ OR $V_2$ OR $\cdots$, means ANDing for constructing a prior, so $f^0$ can produce $V_1$ AND $V_2$ AND $\cdots$.

Consider properties with $0 \leq C_i \leq m$, where we may allow the limit $m \to \infty$. We use here a 'distance interpretation' of the values 0 and $m$ where we interpret the value $m$ as complete absence ('far' or 'False') and the value zero as complete presence ('near' or 'True') of that property (deviation property). Monotonic function from (bounded) distances to some templates are examples. The interpretation of 0 and $m$ can be reversed (similarity property), like in the usual convention in logic or for properties like probabilities where 0 means 'False' or 'impossible' and $m = 1$ 'True' or 'sure'. Then, just the definitions of AND and OR in the following have to be exchanged.

The following are examples of a function equal to a binary logical function when the arguments have values 0 or $m$:

$$
\begin{aligned}
C_{(A \text{ AND } B)} &= C_A + C_B - C_A C_B / m, \\
C_{(A \text{ OR } B)} &= C_A C_B / m, \qquad\qquad (5) \\
C_{(\text{NOT } A)} &= m - C_A.
\end{aligned}
$$

These operations represent a Boolean Algebra (see for example Whitesitt, 1995) and therefore for variables taking only the values 0 and m for example DeMorgan's law $A$ OR $B = \text{NOT}( (\text{NOT } A) \text{ AND } (\text{NOT } B))$ is valid.

The functions are extended to real variables by allowing values $0 \leq C \leq m$ in the above formula. Of course, this extension to real values cannot be unique. So AND can also be implemented as

$$
C_{(A \text{ AND } B)} = c_m(C_A + C_B).
$$

where $c_m$ is a cutoff function with $c_m(x) = id$ for $0 \leq x \leq m$ and $c_m(x) = m$ for $x \geq m$ and therefore always $c_m(C_A) = C_A$. With this, DeMorgan's law gives for OR

$$
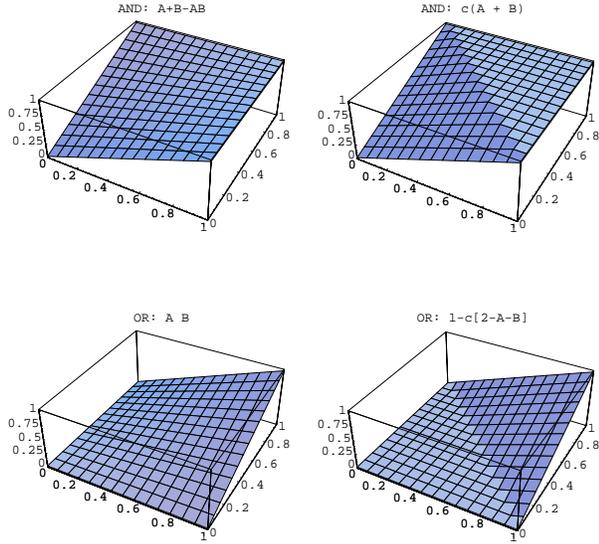C_{(A \text{ OR } B)} = m - c_m(2m - C_A - C_B).
$$

24

Figure 4: Different realizations of AND and OR, in a 'distance interpretation' where 0 stands for 'True' ('near') and 1 for 'False' ('far'). $C_A$, $C_B$ are abbreviated $A$, $B$ and the cutoff function $c_1$ as $c$.



Figure 5: Examples for real valued extensions of tautologies. The last one has the formula $1 - (1 - A^2(1-A)^3 + 1 - B^2(1-B) - (1 - A^2(1-A)^3)(1 - B^2(1-B)))^{100}$.

Not differentiable near our main point of interest at $(0,0)$ where the 'good' $f^0$ should be located but $m$–independent are the following realizations

$$C_{(A \text{ AND } B)} = \max(C_A, C_B),$$
$$C_{(A \text{ OR } B)} = \min(C_A, C_B).$$

This representation of AND and OR by the maximum resp. minimum operation are the standard fuzzy operations used in fuzzy logic.

A variable $C_A$ can always be written $C_A = C_A^2/m$, $C_A = 2C_A - C_A^2/m$, $C_A = c_m(C_A)$. Any (usually monotonic) interpolating function $\sigma(x)$ with $\sigma(0) = 0$ and $\sigma(m) = m$, allows to replace $A$ by $\sigma(A)$ without changing the limit of binary logic. For example, OR could be defined as $\sigma_{OR}(\sigma_A(C_A)\sigma_B(C_B)/m)$ choosing some (monotonic) functions $\sigma_{OR}$, $\sigma_A$, $\sigma_B$. Rules like the law of DeMorgan are valid for variables with values 0 and $m$ but not in general for real values. Fig.4 shows two real valued extensions of AND and OR and Fig.5 some tautologies which can be added to every function without changing the binary limit. Of course, any real valued extension of logical functions is arbitrary except at the four corners. Thus, one might add additional conditions for such functions, like monotonicity or smoothness conditions, or the requirement that certain laws of the Boolean algebra, valid for the binary limit also hold for real values, requiring for example an associative, commutative, and distributive AND and OR. Besides having the correct boundary conditions, usually monotonicity, commutativity, and associativity are required for fuzzy operations. (Such operations are called $t$-norm or $t$-conorm for boundary conditions corresponding to AND or OR.)

### 5.2.3 General combinations

Every logical formula can be expressed by NOT, OR, AND. In particular, one could use the disjunctive or con-
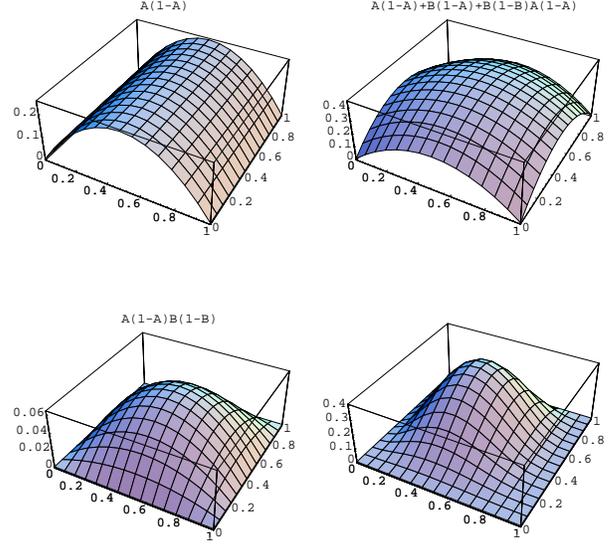
junctive normal form. Also, AND or OR could be eliminated using DeMorgan's law. For example, the exclusive OR is defined as XOR = ($A$ AND NOT $B$) OR (NOT $A$ AND $B$). For disjunct events XOR is equivalent to OR. Accordingly, real extensions of the exclusive OR are

$$C_{XOR} = \frac{1}{m}c_m(C_A + m - C_B)c_m(m - C_A + C_B).$$

or

$$C_{XOR} = \frac{2C_A C_B}{m} - C_A - C_B + m.$$

The latter result can be connected to Eqs.(5) if using $C_i^2/m = C_i$. Fig.6 shows some possible variations of XOR.

Another important relation is IF $C_1$ THEN $C_2$ = NOT $C_1$ AND $C_2$. One way to extend this to real values is

$$C_{IF} = c_m((m - C_1) + C_2).$$

The limit of binary logic requires function values 0 or $m$ at the four corners where the input variables are 0 or $m$. Any function having more than two different values at those four points cannot correspond to a logical function. The only possible linear functions are for example 1, $C_A$, $C_B$ and their negations. Therefore combinations like $aC_A + (1-a)C_B$ (LIN = linear) are more general and do not correspond to logical combinations at the four corners. However, setting $a = C_C$ LIN appears to be a combination of three logical properties with the value of $C_C$ fixed (for all $f^0$).

Thus, such functions can be seen as parts (of combinations) of logical functions with certain values of $C_i$ excluded (e.g. with the value fixed). Specifically, in the limit $m \to \infty$ and $C_i$ finite we obtain for the above rules for $C_A \neq \infty \neq C_B$

$$\begin{aligned} C_{(A \text{ AND } B)} &= C_A + C_B, \\ C_{(A \text{ OR } B)} &= 0, \quad\quad (6) \\ C_{(\text{NOT } A)} &= \infty. \end{aligned}$$
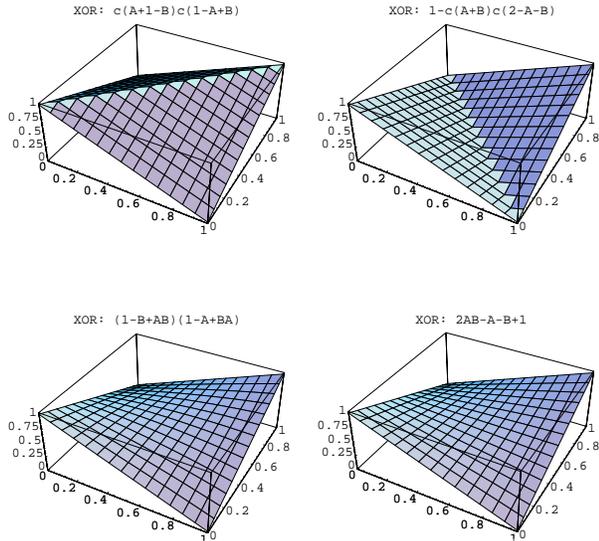
Figure 6: Different realizations of XOR

These are only linear functions. To allow the value $\infty$ for the $C$ one has to extend the definition of OR and NOT according to

$$C_{(A \text{ AND } B)} = \infty \quad \text{for } C_A = \infty \text{ or } C_B = \infty,$$

$$C_{(A \text{ OR } B)} = \begin{cases} C_A \text{ for } C_B = \infty \neq C_A, \\ C_B \text{ for } C_A = \infty \neq C_B, \\ \infty \text{ for } C_A = C_B = \infty. \end{cases}$$

$$C_{(\text{NOT } A)} = 0 \quad \text{for } C_A = \infty.$$

and one finds that the functions on the whole interval $[0, \infty]$ are nonlinear. This representation is interesting in as far as the AND is linear and the nonlinearity of the OR can be implemented by just skipping functions with final $C(f^0) = \infty$ from $F^0$. We will call this a hard implementation of the OR and other implementations soft. Then for all functions $f^0$ under consideration, i.e. in $F^0$, $C(f^0)$ is a linear function of its constituents $C_i(f^0)$. But note that unlike for finite $m$ here $C_{(\text{NOT NOT } A)}$ is not equal $C_A$ for $0 \neq C_A \neq \infty$.

LIN in contrast to AND is strictly monotonically in all $C_i$ even if one $C_{i*} = m$. This allows also if $C_{i*} = m$ cannot be changed (momentarily), the use of gradient information to improve the other $C_i$. To include information related to OR one can combine LIN with a multiplicative implementation of OR:

$$C_{(A \text{ LIN } B)} = aC_A + (1-a)C_B, 0 < a < 1,$$
$$C_{(A \text{ OR } B)} = C_A C_B / m. \tag{7}$$

## 5.3 Special Properties: probabilities, logprobabilities, distances, and averages

We shortly discuss some special properties: probabilities (or partition sums, if not normalized), logprobabilities (related to free energies), distances (related to scalar products) and averages (expectations, related to energies).

### 5.3.1 Probabilities

For probabilities (not densities) the logical combination of events are defined

$$p(A \text{ AND } B) = p(A)p(B|A)$$
$$p(A \text{ OR } B) = p(A) + p(B) - p(A)p(B|A)$$
$$p(\text{NOT } A) = 1 - p(A).$$

(We can understand corresponding expressions for densities formally to be defined by $p(F(A, B)) = \int_{F(A,B)} p(x)dx = \int_A \int_B p_{F(A,B)}(x_A, x_B)dx_A dx_B$, and interpret the rules as $p(x_A AND x_B)dx_A dx_B = p_A(x_A)p_B(x_B|x_A)dx_A dx_B$ so the integral factorizes for independent components $x_A$, $x_B$. With a function $\delta$ fulfilling $\int_A \delta(A)dx_A = 1 \int_B \delta(B)dx_B = 1$, we have $p(x_A OR x_B)dx_A dx_B = (p_A(x_A)\delta(B) + p_B(x_B)\delta(A) + p_{AB}(x_A, x_B))dx_A dx_B$, i.e. $(p_A(x) + p_B(x))dx$ for disjunct events $x_A$, $x_B$, and $p(NOT x_A)dx = (\delta(A) - p_A(x_A))dx_A$. So integration of densities to obtain probabilities $p(A) = \int_A p(x)dx$ has the form of an OR for disjunct events.) With respect to the two (similarity) properties $P(A)$ and $P(B)$ the probabilistic AND, i.e. $p(A)p(B)(p(B|A)/p(B))$, is not one function but a whole family, parameterized by $p(B|A)$, but all coinciding at the four binary corners and the quantitative form of a probabilistic AND changes with the dependency of A on B.

To take advantage of independence in the case of OR one can use DeMorgan's law and write for a set of independent $A_i$, $p(OR_i A_i) = 1 - \prod_{i=0}(1 - p(A_i))$, where $p(NOT(AND_i A_i))$ factorizes. This form is also known as 'noisy OR' (Pearl, 1988, Jensen, 1996), and one may set one $A_i$ constantly equal to one (or zero) to have a nonzero baseline if all other $A_i$ are zero (or one).

### 5.3.2 Log-probabilities

Assume $0 \leq p \leq 1$ is a probability (or for AND and OR a bounded density $0 \leq p \leq c$) then the log-probability $L = \ln p$ is in the intervals $[0, -\infty]$ (or $[\ln c, -\infty]$, respectively) and the rules for probabilities become

$$L(A \text{ AND } B) = L(A) + L(B|A)$$
$$L(A \text{ OR } B) = \ln\left(e^{L(A)} + e^{L(B)} - e^{L(A \text{ AND } B)}\right)$$
$$L(\text{NOT } A) = \ln(1 - e^{L(A)}). \tag{8}$$

Especially interesting is the OR for disjunct events, where it is equivalent to an XOR. Fig.6 (where for probabilities the role of one and zero have to be exchanged) shows that XOR is prototypical for a situation with two clearly separated degenerated maxima. The two degenerate global maxima remain under slight perturbations both still local maxima. It needs a 'considerable' deformation of the probability surface, or, equivalently, an in comparison to other influences relatively flat XOR probability surface, to let one of the two local maxima disappear. At this point the solution to the problem of finding the maximal probable state shows a bifurcation. We will see below that this is related to phase transitions and may remark at this point that the term 'temperature' is related to the relative flatness of the XOR.

26

In contrast, OR for independent events, shown for example in Fig.4 has a continuous line of degenerate maxima, and already a small perturbation can lead to a unique maximum. Indeed, any OR can be expressed as an OR for independent events. Choosing for example $\omega_1 = A$ AND $B$, $\omega_2 = A$ AND NOT $B$, $\omega_3 = $ NOT $A$ AND $B$, $\omega_4 = $ NOT $A$ AND NOT $B$, we have $A$ OR $B = \omega_1$ XOR $\omega_2$ XOR $\omega_3$, leading to the degenerated maxima.

Therefore, we look in a bit more detail to the OR for disjunct events. Expanding the log-probability for disjunct $A$, $B$ around some $L_0$ we obtain in second order in $\Delta_A = L(A) - L_0$, $\Delta_A = L(B) - L_0$

$$L(A \text{ OR } B) \approx L_0 + \ln 2 + \frac{1}{2}\left(\Delta_A + \Delta_B\right) + \frac{1}{8}\left(\Delta_A - \Delta_B\right)^2$$

$$= \ln 2 + \frac{L(A) + L(B)}{2} + \frac{1}{8}\left(L(A) - L(B)\right)^2 \qquad (9)$$

$$= \frac{1}{2}\left(\ln\frac{e^{L(A)}}{2} + \ln\frac{e^{L(B)}}{2}\right)$$

$$+ \frac{1}{8}\left(L(A)^2 - L(A)L(B) + L(A)^2\right),$$

independent of $L_0$, and exact if $L(A) = L(B)$. The last line shows that the first order terms gives the *annealed approximation* where a nonlinear function of the average is replaced by the average of the nonlinear functions $g(<x>) = g(\sum_i p_i x_i) \approx \sum_i p_i g(x_i) = <g(x)>$, with $\sum_i p_i = 1$ (See for example, Seung, 1995). For convex functions the inequality $g(\sum_i p_i x_i) \leq \sum_i p_i g(x_i)$ holds for convex combinations, i.e. $\sum_i p_i = 1$, with equality if all $x_i$ are equal (Jensen's inequality). For the concave logarithm, concave $f$ meaning $-f$ is convex this reads, for example, $\ln(\sum_i p_i x_i) \geq \sum_i p_i \ln(x_i)$.

We may also consider a weighted OR, for disjunct events $A_i$

$$L(OR_i^n A_i) = \ln P(OR_i^n A_i) = \ln\left(\sum_i^n a_i e^{L_i}\right),$$

with $L_i = L(A_i)$. Disjunct $A_i$ may or may not be elementary events $\omega_i \in \Omega$ of the model under study, but they can always be seen as (effective) elementary events for $\bigcup_i A_i$ with respect to a specific OR. A weighted OR, corresponds to an unweighted OR in an enlarged event space, as can be seen by writing

$$a_i e^{L_i} = e^{L_i + L_i^a}, \quad L_i^a = \ln a_i.$$

This may be interpreted as a situation where the $A_i$ have an additional independent dimension $a$ with values $i$, i.e of the labels of the mixture components, so the complete event $A_i$ has log-probability $L(A_i) + L^a(i)$ corresponding to an AND for independent events with $p_i = a_i/Z_a$ and $Z_a = \sum_i a_i$. Thus, $L_i = L(A_i)$ can be seen as a special random variable on the events $i$.

Especially, a (weighted ) OR has the structure of a *cumulant generating function* with respect to the normalized weighting factors $a_i/Z_a$, representing probabilities. (See also Section 5.3.4.) To see this, we look at the Taylor expansion of $P$ around $L_0$

$$P(OR_i^n A_i) = \sum_i^n a_i e^{L_i} = \sum_i^n a_i e^{\Delta_i} e^{L_0}$$

$$= \sum_k^\infty \frac{1}{k!}\sum_i^n a_i \Delta_i^k e^{L_0} = Z_a e^{L_0}\sum_k^\infty \frac{1}{k!} <\Delta_i^k>_a,$$

which contains in the expansion coefficients all $k$th moments of the differences $\Delta_i$ (or of the $L_i$ itself for $L_0 = 0$) with respect to the $A_i$ and $p_i$,

$$M_k^a(\Delta) = <\Delta^k>_a = \frac{1}{Z_a}\sum_i^n a_i \Delta_i^k$$

$$= \frac{d^k}{d\beta^k}\frac{1}{Z_a}\sum_i^n a_i e^{\beta\Delta_i}\Big|_{\beta=0} = \frac{d^k}{d\beta^k}<e^{\beta\Delta_i}>_a\Big|_{\beta=0},$$

with $M_0^a = 1$ due to $\Delta_i^0 = 1$.

We can also introduce a common scaling factor $\beta$, often called *inverse temperature* for all $L_i$. Then because

$$\frac{df(\beta L_i)}{d\beta}\Big|_{\beta=0} = L_i\frac{df(\beta L_i)}{d(\beta L_i)}\Big|_{\beta L_i=0}$$

$$= L_i\beta^{-1}\frac{df(\beta L_i)}{dL_i}\Big|_{L_i=0,\beta=1} = L_i\frac{df(L_i)}{dL_i}\Big|_{L_i=0},$$

both derivatives generates the coefficients of the Taylor expansion. Indeed, one sees directly

$$M_k^a(\Delta) == \frac{d^k}{d\beta^k}\sum_i^n \frac{a_i}{Z_a}e^{\beta\Delta_i}\Big|_{\beta=0} = \frac{d^k}{d\beta^k}<e^{\beta\Delta_i}>_a\Big|_{\beta=0},$$

Thus, while the moments $<\Delta^k>_a$ are the Taylor coefficients of the *moment generating function* $<e^{\beta\Delta}>_a$ for $\beta = 1$ they can also be calculated as their derivatives at $\beta = 0$. Thus, they are up to the factor $Z_a e^{L_0}/k!$ the Taylor coefficients of a *high temperature expansion* around $\beta = 0$. Notice that we expanded every $L_i$ around the same $L_0$ so $E^{L_0}$ could be factored out. In general we could use different origins $L_0^i$ for different $L_i$. Indeed if the $e^{\beta L_i}$ are becoming relatively more separated (large $\beta$ or low temperature case) a common origin for the Taylor expansion becomes a less good choice. (For finite systems, i.e. a finite $\sum_i$, the sum of exponentials is an analytic function so the convergence radius of the corresponding power series is infinite. In contrast for example for the logarithm function (see below), the convergence radius of a Taylor expansion is limited.) These problems are typical for phase transitions (See Section 10) where the log-probability of $L(OR)$ changes its number of maxima, and the behavior of the $L(OR)$ at $L_i = L_0, \forall i$ is drastically different (being e.g. a minima) than that of the components at $L_i = L_0$ (where it may be a maxima).

Analogously, we can expand $L = \ln P$ in the parameter $\beta$ around $\beta = 0$, and obtain the high temperature expansion of $L(OR)$

$$L(OR_i^n \beta A_i) = \ln P(OR_i^n \beta A_i)$$

$$= \ln Z_a + \beta L_0 + \sum_{k=0}^\infty \frac{\beta^k}{k!}C_k^a(\Delta). \qquad (10)$$

where the coefficients

$$C_k^a(\Delta) = \frac{d^k}{d\beta^k}\ln\left(\frac{1}{Z_a}\sum_i^n a_i e^{\beta\Delta_i}\right)\Big|_{\beta=0},$$

27

are called *cumulants*. For example one finds, $C_0^a = 0$, $C_1^a = M_1^a$, $C_2^a = M_2^a - (M_1^a)^2$, $C_3^a = M_3^a - 3M_2^a M_1^a + 2(M_1^a)^3$.

As the origin of expansion $L_0$ enters through the linear term only all cumulants of order $k \geq 2$ are $L_0$ independent (i.e. independent of the mean), and the first cumulant, i.e. the mean, compensates for the constant because $C_1^a(\Delta) = \frac{1}{Z_a} \sum_i a_i \Delta_i = C_1^a(L) - L_0$. Thus, for all $L_0$ the expansion looks like for $L_0 = 0$

$$L(OR_i^n \beta A_i) = \ln Z_a + \sum_{k=0}^{\infty} \frac{\beta^k}{k!} C_k^a(L).$$

For example up to second order this gives ($a_i = 1 \Rightarrow Z_a = n$)

$$L(OR_i^n A_i) = \ln P(OR_i^n A_i) = \ln \left( \sum_i^n a_i e^{L_i} \right)$$

$$= \ln Z_a + <L>_a + \frac{1}{2} \left( <L^2>_a - <L>_a^2 \right). \quad (11)$$

It may be interesting to note, that for example the mean (first moment and first cumulant) for our special random variable $L = \ln p$ is just the negative *average information* (entropy) $<L>_a = <\ln p>_a = -I_a(p)$ with respect to the $p_i = a_i/Z_a$. The negative log-likelihood $-L$ is also sometimes called bit–number with the corresponding bit–cumulants with respect to the $a_i/Z_a$ generated by $\ln <e^{-\beta L}>_a$, related to the Rényi information $I_\beta$ by $\ln <e^{-(1-\beta)L}>_a = (\beta - 1)I_\beta$ (Beck & Schlögl, 1993).

We can also split $L_i$ for all mixture components $i$ into several parts, $L_i = \sum_j L_{i,j}$, (which means ANDing), and define a corresponding set of $\beta_j$ so we that have

$$L(OR_i^n(AND_j A_{i,j})) = \ln \sum_i a_i e^{\sum_j \beta_j L_{i,j}}.$$

Then the mixed derivatives of $L$ gives the multidimensional cumulants of $L_{i,j}$ or $\Delta_{i,j} = L_{i,j} - L_{0,j}$ according to

$$C_{k_1, k_2, \cdots, k_m}^a (\Delta_{j_1}^{k_1}, \Delta_{j_2}^{k_2}, \cdots, \Delta_{j_m}^{k_m})$$

$$= \frac{d^{k_1}}{d\beta_{j_1}^{k_1}} \frac{d^{k_2}}{d\beta_{j_2}^{k_2}} \cdots \frac{d^{k_m}}{d\beta_{j_m}^{k_m}} \ln \left( \frac{1}{Z_a} \sum_i^m a_i e^{\sum_j \beta_j \Delta_{i,j}} \right) \Bigg|_{\beta=0}.$$

The terms $e^{\beta L_i}$ become for large $\beta$ (low temperature) the smaller the smaller $L_i$. In the limit $\beta \to 0$ only events with maximal probability $p(A_{i*}) = \max_i p(A_i)$ ('ground states') survive. Therefore, skipping from the sum $\sum_i P(A_i) = \sum_i e^{L(A_i)}$ for disjunct $A_i$ the smaller terms $P(A_i) \leq \Theta$ (low probability events) and keeping only the larger ones ($p(A_i) > \Theta$) (high probability events) is also called *low temperature expansion*. For continuous variables $i$ where the sum is replaced by an integral this approximation appears as saddle point approximation (see Section 7.2) and their higher order variants.

We also take a short look to an OR for independent events. To make use of the factorization of probabilities for independent events we apply DeMorgans law to write

for discrete events, corresponding to the noisy OR for probabilities

$$L(OR_i A_i) = \ln(1 - \prod_i^n (1 - e^{L_i})), \quad (12)$$

and

$$L(NOT(OR_i A_i)) = \ln(\prod_i^n (1 - e^{L_i}))$$

$$= \sum_i^n \ln(1 - e^{L_i}).$$

For dependent $A_i$ the conditioned factors have to be used. Jaakola & Jordan (1996) give an expansion of the function

$$\ln(1 - e^L) = \sum_{k=0}^{\infty} \ln g(-2^k L),$$

in terms of the logistic function $g(z) = 1/(1 + e^{-z})$ and use it for efficient, approximate calculations in graphical models.

Partitioning the events $A_i$ into non–overlapping subsets, i.e. into disjunct (effective) elementary events $\omega_j$ we have $p(A_i) = p(OR_j^A \omega_j) = \sum_{j \in A} p(\omega_j)$. Then we see that expressing the OR for non–disjunct $A_i$ by disjunct $\omega_j$ gives

$$L(OR_i A_i) = L(OR_j^N \omega_j) = \ln(\sum_j N_j e^{L(\omega_j)}),$$

which reproduces for the smallest of those partitions the product terms $p(\omega_j) = e^{L_{i_1}} e^{L_{i_2}} e^{L_{i_3}} \cdots$ in Eq.12, reweighted with the number $N_j$ of $A_i$ which contain $\omega_j$, and without the intermediate terms with oscillating sign. One can now, for example, apply the high temperature expansion with $Z_a = Z_N = \sum_j N_j$.

### 5.3.3 Distances

Monotonic functions of negative *distances* $\|y^q(f^0) - T_{y^q}\|^2$ are often used for log-probabilities. The definition of $q$ can always be changed in such a way that $T_{y_q} = 0$. The most common example are data terms $d^2 = \sum_i \frac{1}{2\sigma_{q_i}^2} (y_i^{q_i} - \bar{y}_{q_i}(f^0))^2$ where the $y_i^{q_i}$ represent the template vector and $\bar{y}_{q_i}$ the deterministic answer $y^{q_i}(f^0)$.

In spaces where the distances fulfill

$$\|g - h\|^2 = 2\|g\|^2 + 2\|g\|^2 - \|g + h\|^2$$

there exists a scalar product, written as $< \cdot \,|\, \cdot >$, related to the distance by

$$\|h\|^2 = <h \,|\, h>.$$

Different positive definite kernels $\mathcal{O}$ define different scalar products

$$<g \,|\, h>_{\mathcal{O}} = <g \,|\mathcal{O}|\, h> = \int dx \int dx' \, g(x)\mathcal{O}(x, x')h(x').$$

Minimizing squared distances

$$\|y^q(f^0) - T_q\|^2 = <y^q(f^0)|y^q(f^0)> -2 <y^q(f^0)|T_q> +c,$$

(written for real $< y^q(f^0) \,|\, T_q > \; = \; < y^q(f^0) \,|\, T_q >^*$ $=<T_q \,|\, y^q(f^0)>$ with $c$ an $f^0$–independent constant and $*$ indicating complex conjugation) is equivalent to maximizing scalar products (*overlaps*), like $\sum_i \frac{1}{2\sigma^2_{q_i}} T_{q,i} \bar{y}_{q_i}$, for normalized $y$ and $\bar{y}$. (But normalization is a maximal nonlocal condition.) For overlaps one has to use rules with zero representing false and one representing true so the definitions of AND and OR are exchanged.

For example, properties which have positive values bounded by $m$ can be obtained from log-probabilities or distances by $C_A = g^{-1}(d^2(A)) = m(1 - e^{-d^2(A)}) = m(1 - e^{L(A)}) = m(1 - p(A))$ with inverse $L = g(C) = -\ln(1 - C/m) = \ln(m/(m-C))$. This gives for independent $A$ and $B$ the equations (5) for properties. Thus, using $g(C) = -\ln(1 - C/m)$ directly relates properties to probabilities according to $C_i = m(1 - p_i)$.

Writing OR in terms of Euclidean square distances by choosing $L_i = -d_i^2/(2\sigma_i^2) + \ln c_i$, with some constants $\sigma_i$ and $c_i$, one obtains for disjunct events a *Gaussian mixture model*

$$p(f^0) \propto \sum_i c_i e^{-d_i^2(f^0)/(2\sigma_i^2)}.$$

For non–disjunct events product terms like

$$\prod_i e^{-d_i^2/\sigma_i}$$

must be subtracted (or approximately a cut-off function can be included).

### 5.3.4 Averages
**Unnormalized probabilities**

Averages or expectations can be related to unnormalized probabilities $Z(A, \beta)$, (with $\beta$ a fixed parameter we will discuss below) also called *partition sums*

$$p(A, \beta) = \frac{Z(A, \beta)}{Z(\beta)}.$$

Thus,

$$Z(A, \beta) = \sum_{\omega \in A} Z(\omega, \beta) \qquad (13)$$

$$Z(\beta) = Z(\Omega, \beta) = \sum_{\omega \in \Omega} Z(\omega, \beta)$$

for a complete set of disjunct (elementary) events $\omega$, and $p(\Omega) = 1$. If we define $Z(A|B)$ by $p(A|B) = Z(A|B)/Z$ partition sums transform under logical operations like probabilities with an additional normalization factor $Z$ analogously to $m$ in the rules (5) (with AND and OR exchanged). If we choose $Z(A|B) = Z(A, B)/Z(B) = p(A|B)$ (what one can do for averages, see below) then $Z$ appears only in the NOT. Defining shifted log-probabilities $-\beta F$ by

$$Z(\beta) = e^{-\beta F(\beta)}, \quad Z(A, \beta) = e^{-\beta F(A, \beta)},$$

the $F$, also called *free energies*,

$$F(\beta) = -\frac{1}{\beta} \ln Z(\beta), \quad F(A, \beta) = -\frac{1}{\beta} \ln Z(A, \beta),$$

transform[30].

$$F(A \text{ AND } B) = F(A) + F(B|A) - F$$

$$F(A \text{ OR } B) = -\frac{1}{\beta} \ln \left( e^{-\beta F(A)} + e^{-\beta F(B)} \right.$$
$$\left. - e^{-\beta F(A \text{ AND } B)} \right) \qquad (14)$$

$$F(\text{NOT } A) = -\frac{1}{\beta} \ln(e^{-\beta F} - e^{-\beta F(A)}).$$

For the sake of simplicity the $\beta$–dependence of $F$ is here not written explicitly.

We choose as family $A$ a set of disjunct events $\omega \in \Omega$. For the corresponding $\omega$ we define the free energies

$$F(\omega, \beta) = E(\omega)$$

to be $\beta$–independent. We will call them *energy* of $\omega$ and can then write

$$p(\omega, \beta) = \frac{Z(\omega, \beta)}{Z(\beta)} = \frac{e^{-\beta E(\omega)}}{Z(\beta)} = e^{-\beta(E(\omega) - F(\beta))}. \quad (15)$$

[30]Clearly, the OR looks quite complicated, even for disjunct events, and indeed, it is the source of many difficulties. The summation corresponds for example to the calculation of partition sums in statistical physics. Another summation outside the logarithm is added for disordered systems, like for spin glasses, where the (shifted) log-likelihood (or energy) function governing the system is also in reality not exactly known. In principle this corresponds to adding additional components (AND) to the elementary events $\omega$ with possible realizations $i$. Combining different possible realizations $L^i$ by $OR$ corresponding to $\ln \sum_i p_i e^{L_i}$ with $p_i$ a probability distribution over realizations of the energy function. Alternatively, one can restrict to averages (weighted AND) of observables (and therefore the partition sum) over realizations $i$. That means one integrates over part of the components of $\omega$ and considers $\sum_i p_i \ln e^{L_i} = \sum_i p_i L_i$. Including the thermical OR over complete sets of disjunct events (states $L_{ij_i}$), in the average (over different 'replicas' $i$ of a system with differently 'quenched' interactions), yields $\sum_i^n \ln(\sum_j^m e^{L_{ij}}) = \ln(\prod_i^n \sum_j^m e^{L_{ij}}) = \ln(\sum_{j_1}^m \cdots \sum_{j_n}^m \prod_i^n e^{L_{i,j_i}})$, where $i$ can be called replica index. The product of sums creates all kind of product terms for different systems $j_i$ and can be huge or infinite. The replica approach uses the identity $\ln Z = \lim_{n \to 0} \frac{Z^n - 1}{n}$ to substitute this large product by a product with number of factors going to zero, i.e. $n \to 0$ (not 1 !). In the corresponding asymptotic mean field approximation also correlations between different replicas of the systems enter the theory. However it requires more assumptions and calculation tricks than a standard saddle point approximation. For example analytical results for an integer $n$ have to be analytically continued to real $n$ to obtain the limit $n \to 0$. (Mezard, Parisi, & Virasoro, 1987). (Special observables are those which do not fluctuate with $i$, like for large systems for example variables with fluctuations vanishing fast enough with the system size. The values of such observables $x^{nf}$ is then the same for every individual system of the ensemble. Restricted to only non–fluctuating observables $x^{nf}$, i.e. $\omega_i(x) = \omega_i(x^{nf})$ and therefore $E(\omega) = E(x^{nf})$ the probability distribution $p_i$ becomes deterministic and the $\sum_i$ disappears. This is the reason one is especially interested in the (thermodynamic) limit of infinite system size where for some (self–averaging) observables the fluctuations around the average can disappear, e.g. for uncorrelated variables the Gaussian limit theorem applies.)

We see that the log-probability $L = -\beta(E(\omega) - F(\beta))$ is hereby written as the (negative, $\beta$-scaled) difference of a $\beta$-independent energy which describes the system, i.e. the variation of the probability between the $\omega$, and a $\beta$-dependent 'scaling shift' containing the $\beta$ (temperature) dependence.

Sometimes one may also wish to consider *effective energies* $E^{eff}(\omega, \beta)$ which are $\beta$-dependent and might be free energies with respect to a finer set of elementary events. They are however most times used in the range where they are approximately independent of $\beta$ (temperature), and we will, if not stated explicitly otherwise always assume that energies are $\beta$-independent. We can also take the random variable $E$ being dependent on other random variables $E_j(\omega)$ so that $\beta F(\omega, \beta) = g(\{E_j(\omega)\})$. Useful will be a linear combination, for which we write[31]

$$\beta F(\omega) = \beta E(\omega) = \sum_j \beta_j E_j(\omega),$$

with $E_j$ independent of $\beta_j$ and $\beta$. For the sake of simplicity we will understand $F(A, \beta)$ to mean $F(A, \beta, \{\beta_j\})$ in that case, and analogously for other variables like $Z$ and $p$.

### Averages and generating functions

Partition sums can be seen as averages under a uniform probability distribution $p(\omega)$. In general, *averages* or *expectations* of a function $h(\omega)$ over a family $A$ of disjunct events $\omega \in A$ are defined as

$$h(A, \beta) = < h(\omega) >_{A,\beta} = \frac{\sum_{\omega \in A} p(\omega, \beta) h(\omega)}{p(A, \beta)},$$

and we can include $h(A|B) = h(A, B)/h(B)$. We can extend the definition of $E(\omega)$ to $E(A, \beta)$ for all events $A \in \Omega$ by defining the energy to transform like (i.e. to be) an average.

Thus, $E$ is a random variable bounded from below, defined by the number (or vector) $E(\omega) > -\infty$ related to every elementary event $\omega \in A$. We may call $p(\omega)$ the distribution generated by $\beta E$ which can have the form $\sum_j \beta_j E_j$. On the other hand every (vector) $\beta E$, with components bounded from below, is a 'generating random variable' of some $p$, with $Z(A, \beta)$ the normalization constant of $e^{\beta E(\omega)}$ on $A$.

We already encountered unnormalized probabilities and shifted log-probabilities in the discussion of the (weighted, but also unweighted) OR in Section 5.3.2. There the $L$ could be seen energies for (effective) elementary events $A_i$, i.e. for $A_i$ they are also shifted log-probabilities with respect to the unnormalized probabilities $a_i$ corresponding to $Z(\omega, \beta = 1)$. We also introduced auxiliary variables $\beta_j$ and used a splitting of $L_i$ into components $\sum_j \beta_j L_{i,j} + L_i^a$ equivalent to $\sum_j \beta_j E(\omega)$. And therefore we can here relate energies (averages) and free energies (OR over $\Omega$) the same way by derivatives of generating functions as we did in 5.3.2. We briefly formulate this principle again in terms of $E$ and $F$.

---

[31] Then in physics usually only one subgroup $E_j$ is called energy, another might be called particle number.

Now we want to calculate the expectations of the generating random variable $E_j$. For that purpose we use the required $\beta$-independency of $E$ and $E_j$ to write

$$E_j(\omega) = \frac{d}{d\beta_j} \sum_j \beta_j E_j(\omega) = -\frac{d}{d\beta_j} \ln Z(\omega, \beta),$$

which includes the case of a one component $E$. For effective, i.e. $\beta$ dependent energies $E_i^{eff}(\omega, \beta)$ this would give

$$E_j(\omega, \beta) + \sum_i \frac{d}{d\beta_i} E_i(\omega, \beta) = -\frac{d}{d\beta_j} \ln Z(\omega, \beta).$$

Analogously, using Eqs. (13) and (15) we find the same relation for general $A$

$$E_j(A, \beta) = < E_j(\omega) >_{A,\beta} = -\frac{d}{d\beta_j} \ln Z(A, \beta),$$

or for effective energies

$$E_j^{eff}(A, \beta) + \sum_i \frac{d}{d\beta_i} E_i^{eff}(A, \beta) = -\frac{d}{d\beta_j} \ln Z(A, \beta).$$

Assuming these averages to be measurable, this relation might be used to test the range of validity, i.e. the range of $\beta$-independence of an effective energy.

Linearity of the expectation allows to use this also to calculate averages of random variables which are linear functions of $E(\omega)$. To calculate expectations and higher order moments or cumulants of general random variables $h(\omega)$ which are either nonlinear functions of $E(\omega)$ or which are no functions $h(E(\omega))$ of $E(\omega)$, like e.g. if $E(\omega) = E(h(\omega)^2)$, one can extend the partition function by adding $-\lambda h(\omega)$ ('auxiliary field') to the exponent

$$Z(\omega, \beta, \lambda) = e^{-\beta E(\omega) - \lambda h(\omega)},$$

so that

$$h(A, \beta) = -\frac{d}{d\lambda} \ln Z(A, \beta, \lambda)\big|_{\lambda=0}.$$

In general we find

$$\frac{Z(A, \beta, \lambda)}{Z(A, \beta)} = < e^{-\lambda h(\omega)} >_{A,\beta} = \sum_k^\infty \frac{(-\lambda)^k}{k!} < h^k(\omega) >_{A,\beta}$$

which therefore is the *moment generating function*, i.e. the $k$th-moments $M_k(h, A, \beta) = < m^k(\omega) >$ can be found by differentiation

$$M_k(h, A, \beta) = \frac{d^k}{d(-\lambda)^k} \frac{Z(A, \beta, \lambda)}{Z(A, \beta)}\big|_{\lambda=0}.$$

Cumulants are defined as derivatives of the logarithm of the generating function, also called *cumulant generating function*,

$$C_k(h, A, \beta) = \frac{d^k}{d(-\lambda)^k} \ln Z(A, \beta, \lambda)\big|_{\lambda=0},$$

where we skipped the $\lambda$-independent term $\ln Z(A, \beta)$ which is not relevant after differentiation. It is easy to see that cumulants have the nice property of being additive if

the probability $p(\omega) = \prod_i p_i(\omega)$ factorizes into independent subsystems[32] (Beck & Schlögl, 1993). The second cumulant is the well known variance $< h^2 > - < h >^2$.

For the case of a Gaussian $p(\omega, \beta)$, i.e. $E(\omega)$ quadratic in $\omega$, (multidimensional for vector $\omega$) *Wicks theorem* gives higher order moments, often represented by diagrams (see for example Negele & Orland, 1988, Zinn–Justin, 1989, Itzkyson & Drouffe, 1989 ).

### Boltzmann–Gibbs–distributions

Now we want to relate the shifted log-probabilities $-\beta F(\omega)$ to the expectations of the generating random variable $E_j$. In general, for $A \neq \omega$ the difference between energy and free energy (i.e. between averages and shifted log-probabilities, coinciding on the elementary events $\omega$) reads

$$
\begin{aligned}
H(A, \beta) &= \beta E(A, \beta) - \beta F(A, \beta) \\
&= -\sum_{\omega \in A} \frac{p(\omega, \beta)}{p(A, \beta)} \ln \frac{p(\omega, \beta)}{p(A, \beta)} \\
&= -< \ln \frac{p(\omega, \beta)}{p(A, \beta)} >_{A, \beta} .
\end{aligned}
$$

Here $H(A, \beta)$ is the *average information* (in nats, not bits) or *entropy*. The free energy can be thus be expressed by the average energy and entropy as

$$
F(A, \beta) = E(A, \beta) - \frac{1}{\beta} H(A, \beta).
$$

It is well known, that the distribution generated by $\beta E$ maximizes the entropy under the constraints that $p(\omega)/p(A)$ is normalized and the expectations of all expectations $E_j(A, \beta)$ are fixed (For given vector $\beta$ also formulated as principle of minimum free energy $F$. See e.g. Balian, 1991, or in connection with path integrals Roepstorff, 1991). Indeed, the stationarity equations read, introducing the corresponding Lagrange multipliers $\beta_j$, $\alpha$

$$
\frac{d}{dp(\omega, \beta)} \left( H(A, \beta) - \sum_j \beta_j E_j(A, \beta) - \alpha < 1 > \right) = 0,
$$

giving the *Boltzmann–Gibbs distribution*:

$$
p(\omega, \beta) = \frac{e^{-\sum_j \beta_j E_j(\omega, \beta)}}{Z(\beta)},
$$

with $Z(\beta) = e^{\alpha(\beta)+1}$. We recognize the general form (15). Thus we can say a probability distribution is the Boltzmann–Gibbs distribution of its generating variable(s) $\beta E$. The distribution for $E$ is

given by summing over $\omega$ with energy $E$ (which could be a vector), i.e. $p(E) = p(\omega_E) = Z(E, \beta)/Z(\beta) = \int d\omega \, \delta(E(\omega) - E) e^{-\sum_j \beta_j E_j(\omega) - \ln Z(\beta)}$
$= n(E) e^{-\sum_j \beta_j E_j - \ln Z(\beta)} = e^{-\sum_j \beta_j E_j - \ln Z(\beta) + \ln n(E)}$,
with $\omega_E = \{\omega | E(\omega) = E\}$ and $n(E)$ the energy density at $E$.

Families of distributions generated by varying (the vector) $\beta$ within some parameter space $B$ are called *exponential families*, with canonical parameter vector $\beta$, canonical statistics $(-E)$ and cumulant (generating) function $\ln Z$. (See for example Barndorff–Nielsen, 1978, Amari, 1985, 1995, or Appendix D in Lauritzen, 1996 and references therein).

The model from Section 2 defines an elementary event $\omega = (q, y, f^0)$ or, if we include the internal variables $y^q z^q$ of questions $q$, $\omega' = (q, X^q, Y_{X^q}, Z^q, f^0)$, with $q \in Q^l$, $X^q$ the part of the basis used in $Q^l$, $Y_{X^q}$ the corresponding answers for $x \in X^q$, and $Z^q$ the set of internal noise variables for $q$. We are for example (in the case of deterministic $\hat{l}$ and of $p(q)$ not dependent on other variables) interested in calculating the expected risk $r = < l(\omega, \hat{f}) >_{p, \Omega}$ under the total posterior probability

$$
p(\omega) = p(q) p(y|q, f^0) p(f^0|D) = e^{L(q) + L(y|q, f^0) + L(f^0|D)},
$$

in terms of log-probabilities $L(\omega) = L(q) + L(y|q, f^0) + L(f^0|D)$. Thus we can interpret the Bayesian posterior distribution as Boltzmann–Gibbs distribution arising from a maximum entropy procedure for the total log-posterior $L(\omega)$ with an average so that $\beta = -1$ and shifted so that $Z(\Omega) = 1$. It seems simpler, to use the traditional Bayesian formulation than the equivalent maximum entropy formulation. However, there are cases where averages are directly measurable, with measurement error nearly zero.[33] Besides deterministic variables, self–averaging variables in large systems belong to this class. Examples, are macroscopic observables in physics, like energy, which lead to the (generalized) canonical ensembles used in statistical physics. An indeed, in these cases the various implementations (i.e. models $F^0$, like microcanonical, canonical, grandcanonical ensemble) are asymptotically equivalent for large systems.

### Approximation and Kullback–Leibler entropy

Consider another probability distribution $p'(\omega)$ generated by $E'(\omega) = -\frac{1}{\beta'} \ln p'(\omega, \beta') + F'(A, \beta') = -\frac{1}{\beta'}(\ln p'(\omega, \beta') - Z(A, \beta'))$, with the same normalization, i.e. $Z(A, \beta) = Z'(A, \beta')$, so the expectation of $E'$ under $p(\omega)$ reads

$$
E'(p, A, \beta) = < E' >_{(p, A, \beta)}
$$

$$
= - \left( \frac{1}{\beta'} < \ln p'(\omega, \beta') >_{(p, A, \beta)} + F(A, \beta) \right) p(A)^{-1},
$$

---

[32] For example spatial systems with a weakly enough short range interaction, so they can be seen in the large $n$–limit as a collection of independent local subsystems, have cumulants which are asymptotically additive in these subsystems. In these cases, besides the $p$–generating energy, the free energy $F \propto \ln Z$ and not $Z$ is an additive variable. As the variance of a sum of $n$ independent random variables scales with $1/n$ the cumulants have then for large systems low variance, hence they are proportional to the volume (extensive variables) and possible candidates for macroscopic observables of the system.

[33] Therefore, for example Jaynes, 1996, sees the two approaches as two different methods, applied in different situations: Maximum entropy to fix averages which are known without much computation, and Bayesian methods dealing with models to calculate the relevant average.

where $E'$ is averaged over $p(\beta)$ and not over $p'(\beta')$ generated by itself. Concavity of the logarithm function allows to apply Jensen's inequality to the difference $\beta'E'(p, A, \beta) - \beta E(A, \beta) = <\ln p> - <\ln p'> = K(p, p')$, i.e. the *Kullback–Leibler entropy* and we obtain

$$\beta'E'(p, A, \beta) \geq \beta E(A, \beta).$$

Thus, the $\beta E$ generating the averaging distribution has the minimal average under this distribution. $\beta'E'$ represent different loss functions. We define an *approximation problem* to be the problem of minimizing the expectation (expected risk) over $A$ under $p$ for a family of loss functions, defined on the same set $\omega$ and having all the same normalization $Z$ over $A$. Those loss functions can be parameterized in the form $\beta'E' = -c_1 \ln p'(\omega) - c_0$ with $p'$ a normalized probability density. We will call this a family of *approximation losses*. Then, for any parameterized family of $\beta'E'$ the probability related to the solution with minimal expected risk has minimal Kullback–Leibler distance to the actual averaging probability distribution.

Eq. (5.3.4) defines the true expected risk $\beta E$ as the solution of a minimization problem. This corresponds to a *variational method* for calculating the true expected risk (see e.g. Balian, 1991, Neal & Hinton, 1993, Dayan, Hinton, Neal, & Zemel, 1995):

1. Besides the effects of not included parameters, the difference between the true expected risk $\beta E$ and an expected risk $\beta^* E^*$, minimal in a parameterized subspace, has only contributions of second order in the parameters used for minimization. 2. The true expected risk is bounded by the approximated risk.[34]

## 5.4 Including human knowledge and other available preprocessors

A human interface for constructing fuzzy priors consists of two steps:

1. Subsymbolic level:

    i. Defining properties $C_i \in \mathcal{C}$, i.e. functions of a parameterization of $f^0$, as correlates to internal concepts $\tilde{C}_i$. For example, a property can be specified as a typical variant of an eye or a typical structure of an electrocardiogram, but as well by the output of another available (e.g. approximation) algorithm. The properties $C_i$ are in some contexts also called linguistic variables (Zadeh, 1996).

    ii. Definition of possible combinations of properties, i.e. mappings $\mathcal{C} \times \mathcal{C} \to \mathcal{C}$ or linguistic functions, approximating internal mappings, i.e. the structure of concepts. Specifically, a linguistic 'not', 'and', 'or' can be related

to some real valued extension of the logical NOT, AND, OR.

2. Symbolic level: Combining simple properties to create complex ones, i.e. applying linguistic functions (rules) to linguistic variables according to the communicated symbolic structure of the prior. For example the property being a face like object can be build up from properties of having two eyes, a nose and a mouth like object. Probabilistic rules are a special set of linguistic functions which depend on a (communicated) dependence structure of the variables.

Consider we want to construct a property $C$ (deterministic question) of $f^0$, describing prior knowledge about images of faces, out of subproperties $C_i$ by fuzzy operations. We can choose for example a log-prior

$$L \propto -||C - T_C||^2 + c.$$

Also, the subproperties $C_i$ can be deviation properties

$$C_i \propto -||C_i' - T_i||^2 + c_i.$$

For probabilistic questions one might wish to integrate over $y_q$ with $p(y_q|q, f^0)$

$$C_q \propto -\int dy_q\, p(y_q|q, f^0)||y_q - T_{y_q}||^2 + c.$$

The templates $T_C$, $T_i$, $T_{y_q}$ can be chosen as output of an available approximation $\hat{y}_q$, e.g. $T_{y_q} = \hat{y}_q$. The approximation $\hat{y}_q$ may be produced from a previously trained artificial neural network, or any other statistical approximator, as well as from an expert. A construction according to fuzzy methods can be as follows: Let $p(y|q, f^0)$ be the probability that the given image $q$ is a face (i.e. y=face). (That means, this is the classification and not the generation probability for faces.) Let $T_{y_q}$ be the answer of some already available face detector. Also, several $T_{y_q}$, i.e. several answers or different approximation methods, can be included at the same time, according to their known or subjectively believed dependency structure. The dependency between approximators may be arbitrary including as special cases approximations which are independent or disjunct (only important for OR not AND).[35] The template $T_{y_q}$ itself can be constructed according to fuzzy methods: Define ($f^0$–independent) $q$–templates $T_q$ for $q$ (Not $y_q$–templates $T_{y_q}$ !). They may correspond to typical constituents like eyes, mouth, nose, and look like $\sum_i ||q_i - T_i^{eye}||^2$ with $i$ denoting the pixel index. Define transformation of templates, including at least translation (represented by a

---

[34]Variational principles for minima (or maxima, respectively) have the advantage of giving such a bound, while variational principles related to saddle points (e.g. for complex functions), do not. On the other hand, for minima (or maxima) all second order corrections have the same sign, while for saddle points second order contributions with different sign can average away (see for example, Lemm, 1995ab). Especially in high-dimensional spaces this may considerably improve the approximation.

[35]The problem of combination of different approximations of $y_q$ (or more general of arbitrary available data) to get a better approximation of $y_q$ is just a version of the usual statistical approximation problem. Thus, a great number of possible algorithms is available to deal with the problem. Some methods refer especially to the situation where the data are approximations of the same value (combination of experts/approximators/classifiers). Tree–like methods construct a local 'domain of responsibility' for each expert. (See for example the mixture of experts, (Jacobs, Jordan, Nowlan, & Hinton, 1991; Jordan, Jacobs, 1994).)

transformation of the index $i$ of the template). Combine the constituents, their possible variations and combinations with fuzzy AND, OR, NOT (which includes more complicated rules like XOR or IF $\cdots$ THEN) to obtain a final property $T_{y_q}$. This defines a fuzzy face classifier, incorporating human prior knowledge.

We now discuss two principal approaches to incorporate preprocessor information, like the fuzzy templates, into a subsequent algorithm.

### 5.4.1 Two coupling principles

The loss minimizing algorithm makes the decision $\hat{f}$. Such an optimizer must have an interface for relevant data, i.e. for the $q^l$ and corresponding answers $y$. This allows to distinguish two principal variants to include the information of a preprocessor (See Fig.7):

1. Feeding the output of the preprocessor in the given interface for relevant data, including the entrance for the $y$ which define the goal of learning. We will call this prior cascade and the preprocessor a data model generator (Fig.7 top).

2. Feeding the preprocessor output to extra input channels of the subsequent optimizer. We will call this case input cascade. Here the preprocessor has not necessarily to produce a data model. It can be implemented in two variants:

   a. asynchronously with the data $D^l$ (Fig.7 middle),

   b. synchronously with the data $D^l$ (Fig.7 bottom).

   $q$-variables which enter the loss function only indirectly through $\hat{f}(q)$, can then be skipped. This allows to effectively replace $q$ by preprocessor output.

We now discuss the two approaches in more detail.

### 5.4.2 Prior cascade (data modeling preprocessor)

In a prior cascade, a data modeling preprocessor intends to produce a (fuzzy) model $f$ for the input of the optimizer, i.e. of the relevant data $D^l$. A Bayesian model is a special model formulated in terms of probabilities. Accepting its validity, the model can be used to create new virtual examples $D^{V,l}$ for relevant questions $q \in Q^l$. For the corresponding relevant data $D^l$ a loss function is defined fixing their 'interpretation, for the optimizer. For virtual examples $D^{V,l}$ it is assumed that the same interpretation, i.e. loss function, applies, so they can use the $q$ and the $y$ entrance of the optimizer. Thus, the flexibility of this approach consists in the fact that the $D^{V,l}$ fit a given input format and a given interpretation, so they do not require adaption of the optimizer. Such a preprocessor may be seen as minimizing (implicitly or explicitly) an approximation loss for the data. This does not necessarily coincide with the loss used by the optimizer, which may, for example, include additional complexity penalties. Indeed, if the optimizer would have to minimize the same loss as the preprocessor, which would include having the same architectural restrictions, it cannot produce new results, and it is necessarily $f = \hat{f}$, as long as no
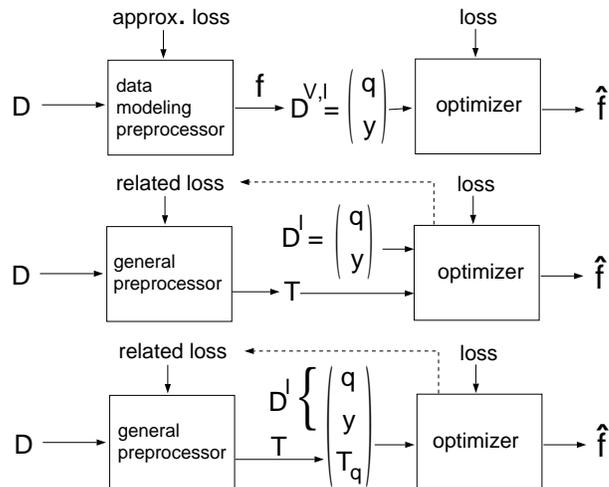


Figure 7: Top: Prior cascade (data modeling preprocessor). The preprocessor uses the data interface of the optimizer. So the interpretation of its information as (virtual) data is implicit. Thus, 1. the preprocessor must represent a data model, 2. the optimizer has not to be adapted. (approx. loss = approximation loss)

Middle: Input cascade (general preprocessor). The preprocessor uses not the whole interface for loss relevant data, i.e. at least not the entrance for $y$ which are the variables for the loss function defining the goal of learning. It can, however, for loss which is not explicitly $q$ dependent (but always indirect over $\hat{f}$), replace part or all of $q$. Thus, 1. the preprocessor does not need to represent a data model, therefore its loss can also be adapted according to the actual needs of the optimizer, 2. however, the optimizer must in general be adapted if the preprocessor adds information.

Bottom: Special case of an input cascade: The input $q$ is changed, either by replacing $q$ by preprocessor output (only if the loss is $q$ independent, e.g. multilayer neural network) or by adding new dimensions (an extreme case is the support vector machine, with the preprocessing implicit in the kernel). Those additional dimensions can also be the output of a data modeling preprocessor and represent an approximation for $y$ or for the optimal reaction $\hat{f}(q)$. (Which is the same in a pure approximation problem.) For a more detailed explanation see text.

new data are included. If new data are available the pre-processor generates, only part of the optimizer's input, weighted according to its prior probability. Approximation loss is related to any non–approximation loss which depends on the relevant data.

Consider now a fuzzy preprocessor, like the fuzzy face detector we described in principle. There, new training examples can be generated sampling according to a term $C_q \propto -\int dq p(q) \|y_q(f^0) - T_{y_q}\|^2 + c$ in the log-prior (for probabilistic $f^0$ we integrate also over $p(y_q|q, f^0)$). The approximation $T_{y_q}$ for $\hat{y}_q$ itself can be based on approximations implemented by virtual examples. Note, that virtual $q$ can include questions which are not available as training questions, as long as a loss function is defined for them. For a face detector, those $q$ can correspond to images which appear neither as faces nor as non–faces, i.e. which are no 'natural' images ($= Q^l$). Such 'non–natural' $q$ can be images of faces with some subconcepts missing or exchanged, for example, one eye transformed to a none-ye, e.g. to a square. Such a squared eye face is probably not in the set of 'natural' images $Q^l$ but may be useful to teach the face detector the concept of an eye and its relevance for face detection. Including those examples near the border between faces and (the concept of, not the natural) non–faces, specifies the generalization intended by the trainer.

The $q$–integral in the log-prior can be treated in several ways:

1. the integration may be performed *analytically*, and the integral (infinite sum) replaced by an equivalent term, which is easier to evaluate,

2. by *virtual examples*, i.e. pairs $(q, T_{y_q})$, e.g. fuzzy classified (See e.g. Lin, Kung, & Lin, 1997). In case of many virtual examples this estimation may be called a numerical evaluation of the integral.

3. by *application adapted sampling*, i.e. virtual examples are (also) used during the application phase of the algorithm. That means, in a specific application situation $q^l$ one includes newly sampled case–specific virtual examples, before answering to $q^l$. For example, one can evaluate the corresponding $T_{y_{q^l}}$ (if this can still influence the output and if available), and can also include $T_{y_{q'}}$ for other, e.g. 'similar', $q'$. We remark that this variant

    a. assumes the availability of the approximation $T_{y_q}$ also during the application phase and not only during training,

    b. requires on–line '*application–adapted relearning*', because the newly, case specific generated examples have to be incorporated in the learning process. An trivial example of application specific on–line 'learning' is interpolating between different available answers, a less trivial example are neural networks which are retrained with application specific 'hints' $T_{y_{q^l}}$ before giving an answer to a specific $q^l$. Also, a higher level approximator can include local prior terms according to the output of a parallel working fuzzy face detector for a given image, preclassifing faces and non–faces.

    c. is equivalent to the use of an *infinite* number of virtual examples, because for every $q^l$ (i.e. arbitrarily many) new examples are drawn. That does not mean that a single prediction depends on an infinite number of virtual examples, but the definition of the whole learning machine, and therefore its expected performance.[36]

### 5.4.3  Input cascade (general preprocessing)

In an input cascade, a general preprocessor produces output $T$, for which no loss function, i.e. prewired interpretation, has to exist. This kind of information $T$ cannot use the relevant data entrance of the optimizer. In cases, however, where the loss function is not explicitly $q$–independent (not $y(q)$– or $\hat{y}(q)$–independent!) the $q$–part of the data can be completely skipped and replaced. This does not include the $y$ variables which define the value of the loss function and therefore the goal for the optimizer. Thus, the output of a preprocessor can replace part of the $q$ or can be added as additional input to the optimizer, which has to be adapted to the new format of the input. This internal adaption defines the interpretation of the additional data.

In principle the total (information of the) previous approximator $T$ could be added as new input dimension. However, the increase in complexity by irrelevant data tends to lower the performance in an input cascade (See below). Also a data modelling preprocessor can be used in an input cascade by adding available virtual data pairs to $q$, leaving $y$ unaltered.

Especially, the information of the preprocessor may be feed in synchronously with every data pair $y, q^l$. The transformed and often enlarged input space is also called feature space. Optimizing algorithms can often be adapted relatively easy for higher input dimension of $q$. Thus, a simple form of an input cascade or preprocessing is including part of the information $T$ in every data pair $(q, y)$. This can be the result $f'$ of another optimizer, but also linear independent components of the vector $q$, or distances to prototypes obtained by unsupervised algorithms. We mentioned that in cases where the loss function is explicitly $q$–independent the $q$–part of the data can be completely skipped and replaced by $T_q$, so that the original data entrance formally fits. This is the standard case of an hierarchical optimizer, like for example a neural network.

In practice one can for example add output $T_{y_q}$ corresponding to the same (and maybe a few related) $q$, i.e. $q^l \rightarrow \{q^l, T_{y_q}\}$. (Also aspects of the learning history might be added, but one usually assume this to be well enough represented by the actual internal state of the

---

[36]Again, infinity means nothing else than assuming to be always able to do something if needed, like in this case creating a template for every new input. That is how we defined the algorithm and if there are cases where we cannot create corresponding templates, then we are just not able to apply this algorithm, nevertheless we say the infinite data statement is true for the algorithm by definition. (definition = specifying possibilities of control).

optimizer). If $T_{y_q}$ for a given $q^l$ is not available, one may use a $q$ where a similar $T_{y_q}$ is expected, i.e. complete the template generator by adding another prior assumption. For example assuming smoothness with respect to a distance $||\cdot||$, one can use a $q$ with $\text{argmin}_q ||q^l - q||$ ('nearest' $q$) or interpolate between some of the neighbors. One may encode that a specific value is not available, or its expected correctness. For added as well as for replaced input a function $\hat{f}$ of the augmented and not of the original input is learned. That means, the preprocessor also has to be available during application of $\hat{f}$.

For $T$ to be of any help, the preprocessor, must deal with aspects related to the loss function, i.e. it must have (implicitly or explicitly) a loss function related to the loss function of the optimizer. For example, the preprocessor can explicitly use the same loss function, and represent therefore a previously trained optimizer. As the loss function of the preprocessor does not have to produce a data model, it can be adapted to the state of the optimizer. A typical example of an input cascade with feedback of the loss functions is a multilayer neural network. For all 'on–line' optimizers which do not use the complete available data set $D^l$ at the same time (e.g. backpropagation), the actual parameter values (e.g. weights in a neural network) representing the memory for past data, can be seen as output of a preprocessor. As any architectural restrictions can be related to prior data, any restricted optimizer may be seen as result of some preprocessing.

In a prior cascade the next level algorithm treats the additional information as relevant data, i.e. the Bayesian interpretation of the template data is hardwired. In an input cascade the algorithm is free in its use of the data coming from the preprocessor. More precisely, the preprocessor data have a meaning implicitly implemented in the optimizer's algorithm and its architecture. Consider the extreme case, where the answer, the optimizer is looking for, is already included as an additional input dimension. This is only of any help if the algorithm of the optimizer has with a reasonable probability in its space $\hat{F}$ of possible hypotheses a function which is the projection of the whole input into this dimension. This includes that, even if all other original $q$–dimensions are deleted, the identity (which might be parameterized very complicatedly) must be part of the hypothesis space to allow the algorithm to find this 'simple' solution. Most practically used algorithms can probably find the identity or a projection easily, but one can also easily construct a model, which cannot. Formulated more tautological, additional information is helpful if it increases the probability of finding a better solution. (The meaning of 'better' can be specified in many variants.) The solution found by an algorithm must be part of a (learning history dependent) space $\hat{F}^c$ of hypotheses, which have been actually considered as possible solutions by the algorithm. A preprocessor can transform the input so that the transformed optimal solution $\hat{f}^*$ is (averaged over possible learning histories, if the templates are added before seeing all $q^l$, and random variables of the algorithm) with high probability in the set of consid-

ered hypotheses $\hat{f}^* \subset \hat{F}^c$. For example, for algorithms for which the projection is easily learnable, adding templates for the output as additional input variable might be a good choice.

If we assume the space of considered hypotheses $\hat{F}^c$ to be bounded, with respect to some resources, then an input transformation effectively interchanges (maybe probabilistically) some considered with non–considered hypotheses. Thus, changing the considered hypotheses by additional information, means changing one factor within a multiple OR (of all $\hat{f}^c$). Realizations of a multiple OR are usually very flat functions within the set $\hat{F}^c$ (e.g. constant in the deterministic case) with (more or less) sharp transition to zero at the boundary. Thus, an input cascade is expected to have a very flat prior being only effective at the cutting edge. As a (highly) multiple OR is difficult to implement, information $(D)$ which seems meaningful, without being easily related to $p(f^0|D)$ can easier be included in an input cascade, (i.e. used for preprocessing) than in a prior cascade (i.e. data modeling). This is usually the case in non–approximation problems when information is related to the optimal answer and not to the possible state (See also Section 8).

Notice, that in a prior cascade we can change the importance of priors by changing their weight factor. Because this is a parameter of a Bayesian model, there is usually no direct analogous parameter for the input cascade, so the importance of information is encoded implicitly in (the parameters describing) the algorithm. Assuming the Bayesian model to be correct, the prior cascade should give better results. In cases the Bayesian model is not correct, the input cascade has less bias, as it is similar to a prior cascade with smaller weights for the non–correct priors. However, adding input variables makes the problem more complex by increasing the space of possible solutions. For example the VC dimension may increase, or the (algorithm specific, implicit) prior for the true state of nature might be smaller in an enlarged space. Alternatively, the complexity of the approximator may be reduced, for example by requiring a stronger smoothness in higher dimensions. So for example the support vector machine increases enormously the input dimensions (feature space), while on the other hand their VC dimension is expected to stay approximately constant. Compensating the addition of new variables by restriction of the search space $\hat{F}$, is only successful if the restrictions reflect correct priors.[37] They should therefore be the result of a prior cascade, which suggests the possibility to include this information also without increasing the input dimension.

Thus, adding input variables, is only expected to help,

---

[37] In the support vector machine, the restrictions imposed by constructing an optimal hyperplanes in the feature space, can, choosing the appropriate kernels, lead to arbitrary restrictions in the original input space. Thus, the selection of the kernel includes the necessary prior knowledge. In the extreme case were the optimal solution (a one dimensional binary variable in binary classification) is presented as template, a feature space with the dimension of that solution is sufficient to find the optimal hyperplane.

if they contain enough information to compensate either for the higher dimensionality or for skipped variables in the problem. In contrary, for a prior cascade also including variables with only small effects should improve the performance, as long as those effects are correctly modelled. An input cascade implemented by adding (or replacing) components to the relevant data is an application adapted sampling method which requires on–line availability of the preprocessor producing $T_{y_q}$, and therefore corresponds for an infinite set $Q^l$, to a formally infinite amount of data.

These methods of an integration of available approximators, allow to use information contained in other learning systems, and are therefore general methods of knowledge transfer.

## 6 Decision problems

### 6.1 Definition

In this section we study questions conditioned on data. Assume we have to choose one question (alternative) $\hat{f}$ out of a set of questions $\hat{F}$[38]. We will now use the symbol $l$ (actual loss) for the answers of $\hat{f}$ and assume that it contains all decision relevant information. Then a decision should only depend on the loss distributions of the various $\hat{f}$. We choose a minimum on the set $\hat{F}$ by defining a (risk) functional $r[p(\cdot|\hat{f}, f)]$ mapping probability densities of answers $p(l|\hat{f}, f(D))$ into a subset of the real numbers, bounded from below. A *decision problem* consists in finding the answer to a question $q_r$ with

$$\hat{f}^* = q_r(f(D)) = \operatorname{argmin}_{\hat{f} \in \hat{F}} r[p(\cdot|\hat{f}, f)].$$

If we want to emphasize the data dependency of the decision we speak of a *learning problem*. Approximation and classification problems are special cases of decision problems.

We will now have a closer look to the questions $\hat{f}$. Using $p(y^q|q, f) = \int df^0 p(f^0|f) p(y^q|q, f^0)$ we will from now on write the formulas for states of knowledge $f$ instead for pure states $f^0$. According to Eq.(3) there exist some $q$ so that the probability of suffering loss $l$ in a given state of knowledge $f$ can be written

$$p(l|f, \hat{f}) = \int dq \int dy \int dz \, p(q|y_c, z_c, \hat{f})$$

$$\times p(y|q, f) p(z|q, y, \hat{f}) \delta(l(q, y, z, \hat{f}) - l).$$

This implies that noise variables $z$ of different questions are independent and one specific realization of $y$ according to $p(y|q, f^0)$ can only appear multiple times within one question. The situation can be represented by influence diagrams with decision and value nodes (Pearl, 1988) We can define the decision relevant basis

set $X = X^l$ by the set $Q^l$ of test data $q^l = q$. Then this formula is in the form usually given in Bayesian decision theory. Starting from an factorial state with respect to $X^l$ structural information relates generalized questions $q^D$ from the training data to the different $q = x$ and all information enabling generalization has to come from nonlocal $q^D$.

The explicit $\hat{f}$–dependence comes from three factors which are

   i. the action $z$ (including noise) producing device $p(z|q, y, \hat{f})$,

   ii. the defining (loss) function $l(q, y, z, \hat{f})$,

   iii. the test set generator $p(q|y_c, z_c, \hat{f})$.

One usually uses a formulation where the components in 2. and 3. are chosen to be explicitly $\hat{f}$–independent.

  1. a) $p(q|y_c, z_c, \hat{f}) = p(q|y_c, z_c)$, or

     b) $p(q|y_c, z_c, \hat{f}) = p(q|y_c)$ (fair), or

     c) $p(q|y_c, z_c, \hat{f}) = p(q)$ (static),

  2. $l(q, y, z, \hat{f}) = l(q, y, z)$.

Remark:

  1. We always can fulfill conditions 1 a) and 2 by introducing additional $\hat{f}$–dependent variables $z$ (and a zero dimensional $x_1$) without effectively changing the model, as $\hat{f}$ has always to be parameterized. Therefore we defined $\hat{f}$ in Section 2 by $p(z|q, y, \hat{f})$. To fulfill the stronger condition $p(q|y_c, z_c, \hat{f}) = p(q|y_c)$ also all random processes have to be attributed to the state by including $z_c$ into $y_c$, and for the strongest condition $p(q|y_c, z_c, \hat{f}) = p(q)$ by including $y_c$ and $z_c$ into $q$.

  2. This covers all cases of function approximation and pattern classification.

We will call a decision problem *fair* if condition 1b is fulfilled so different action states $\hat{f}$ are compared in the same situations $q$ independent of their own previous outcomes, and *static* if condition 1c is fulfilled, and $p(q)$ does not depend on any other outcomes.

If interested in the expectation of $l$ the remaining $z$ can for static problems be integrated out defining an integrated loss function $\tilde{l}$

$$\tilde{l}(q, y, \hat{f}) = \int dz p(z|q, y, \hat{f}) l(q, y, z, \hat{f})$$

being a deterministic function of $\hat{f}$.

### 6.2 Parallel decision problems

If we assume $p(z|q, y, \hat{f})$ to be $y$–independent this can be seen as a device producing answer (action) $z$ in situation $q$ when the model (of nature) in state $f$ produces $y$. Thus, $\hat{f}$ could be interpreted as (action) state of a second, independent model of actions capable in producing answers to the same questions as the original model (of nature). We will now write $\hat{y}$ for those $z$ which are produced independent of $y$. This independence can also be

---

[38]We use the letter $\hat{f}$ instead of $q$ for this set of questions to symmetrize the further notation. This corresponds to an interpretation of the selected question as action state $\hat{f}$ in analogy to the model state $f$. The selected action state $\hat{f} = \hat{f}(f)$ can be seen as reaction (e.g. approximation) to the state of knowledge $f$.

seen as definition of the $q$ as that part of the visible variables $(q, y)$ which can influence the answer $\hat{y}$. For example, the answer of the model $y$ can be assumed to be available for the device $\hat{f}$ only after the action $\hat{y}$ is produced. Now let us define a decision problem which compares two models in states $f$ and $\hat{f}$ answering to the same questions $q$ under equal conditions. This can be seen as comparing the action state $\hat{f}$ with the state of nature $f$ in situations $q$ according to the criteria $l(q, y, \hat{y})$.

If we want to model a causal dependence of the action $\hat{y}$ from 'previous' $y$, we choose the general formulation as

**PM**: a *parallel decision problem with memory*, defined as a model with $p(z|q, y, \hat{f}) = p(z|q, y_c, \hat{f})$. with the causal structure $y_c$ being the same as in $p(q|y_c, \hat{z}_c)$. Then we write $\hat{y}$ for $z$, i.e. $p(\hat{y}|q, y_c, \hat{f})$. (The notation allows dependencies within the components of $\hat{y}$.) Thus, the action producing device outputs $\hat{y}$ in action state $\hat{f}$ with its components $\hat{y}_i$ independent of the corresponding answer component $y_i$ of the model of nature $f$. However, it can use the values of previously determined components of $y$ and choose its actions according to its 'success' in the past.

We call a decision problem $\hat{F}$

**P** : a *parallel decision problem without memory* (or simply a parallel decision problem) if, writing $\hat{y}$ for $z$, the action model produces $\hat{y}$ in action state $\hat{f}$ independent of the answer $y$ of the model of nature $f$: $p(\hat{y}|q, y, \hat{f}) = p(\hat{y}|q, \hat{f})$ (Dependencies within components of $\hat{y}$ are allowed).

This leads to a loss probability, written for the more general case:

$$p(l|f, \hat{f}) = \int dq \int dy \int d\hat{y}\, p(q|y_c, \hat{y}_c)$$

$$\times p(y|q, f) p(\hat{y}|q, y_c, \hat{f}) \delta(l(q, y, \hat{y}) - l),$$

Note the difference between not yet seen (expected) test data $\hat{D}^l = (q, y, \hat{y})$, which are not yet determined and correspond to the integration variables, and the known (training) data $D$, which determines our actual state of knowledge $f = f(D)$ and do not appear explicitly in the above formula. Parallel decision problems with memory show an asymmetry between $f$ and $\hat{f}$ as we allow a dependence of $\hat{y}$ from past values of $y$ while we defined $y$ only to be depending on the question $q$ and pure state $f^0$ and not from $\hat{y}$. All $\hat{y}$–dependency of real world measurements $p(y|q, f^0)$ by definition has to come from changing the question according to $p(q|y_c, \hat{y}_c)$ or active changes of $f^0$ by $\hat{y}$ which we do not allow. We can parallel the treatment of $\hat{f}$ and $f$, by introducing a set of $y$–independent basis action states $\hat{f}^0$ and an 'algorithm' $p(\hat{f}^0|q, y_c, \hat{y}_c, \hat{D})$. Here we included possible additional data $\hat{D}$ in the formalism, which also determine the available action states $\hat{f}$, and we can write

$$p(\hat{y}|q, y_c, \hat{f}, \hat{D}) = \int d\hat{f}^0 p(\hat{y}|q, \hat{f}^0) p(\hat{f}^0|q, y_c, \hat{y}_c, \hat{D}).$$

For $p(\hat{D}|D) = p(\hat{D})$ we can determine $\hat{f}$ before the training begins, otherwise, the set $\hat{F}$ of available choices has also to be updated with the training data $D$.

*Approximation problems* are parallel decision problems with approximation loss (see Section 5.3.4), parameterized in the form $\tilde{l} = -a \ln p(y|q, \hat{f}) + c$ with $\int dy/, p(y|q, \hat{f}) = 1, \forall q, hatf$. With a temporal connotation they may also be called prediction problems. In such problems, we may wish to include the structure of the $q$ into the action $\hat{y}$–producing devices. For the set of $\hat{f}^0$, which in a parallel decision problem without memory already are the $\hat{f}$, we can then define an action model $p(y^x|x, \hat{f}^0)$ for the basis questions $x$ use for a question

$$p(y|q, f) = \int dx \int dy^x \int dz\, p(x|y_c^x, z_c, q)$$

$$\times p(y^x|x, f) p(z|x, y^x, q) \delta(q(x, y^x, z) - y).$$

an action state

$$p(\hat{y}|q, \hat{f}^0) = \int dx \int dy^x \int dz\, p(x|y_c^x z_c, q)$$

$$\times p(y^x|x, \hat{f}^0) p(z|x, y^x, q) \delta(\hat{q}(x, y^x, z) - \hat{y}).$$

Note that also a $\hat{f}$–independent loss function can penalize complexity or requirement of resources of the $\hat{f}$ simply by including corresponding variables as components into $\hat{y}$.

While the $\hat{f}$ define the possible alternatives of loss distributions, the decision has to be made with respect to some risk functional $r$ applied to this distributions. Choosing, for example, in the case of a real loss function for parallel decision problems the expectation as functional $r$ we have to minimize the 'expectation risk' or expected risk

$$r(f, \hat{f}) = r[p(\cdot|f, \hat{f})] = \int dl\, l\, p(l|f, \hat{f})$$

$$= \int dq \int dy \int d\hat{y}\, p(q|y_c, \hat{y}_c) p(y|q, f) p(\hat{y}|q, y_c, \hat{f}) l(q, y, \hat{y}).$$

*Density estimations* are a special form of parallel decision problems having only one $q$ with the meaning 'get next data'. Choosing an integrated loss which can for every $\hat{f}$ be interpreted as a log-probability, i.e. $\tilde{l}(y, \hat{f}) = \ln p(y|\hat{f}) + c$ this gives

$$r(f, \hat{f}) = \int df^0 \int dy\, p(f^0|y^D) p(y|f^0) \ln p(y|\hat{f}).$$

One may well include other $\hat{f}$–specific aspects in the loss function, like $\hat{f}$–specific complexity costs, e.g. a term $\int dy\, p(y|\hat{f}) \ln p(y|\hat{f})$ related to the encoding costs of $y$ given $\hat{f}$ (Rissanen, 1989). Sometimes the term 'unsupervised' learning is used for such problems, including problems which are defined by algorithms and do not explicitly refer to a set of parameterized models $p(y|f)$, like e.g. self–organizing maps (see e.g. Kohonen, 1995). This term might be misleading, because all variables $y$ enter the (explicit or implicit) loss function, and are therefore 'supervised' variables. The reason for using nevertheless

the word 'unsupervised' is, that density estimation is often applied to variables, which are used in a second step as $q$ in a problem, where the loss function does not depend on $q$ but on $y$. Indeed, the determination of $p(q)$ can in our formulation always be seen as a preprocessing step, because we explicitly require $p(q)$ to be independent of $f^0$.

To adapt to new data, parallel decision problems require an inversion to obtain

$$p(f^0|y^D, q^D) = \frac{p(y^D|q^D, f^0)p(f^0)}{p(y^D|q^D)}.$$

While this corresponds to an interchange of the roles of the variables $y^D$ and $f^0$, we will in the next Subsection refer to another inversion: The interchange of $y$ with $q$.

### 6.3   Inverse decision problems

Here we will discuss *analysis by synthesis*, where an action device $p(z|y, \hat{f})$ is 'analyzed', i.e. its corresponding loss minimized, under a synthesis model $p(y|q, f)$ for $y$, i.e. the question for $\hat{f}$ instead for its answer $z$. For example, instead of solving an approximation problem, which is a parallel problem yielding a predictive model $\hat{f}$ of $f$, we may be interested in approximating the question ('cause') $q$ yielding to $y$. Thus, we may want to approximate the inverse probability $p(q|y, f)$. Notice, however that in our terminology a decision problem with 'inverse approximation loss', parameterized by $\tilde{l} = -a \ln p(q|y, \hat{f}) + c$, with $\int dq\, p(q|y, \hat{f}), \forall y, \hat{f}$ is not called an approximation problem (except for the trivial case $p(y|q, f) = p(q|y, f)$). An approximation problem would require a normalization over $\forall q, \hat{f}$, i.e. in this case where $q$ and $y$ are interchanged of $\int dq/, p(q|y, \hat{f})$. Indeed, only for (parallel) approximation problems the inequality (5.3.4) can be applied for given $q$, and therefore, as we will see in Section 7 a maximum posterior approximation can be sufficient.

Often, the variable describing the desired output of the action device is easier to choose as condition $q$ to build possible models $p(y|q, f^0)$ of nature, than its input variable. Consider the following example: Let $y^a$ have the values 'face' and 'non–face' and $q^a$ be corresponding images of faces and non–faces. The task is to build a face detector, i.e. a device $p(\hat{y}^a|q^a, f)$ which outputs an approximation $\hat{y}^a$ of $y^a$ for a given image $q^a$. That would mean we need a model $p(y^a|q^a, f)$, i.e describe for a given image possible probability distributions being a face. This might be done using fuzzy priors, but might not be very reliable, as the probability for a given image to be a face is related to the kind of non–face objects included in $y^a$. Every change in the set of non–faces, would require an update of the fuzzy prior. Much more natural is it — as probably more related to the subjective process of generating priors by contemplating about typical features of faces — to approximate the inverse probability of images $q^a$ given it is a face $y^a$. This can be done (as well as for the inverse direction) by a fuzzy decomposition of the image $q^a$ into feature components $q^{a'}$ (e.g. eye template) according to the methods we discussed earlier. Thus, we built a generative model for

images given faces. Also priors within the class of faces can now be formulated, independent of the non–face objects, in terms of face generative parameters, e.g. specific individuums, illumination conditions, rotation angels, emotional expressions. Also, the distribution for relevant questions $p(q^a)$ for the low dimensional variable face/non–face is more easily adaptable and measurable than for very high-dimensional images. For example, one can measure or control the average number of people passing a camera (i.e. $p(q^a)$) much easier than to change or even approximate the probability of the corresponding images $p(y^a)$.

However, according to our convention denoting by $q$ the generalized questions to nature and not the input to the action producing device, we have to exchange the letters in the notation and write $q = y^a$ for the face variable and $y = q^a$ for the image. Then the face detector $p(\hat{q}|y, f)$ gives an approximate classification $\hat{q}$ for the variable $q$.

We study the situation, where an action device inverts an available generative model, now more formally and define an

**IM**: *inverse decision problem with memory* as a model with $p(z|q, y, \hat{f}) = p(z|y, q_c \hat{f})$, i.e. where the the action model produces answers $\hat{q}$ depending not on the index of a generalized question $q$ but on their answer $y$. In this situation we will use for $z$ the notation $\hat{q}$. Dependencies within components of $\hat{q}$ are allowed.

In the same way we define a

**I**: *inverse decision problem without memory* as a model with $p(\hat{q}|q, y, \hat{f}) = p(\hat{q}|y, \hat{f})$, with $z = \hat{q}$ i.e. where the the action model produces answers $\hat{q}$ dependent on $y$, however independent of its history. Dependencies within components of $\hat{q}$ are allowed.

Because it also seems quite natural to work with variables representing input and output of the action producing device we will call

**AR** : a problem to be formulated in *action representation*, if the action producing device is written in a form $p(z|q, y, \hat{f}) = p(\hat{y}^a|q^a, \hat{f}) = p(\hat{y}^a|q^a, y^a, \hat{f})$. This defines (action) output $\hat{y}^a = z$, requires the $i$–component of the (action) input $q_i^a$ to include at least that part of $q_i, y_i$ on which $z_i$ depends. The remaining variables are called $y^a = \{q, y\} \setminus q^a$. Analogously, we can define an action representation with memory to be of the form $p(z|q, y, \hat{f}) = p(\hat{y}^a|q^a, y_c^a \hat{f}) = p(\hat{y}^a|q^a, y^a, \hat{f})$.

**MR** : We will call our original formulation with $q$ representing questions and $y$ answers a *measurement representation*.

Table 2 summarizes the definitions, our convention used in the measurement notation (MR) and the action representation (AR).

Now we show that every inverse model looks 'pseudo parallel' in action representation without being necessarily equivalent to a parallel decision problem. The reason is that the roles of $q$ and $y$ are not freely exchangeable in

**Decision problems**

| $p(z|q,y,\hat{f})$ | parallel | inverse |
|---|---|---|
| without memory | $p(z|q,f)$ | $p(z|y,f)$ |
| MR: | $p(\hat{y}|q,\hat{f})$ | $p(\hat{q}|y,\hat{f})$ |
| AR: | $p(\hat{y}^a|q^a,\hat{f})$ | $p(\hat{y}^a|q^a,\hat{f})$ |
| with memory | $p(z|q,y_c,f)$ | $p(z|y,q_c,f)$ |
| MR: | $p(\hat{y}|q,y_c,\hat{f})$ | $p(\hat{q}|y,q_c,\hat{f})$ |
| AR: | $p(\hat{y}^a|q^a,y_c^a,\hat{f})$ | $p(\hat{y}^a|q^a,y_c^a,\hat{f})$ |

Table 2: Classification of decision problems

a decision model. We defined a model to be completely determined by $q$, $y$ and $f$, the $f$–independent part of variables $q$, i.e. $p(q|f) = p(q)$, and all $f$–dependent variables $y$. That means, the distribution of $q$ is always known, and can only depend on a state of nature over already observed data, i.e. $p(q|y_c,\hat{q}_c,f) = p(q|y_c,\hat{q}_c)$ independent of $f$. Then the decomposition $p(y|q,f) = \int df^0 p(y|q,f^0)p(f^0|f)$ does not need to include a factor $p(q|f^0)$, and the joint probability $p(q,y|\hat{q}_c f)$ factorizes into a $f$–dependent and a $f$–independent factor

$$p(q,y|\hat{q}_c f) = (q|y_c,\hat{q}_c)p(y|q,f). \qquad (16)$$

Accordingly, every problem in action representation can be written in at least a 'pseudo parallel' form. For the example of an inverse decision problem we can factorize the probability in the action representation

$$\begin{aligned} p(q,y|\hat{q}_c,f) &= p(y|q_c,\hat{q}_c,f)p(q|y,\hat{q}_c,f) \qquad (17) \\ &= p(q^a|y_c^a,\hat{y}_c^a,f)p(y^a|q^a,\hat{y}_c^a,f) \end{aligned}$$

We remind, that the first line has to be read as $p(y_1|f)p(q_1|y_1,f)p(y_2|y_1,q_1,\hat{q}_1,f)p(q_2|y_2,y_1,q_1,\hat{q}_1,f)\cdots$, where $p(q_2|y_2,y_1,q_1,\hat{q}_1,f)$ can not necessarily be simplified to $p(q_2|y_2,f)$ if $p(q|y_c,\hat{q}_c) \neq p(q)$, i.e. for a nonstatic model. The factors of this inverse picture are related to the original model by[39]

$$p(q|y,\hat{q}_c,f) = \frac{p(q|y_c,\hat{q}_c)p(y|q,f)}{p(y|q_c,\hat{q}_c,f)},$$

and

$$p(y|q_c,\hat{q}_c,f) = \int_c dq\, p(q|y_c,\hat{q}_c)p(y|q,f).$$

where the symbol $\int_c$ denotes a 'causal' or 'conditional' integral defined as $\int_c dq\, p(q) = \prod_i \int dq_i p(q_i|q_{c(i)})$ (not $\int \prod_i dq_i\ p(q)$ !) with $q_{c(i)} = \{q_j, j < i\}$, so that for every factor the $q_{c(i)}$–dependency remains.[40] Notice, that

---

[39]See Saul, Jaakola, & Jordan, 1996, Jaakola & Jordan, 1996, for a variational method for such inversions in sigmoid belief networks.

[40]Jordan, 1995, shows the interesting fact that Gaussians and the logistic function $\frac{1}{1+e^{-z}}$, which is often used as activation function in neural networks, are for binary classification related by such an inversion. For binary $y_i \in \{y_1, y_2\}$ and Gaussian $p(x|y_i)$ with equal variance (or covariance matrices) $p(y_i|x) = \frac{p(x|y_i)p(y_i)}{\sum_j p(x|y_j)p(y_j)} = \frac{1}{1+e^{-\ln\frac{p(x|y_i)}{p(x|y_{j\neq i})}-\ln\frac{p(y_i)}{p(y_{j\neq i})}}}$

this is not the inversion necessary to obtain $p(f^0|D)$ from $p(y|q,f^0)$. For a static model, i.e. for $p(q|y_c,\hat{q}_c) = p(q)$, the joint probability factorizes also in the inverse notation, $p(q_i)p(y_i|q_i,f) = p(y_i|f)p(q_i|y_i,f)$, with the individual factors $p(y_i|f)$ still being state dependent. If we try to exchange the role of $q$ and $y$ we find in Eq.(17) compared to Eq.(16):

1. History dependent questions. To interpret $y$ to a question in a measurement representation (and $f$ as a state) $f$ and $y$ must completely determine $p(q|y,\hat{q}_c,f) = p(q|y,f)$, like it is the case in the inverse of a static model. In other cases one can define a set of history (i.e. $\hat{q}_c$, $q_c$, $y_c$) dependent questions $y_i'$ (or states) labeled by $\{y_i, y_c(i), q_c(i), \hat{q}_c(i)\}$.

2. State dependent questions. According to our definition of questions $y$ can only be interpreted as question (for nature not for $\hat{f}$ !) if state independent, i.e. $p(y|q_c,\hat{q}_c,f) = p(y|q_c,\hat{q}_c)$. In other cases the distribution of relevant 'questions' $y$ in the inverse picture is state dependent, i.e. not exactly known, and they are therefore effectively part of the unknown variables, which we call answers. It is always possible to enlarge $F^0$ to include $q$ and write $p(q|F_q) = p(q|D^q) = \int df_q^0 p(f_q^0|D^q)p(q|f_q^0)$ depending on data $D^q$ which have been used to determine the probability distribution. Thus, the space $F^0 = F_q^0 \otimes F_{y|q}^0$ also contains hypotheses about different possible parameterizations $p(q|f^0)$. However as long as there is a 'factorial border' between $f_q^0$ and $f_{y|q}^0$, i.e. $p(f^0) = p(f_q^0)p(f_{y,q}^0)$ factorizes, the problems for $q$ and $y$ given $q$ (not y alone) are independent. As soon as they are dependent, we must perform the decomposition into pure $f_{y|q}^0$ states $q$–dependent, i.e. we must treat $q$ as answer. For example, when not only $p(y|q,f)$ but also $p(q)$ is sampled when obtaining data, this is relevant for determining $p(q)$, in any case when it is not yet completely known. In this case conceptually both $q$ and $y$ are part of the answers, let us say $(q,y) = y'$, and the set of $q'$ is reduced to one element with the meaning: 'Get next data pair'. However, when the hypothesis spaces factorize and the data do not induce new correlations, both problems can be treated separately, and $q$ can be used as question for the $y$ problem.

We finally notate the probability to suffer loss $l$ in the measurement representation

$$p(l|f,\hat{f}) = \int dq \int d\hat{q} \int dy\, p(q|y_c,\hat{q}_c)p(y|q,f)$$

$$\times p(\hat{q}|y,\hat{q}_c,\hat{f})\delta(l(q,\hat{q},y) - l)$$

with expected risk

$$r(f,\hat{f}) = \int dq \int d\hat{q} \int dy\, p(q|y_c,\hat{q}_c)p(y|q,f)p(\hat{q}|y,\hat{q}_c\hat{f})l(q,\hat{q},y).$$

and in action representation

$$p(l|f,\hat{f}) = \int dq^a \int dy^a \int d\hat{y}^a\, p(q^a|y_c^a,\hat{y}_c^a,f)p(y^a|q^a,\hat{y}_c^a,f)$$

---

gives a logistic function because for equal variance the term quadratic in $x$ cancels and only the difference of the first moments remain. This holds also if the $p(x|y_i)$ belong to the same exponential family with possibly different first, however equal higher moments.

39

$$\times p(\hat{y}^a | q^a, \hat{y}_c^a, \hat{f}) \delta(l(q^a, \hat{y}^a, y^a) - l).$$

with expected risk

$$r(f, \hat{f}) = \int dy^a \int d\hat{y}^a \int dq^a \, p(q^a | y_c^a, \hat{y}_c^a, f) p(y^a | q^a, \hat{y}^a, f)$$

$$\times p(\hat{y}^a | q^a, \hat{y}_c^a, \hat{f}) l(q^a, y^a, \hat{y}^a).$$

We conclude: An inverse decision problem cannot necessarily be formulated as a parallel decision problem, because its action representation yields in general state dependent $q^a$.

We will see later that for certain static decision problems a one step maximum posterior approximation is sufficient only for approximation problems, i.e. when $p(y|q, f)$ itself shall be approximated. Then this step has numerically the form of an empirical risk minimization. Obviously, approximating the inverse probability $p(q|y, f)$ is not an approximation problem for $p(y|q, f)$, (except when $p(y|q, f) = p(q|y, f)$) and we have just seen that it cannot be transformed into one. Thus, for inverse problems a maximum posterior approximation should in general be completed by a second step. This two step procedure has for inverse problems the following form: In a first step $p(y|q, f)$ is approximated by $p(\hat{y}|q, \hat{f}_{app})$. In the second step the optimal approximating action device $\hat{f}_{app}$ is identified with a (fixed) state of knowledge $f^* = \hat{f}_{app}$, and accordingly $\hat{y}$ with $y$. Then $p(\hat{q}|y, \hat{f})$ is chosen to minimize the loss for the approximating device, i.e. for given state $f^*$. However, for a full Bayesian treatment, and for the justification of an empirical risk minimization the difference between parallel and inverse problems is in no way conceptually important.

## 6.4   Algorithms

Often the process of calculating the answer to $q_r$, that is the optimal learning or decision algorithm, is to difficult to be actually performed. Thus, one has to use a simplified *decision or learning algorithm*, i.e. another question $\hat{q}_r = \hat{a}$, to produce an answer $\hat{f} = \hat{a}(f(D))$. We discuss the following for parallel decision problems, but it applies analogously to inverse and general fair decision problems. If we have to decide between several available algorithms this corresponds to another (higher level) decision problem[41], with the $\hat{y}$ producing device with $p(\hat{y}|q, \hat{f})$ replaced by a $\hat{f}$-producing device $\hat{a}$ with $p(\hat{f}|q, y_c, \hat{y}_c, \hat{a})$, and we have to define the $\hat{a}$ dependency of the loss which could be formulated by extending the set of relevant questions $Q^l \rightarrow Q^{l,a}$ to include loss relevant aspects. If no algorithm specific aspects have to

be included, only an additional $\hat{f}$-integration has to be performed. The same is essentially true if the algorithm specific loss can be represented by a $q$, $y$, $\hat{f}$-independent constant. On the other hand, loss related aspects which depend on more than one of the original relevant questions $q^l$ also depend on the correlations between $q^l$ (for which the original loss, and thus also the risk, is insensitive) and can generate a much more complicated set of relevant questions and associated loss $l$ for the new problem.[42] This expression corresponds to the expression $p(\hat{f}^0 | q, y_c, \hat{y}_c, \hat{D})$ we discussed for parallel decision problems. Thus, from this point of view the $\hat{f}$ corresponds to possible pure action states $\hat{f}^0$, and the algorithm $\hat{a}$ to the available data $\hat{D}$. Interestingly, one notation only refers to data, and assumes the algorithm to be implicit, while the other refers to the algorithm and assumes the data to be implicit. Thus, we have a parallel decision problem with memory where the $\hat{f}$ are part of the internal noise variables $\hat{y}_{\hat{a}} = (\hat{y}, \hat{f})$. The $q^l$ are related to a loss function $l$ which can depend on $\hat{a}$ and $q_r(\hat{a})$ chooses the $\hat{a}$ which minimizes a functional $r$ of the loss distribution

$$p(l|\hat{a}f) = \int d\hat{f} \int dq \int dy \int d\hat{y} \, p(q|y_c, \hat{y}_c, \hat{f}_c) p(y|q, f)$$

$$\times p(\hat{y}|q, \hat{f}) p(\hat{f}|q, y_c, \hat{y}_c, \hat{a}, f) \delta(l(q, y, \hat{y}, \hat{f}, \hat{a}) - l).$$

A loss $l(q, y, \hat{y}, \hat{f})$ can be chosen $\hat{f}$-independent by including additional variables $\hat{y} = \hat{y}(\hat{f})$ if necessary. To include algorithmic specific aspects without using an explicit algorithm dependent loss function we can, for example, include $\hat{a}$-dependent internal variables $\hat{y}_{\hat{a}}$ in the loss function produced with $p(\hat{y}_{\hat{a}}|q, \hat{a})$ measuring for algorithm specific variables like their requirement of resources like calculation time, memory and other aspects.

An algorithm is defined by its $p(\hat{f}|q, y_c, \hat{y}_c, f, \hat{a})$. It produces an answer $\hat{f}$ which should at least be depending on its state of knowledge $f$, that is its prior probabilities and training data. We always assume its knowledge of the $\hat{f}$, i.e. of the $p(\hat{y}|q, \hat{f})$. The dependence of $p(\hat{f}|q, y_c, \hat{y}_c, f, \hat{a})$ from $y_c$ and $\hat{y}_c$ allows the algorithm to adapt, i.e. learn, during the test phase. In that case the loss evaluates the learning curve of algorithms. One can allow the choice of $\hat{f}$ to be dependent from the test question $q$ being equivalent to enlarging the space of available $\hat{f}$. The variable $\hat{f}$ can be a vector, for example, if a decision is required after presenting part of $q$. In situations with $p(q|y_c, \hat{y}_c, \hat{f}_c) = p(q|y_c, \hat{y}_c)$ and $l(q, y, \hat{y}, \hat{f}) = l(q, y, \hat{y})$ algorithms can be seen as improved $\hat{f}$-dependent $\hat{y}$-producing devices with

$$p(\hat{y}|q, y_c, f, \hat{a}) = \int d\hat{f} p(\hat{y}|q, \hat{f}) p(\hat{f}|q, y_c, \hat{y}_c, f, \hat{a}).$$

Choosing for real loss $r$ as the expectation functional we are looking for the $\hat{a}$ which minimizes

$$r(f, \hat{a}) = r[p(\cdot | f, \hat{a})] = \int dl \, l \, p(l | f, \hat{a})$$

---

[41] This problem of finding an optimal algorithm can be a much more complicated problem than finding the optimal decision. So comparison of algorithms can be done in only a few number of (simple enough) cases (See for example, Watkin, Rau, & Biehl, 1993 and references therein for a review). When also using approximations for the high level problem all the same problems appear on the higher level again. But if we assume the decision problems of the different levels to be similar we could at least check for consistency: Does approximation $A$ on level i produce approximation $A$ on level $i - 1$.

[42] Compare the distinction between $l$ and $L$ in (Haussler, 1995).

$$= \int d\hat{f} \int dq \int dy \int d\hat{y}\, p(q|y_c, \hat{y}_c, \hat{f}_c) p(y|q, f) p(\hat{y}|q, \hat{f})$$

$$\times\, p(\hat{f}|q, y_c, \hat{y}_c, f, \hat{a}) l(q, y, \hat{y}, \hat{f}, \hat{a}).$$

## 6.5 Minimization

Solving decision problems requires minimization algorithms. In general, also the minimization problem can be seen as a learning problem: Given data of function values and information about the function (e.g. differentiable, symmetric) give the position of the minimum. In so far our discussion of learning and prior information also applies to optimization. Usually, minimization algorithms perform active queries for local data (for example, taking a new data point in direction of the gradient) to improve a given guess for the location of the minimum and proceed iteratively until a certain convergence criterion is fulfilled. An iteration procedure can be written

$$\hat{f}^{i+1} = G(i, D_i, \hat{f}^i) = G_i(\hat{f}^i),$$

where $\hat{f}^i$ denotes the current guess for the location of the minimum at step $i$ and $D_i$ (which for the sake of simplicity will be skipped from now on from the notation) the new and accumulated previous data ($D_i = (r(f, \hat{f}^j), \hat{f}^j)\}_{j \leq i}$). Besides data points $G$ depends on prior knowledge about the function (in our case $r(f, \hat{f})$). The past iterations provide data for the minimization problem, and $D_i$ indicates that $G$ changes with the amount of available data, i.e. the number of iterations. Thus, one can say $G$ is trained on the available data. This at least implicitly assumes prior knowledge which allows available data to carry information about other parts of the function.

In general, we can allow any reparameterizations $T$ (of the $\hat{f}$, in our case) as long as $T$ is *locally injective at the locations of global minima* $\hat{f}^*$, i.e. $T^{-1}(T(\hat{f}^*)) = \hat{f}^*$. Reparameterizations can be not injective for points which are no minima, as those can be excluded from further search. But also under *globally bijective reparameterizations* minimization problems can look quite different for the transformed variables. Reparameterizations can be nonlinear, differentiable or non–differentiable, or linear transformations for vectors or functions (like $\hat{f}$) i.e. a change of the representing basis if $\hat{F}$ is a vector or Hilbert space, or even a random permutation of the function values. Transformations can create arbitrary neighborhood relations, so the minimization problem can become trivial, like for example when the function values are ordered monotonically (to find such a permutation, however, the minimum problem has to be solved), or arbitrary hard (non–smooth, random).

Sometimes, it is technically helpful to 'linearize' problems, by giving every degree of freedom its own linear dimension. If we define $R$ to be the space of all possible functions $r(f, \hat{f})$ (for fixed $f$) and give the values of $r$ a linear structure, we can expand any $r$, at least formally, into basis functions $r = \sum_i a_i b_i(\hat{f})$. Thus any function $r'(\hat{f}) = \sum_i a_i' b_i(\hat{f})$ on $\hat{F}$ taking values in the linear range of $r(f, \hat{f})$, for example $r'(\hat{f}) = r(T(\hat{f}))$, is

by construction of $R$ a linear transformation $A(\hat{f}', \hat{f})$ of $r(\hat{f})$. Then $r' = Ar$, which is defined by the mapping of the coefficient vectors with $a' = Aa$.[43] The resulting dimensionality may however very soon be intractable huge, e.g. for infinite $\hat{F}$ the resulting space has infinite dimension. To be able to use such a space for calculation, it must be restricted. In Hilbert spaces for example only functions are allowed for which the expansion in a basis can at least be arbitrarily well approximated in some norm by an (arbitrary large, but) finite number of basis functions.

Reparameterizations always change only the arguments not the function (risk) values itself. There is also the possibility to change the function values without changing the location of the minimum. In general, for a function $f$ (in our case the risk, i.e. $f$ does not denote a state here) the positions $x^*$ $(=\hat{f}^*)$ of global minima $f(x^*)$ $(=r(f, \hat{f}^*))$, defined by $y^* = f(x^*) \leq y = f(x)$, $\forall x$, do not change under transformations $h$ which obey $y^* \leq y \Leftrightarrow h(y^*) \leq h(y), \forall y$. We will call such transformations $h$ *strictly monotonically increasing relative to* $y^*$. (Analogously, we define strictly monotonically decreasing relative to $y^*$ $y^* \leq y \Leftrightarrow h(y^*) \geq h(y), \forall y$.) As we do not know $y^* = f(x^*)$ in advance we have to require this relative to all values which are possible candidates for a global minima.

Minimization methods can only return a minimum within some selected finite (sub)set of considered function values. Local methods (see below) for example find local minima. Locations of minima are invariant for all subsets under *strictly monotonically increasing transformations* $h$ defined by $y < y' \Rightarrow h(y) < h(y'), \forall y, y'$, equivalent to strictly monotonically increasing relative to all $y$, i.e. $y \leq y' \Leftrightarrow h(y) \leq h(y'), \forall y, y'$. Analogously defined strictly monotonically decreasing transformations change a minimum to a maximum.

There is a large variety of methods and concepts available for minimization, with many possibilities of combinations. (For a discussion with respect to neural networks see for example Golden, 1996. In Section 9 optimization methods are needed for maximizing the posterior probability.) The fact that the following principles of minimization algorithms can be applied to transformed (strictly monotonically increasing relative to global minima) and reparameterized (bijective at locations of global minima) problems, makes clear, on one hand, that in general a large number of possibilities can exist to attack a specific problem. On the other hand it shows that also the optimization process depends on prior information.

The simplest method, which does not refer to any dependencies of function values, is

1. an *unadapted search* (stochastic, predetermined de-

---

[43]See for example generalized additive models, where interaction terms are added (Hastie & Tibshirani, 1990), the comments in Minski-Papert-1990 (the new edition of the 1969 book) about the general applicability of (linear) perceptrons, the support vector machine, where linear relations in the feature space correspond to nonlinear relations in the input space (Vapnik, 1995), and for a general approach, Smola & Schölkopf, 1997.

terministic, or exhaustive), i.e. with data $D_i$ collected independent of function (risk) values for $\hat{f}$. One iteration step consists in sampling of a new data point, and $G$ compares this point with the current guess $\hat{f}^i$. The sampling distribution does not depend on previous iteration steps.

Deterministic algorithms (deterministic $G$) are special cases of stochastic algorithms (probabilistic $G$). Prior information enters $G$ if the sampling distribution at step $i$ depends on $D_i$. The most common case are local dependencies yielding

2. *local iterative methods.* For smooth functions small values are in the neighborhood of other small values and therefore the search (sampling probability) is concentrated near the current guess of the minima. As differentiable functions have a gradient equal to zero at a minimum this allows to search for those zeros which are possible locations of minima. Sometimes stationary points can be found analytically, but usually nonlinear equations require iterative solutions. (The term 'analytical' is commonly used for solutions where the iteration can be done easily, like determining a certain square root or the numerical evaluation of constants like $e$ or $\pi$.) Technically, the iteration is often implemented in the form of a relaxation method (see Section 9), which includes common algorithms like those based on the gradient, and its stochastic (like on–line learning), restricted (e.g. line search) variants, and may include higher order derivatives (as in Newton methods) (Pierre, 1986, Bazaraa, Sherali, & Shetty, 1993, Bertsekas, 1995). For these algorithms $G$ depends on the value of the function and some of its derivatives at the current location. In a discretized implementation of derivatives, (and similar in a simplex search method) $G$ depends on a set $D_i$ of function values $r(\hat{f})$ which can be called a local population. Those methods find local minima. A nonlocal aspect can be added simply by combining them with unadapted search methods, usually implemented by comparing different local minima.

There can be nonlocal dependencies between function values, which makes it desirable to have $G$ being dependent not only on a local but also on a nonlocal population of function values, leading to

3. *parallel or nonlocal iterative methods.* Here $G$ depends on a nonlocal population $D_i$. like in genetic algorithms, (See Holland, 1975, Goldberg, 1989, Davis 1991, Michalewicz, 1992, Schwefel, 1995, Mitchell, 1996) This allows nonlocal interactions between a possibly large number of function values. The dynamic of the population $D_i \rightarrow D_{i+1}$ corresponds to an iteration for a population vector.

In general, we only have to require that a fixed point of the iteration corresponds to a solution of the minimization problem. This allows transformations of the problem during iteration:

4. *Transformation methods* use transformations of the problem (i.e. of $G_i$), starting with an easy solvable, e.g. one–minimum problem, and slowly transforming to the problem of interest. They are called homotopy or continuation methods if they approximate a smooth family of transformations (Allgower, 1990, Richter & DeCarlo 1983, or, for an application in scattering theory, Giraud & Nagarajan, 1991, Wierling et al. 1994). Parameter, like the step width in gradient algorithms, mutation rate in genetic algorithms, or the temperature in simulated annealing (See Ripley, 1987, Davis, 1987, Aarts & Korts, 1989) can also be seen as such deformation parameters.

An example of transformations which correspond to a strictly monotonic transformation at fixed points are transformations with $h(y) \leq h(y^*) \Rightarrow y \leq y^*, \forall y$, which we will call *minimality sufficent relative to $y^*$*, if $y^*$ is updated (at least from time to time) during iteration. (Accordingly $h(y) \geq h(y^*) \Rightarrow y \geq y^*, \forall y$ will be called maximality sufficient relative to $y^*$.) For example, adding a function $\tilde{r}(\hat{f})$ to the risk $r(f,\hat{f})$ (e.g. a (quasi) distance $D(\hat{f}^i, \hat{f}^{i+1})$) with a minimum at the current guess ($\hat{f}^i$) ensures, that decreasing both terms ($r+\tilde{r}$) also decreases the original function (r). This is, for example, used in the EM–*like algorithms* (Dempster, Laird, Rubin, 1979, Tanner, 1993, Gelman, Carlin, Stern, Rubin, 1995, for an information geometrical interpretation and the related (most times identical) *em* algorithm see Amari, 1985, 1995).

Non–exhaustive search is always restricted to a subspace. Restricting to a subspace a priori is also called

5. a *variational method.* Here the function is parameterized and only a part of the parameters are used for minimization. In a linear variational method the stationarity condition is, for example, expanded into a linear basis of a Hilbert space of functions, and solved in a linear subspace. Examples include the methods of finite elements. Variational methods, including nonlinear ones, are also often used in physics, especially in quantum mechanics. There, for example, finding a bound for the ground state energy of quantum mechanical systems smaller than some instability causing threshold can have drastic consequences. In this context a product ansatz for functions in several variables is also called mean field approach. Variational methods have recently also be applied to general graphical models. (See for example Saul, Jaakola, & Jordan, 1996, Jaakola & Jordan, 1996).

We conclude with the interesting observation, that (non–exhaustive) minimization requires knowledge about nonlocal dependencies also for the risk functional $r(f,\hat{f})$. This suggests the principal possibility of reducing for practically solvable problems the risk functional to independent values, whose number in practice must be finite, corresponding to a finite $\hat{F}$ and finite effective $F^0$.

# 7 Bayesian approach

## 7.1 Inserting model states

The following constituents of a decision problem are assumed to be known: 1.) the action producing device $p(\hat{y}|q, \hat{f})$, 2.) the definition of the test distribution $p(q|y_c, \hat{y}_c)$,[44] 3.) the loss function $l(q, y, \hat{y})$. For evaluating a risk functional $r$ the main problem remains determining the answer probabilities $p(y|q, f^0)$ for the test questions of the realized pure state (of nature). There exists a Bayesian and a Frequentist approach for this problem.

In the *Bayesian* approach model states $f^0$ are inserted as hidden variables. This is the concept we used in this paper defining a state of knowledge $f$ with probabilities $p(f^0|f)$ and is called in the context of decision problems Bayesian decision theory (Berger, 1985). The answer characteristics $p(y|q, f^0)$ of the possible pure states have to be known.[45] According to the Bayesian paradigm the (training) data dependence $p(f^0|f(D)) = p(f^0|D)$ of the state of knowledge $f = f(D)$ can be written

$$p(f^0|D) = \frac{p(D|f^0)p(f^0)}{p(D)} = \frac{p(y^D|q^D, f^0)p(f^0)}{p(y^D|q^D)},$$

with $q^D$, $y^D$ the vectors of questions and corresponding answers in the data. A Bayesian expected risk reads

$$r(f, \hat{f}) = \int df^0 \, p(f^0|D) r(f^0, \hat{f})$$

with

$$r(f^0, \hat{f}) = \int dq \int dy \int d\hat{y} \, p(q|y_c, \hat{y}_c)$$
$$\times p(y|q, f^0) p(\hat{y}|q, \hat{f}) l(q, y, \hat{y}),$$

or for an inverse setting

$$r(f^0, \hat{f}) = \int dq \int d\hat{q} \int dy \, p(q|y_c, \hat{q}_c)$$
$$\times p(y|q, f^0) p(\hat{q}|y, \hat{f}) l(q, \hat{q}, y).$$

The posterior probability $p(f^0|D)$ is the only data dependent term. Introduction of model states $f^0$ makes the treatment of the training data independent of the test set. Both cases, test questions not included in the training data and training data not included in the test data are no conceptual problems.

When $p(y|x, f^0)$ is specified by $p(x|y, f^0)$ and an $f^0$–specific prior $p(y|f^0)$ according to

$$p(y|q, f^0) = \frac{p(q|y, f^0) p(y|f^0)}{p(q|f^0)}$$

one has to calculate the probability of the data under $f^0$

$$p(q^D|f^0) = \sum_{y^D} p(q^D|y^D, f^0) p(y^D|f^0),$$

(see for example Saul, Jaakola, & Jordan, 1996) to get $p(f^0|D)$ and integrate over different states $f^0$.

---

[44]If not under direct control, $p(q)$ can be determined separately, or the $q$ can be included in the set of $y$.

[45]Under the assumption of the chosen model the Bayesian approach is (defined as being) optimal, but in practice the method depends of course on the correctness of the model for the situation in mind.

## 7.2 Maximum posterior approximation

Practical calculations of a Bayesian risk are, if not analytically solvable, in general only possible for a restricted set of $f^0$ and $\hat{f}$. Numerical methods using Monte Carlo integration techniques[46] making the full integration in some cases feasible, are used in the area of neural networks for example by the Boltzmann machine (Hinton, & Sejnowski, 1983, 1986; Ackely, Hinton, & Sejnowski, 1985) and have been applied to Bayesian calculations (See Gelfand & Smith,1990, Gelfand, Hills, Racine–Poon, Smith, 1990, Geyer, 1992, Besag & Green, 1993, Smith & Roberts, 1993, Tierney, 1994, or Gelman, Carlin, Stern, & Rubin, 1995) including Bayesian analysis of neural networks (Neal 1993, 1996).

Those methods perform an plug–in estimate of the expected risk. The difference to the Frequentist method discussed later is that the test data are generated according to the a posteriori distribution. One also has to proof that the Monte Carlo estimate converges without having calculated the exact solution. That means using a finite sample we have to assume or calculate nonlocal knowledge about the risk function. In practice we may check the accuracy by repeating the calculation and estimating the variance of the optimal $\hat{f}$ using for example cross–validation or the bootstrap. Those are methods of classical statistics and a plug–in estimate is technically an empirical risk minimization (see Section 8) with $f$–*generated virtual examples* sampled according to $p(y|q, f)$.

The problems related to the plug–in principle will be discussed below for the Frequentist approach.

Alternatively to Monte Carlo methods, the method of Laplace can be used to approximate the risk integral. This is the real version of the saddle point approximation or method of steepest descent for complex functions (see for example: De Bruijn, 1981; Bleistein, Handelsman, 1986) Taylor expansion of $h(x)$ around its maximal value $x^*$ and performing the resulting integrals gives for a real one–dimensional function $h$

$$\int_{-\infty}^{\infty} df^0 \, e^{-\beta E(f^0)} = \left( \frac{2\pi}{\beta E^{(2)}(f^{0,*})} \right)^{\frac{1}{2}} e^{-\beta E(f^{0,*})}$$

$$\times e^{+\frac{1}{\beta}\left( \frac{5(E^{(3)}(f^{0,*}))^2}{24(E^{(2)}(f^{0,*}))^3} - \frac{E^{(4)}(f^{0,*})}{8(E^{(2)}(f^{0,*}))^2} \right) + \mathcal{O}(\frac{1}{\beta^2})}, \qquad (18)$$

written for a function with one minimum $E(f^{0,*})$ (or maximum for $-E$), $E^{(2)} > 0$, $E^{(i)}$ denoting the $i$th derivative at $f^{0,*}$. Interpreting $1/\beta$ as 'temperature', this expansion in $1/\beta$ is a 'low temperature' approximation. In the multidimensional case $E^{(2)}$ in the square

---

[46]Invented in statistical physics and going back to Metropolis, Rosenbluth, Rosenbluth, Teller, Teller (1953) (Metropolis algorithm) and Alder & Wainwright (1959) (molecular dynamics). For applications and developments in physics see for example Hammersley & Handscomb (1964), Binder (1986, 1987, 1995), Binder & Heermann (1988), or the last chapter in Montvay & Münster (1994), for mathematical background on Markov chains, for example, Seneta, 1981, and for their early use in statistics see Hastings, 1970, Ripley, 1977, Geman & Geman, 1984, Ripley, 1987.

root factor for example has to be replaced by the determinant of the matrix of second derivatives. For positive quadratic $h$, i.e. Gaussian $e^{-\beta E}$, all terms $E^{(i)}$ with $i > 2$ vanish. Higher order terms are obtained via Wick's theorem[47] and include, in a graphical notation, an exponential of all so called 'linked diagrams'. [48] (see for example Negele & Orland, 1988, Itzkyson & Drouffe, 1989).

Generalizations of the saddle point formula for multiple extrema exist which care about overlapping parts belonging to sufficiently close extrema (Berry, 1966, Miller, 1970, and Connor & Marcus, 1971). However, one usually assumes the extrema to be well separated, and in our case, being not so much interested in the actual value of the Bayesian risk than in finding the best $\hat{f}$, only effects varying between different $\hat{f}$ would be important.

To find the expansion point $f^{0,*}$ one has to solve the stationarity conditions

$$\frac{d}{df^0} E(f^0) = 0, \qquad (19)$$

which are nonlinear for non–quadratic $E$, and have then to be solved iteratively to find a self–consistent solution. In a multidimensional case a self–consistent solution $f^{0,*} = \{f^{0,*}_x | x \in X\}$ is only influenced by the values of $f^{0,*}_{x'}$ but no other $f^0 \in F^0$. Thus, $f^{0,*}$ has to incorporate approximately the combined effects of all other $f^0 \in F^0$, so sometimes $f^{0,*}$ is also called a mean field solution. Analogously one sometimes refers to the stationarity condition (19) as mean-field equation and to the saddle point approximation as *mean field approach*.

According to Eq.(18) one might apply the saddle point approximation to the whole $f^0$–dependent integrand $p(f^0|D)r(f^0,\hat{f})$ or $y-$ and $q$–dependent to $p(f^0|D)p(y|q,f^0)$. Both variants are clearly usually too complicated, leading for example to a $\hat{f}$–dependent factor depending on the second derivative.[49] But having a large amount of data or strong nonlocal dependencies it is often reasonable to assume the posterior probability $p(f^0|D)$ to be peaked sharply around one maximum. More formally, we may identify $\beta$ with the number $n$ of training data, and assume that the sample mean

[47]Which is a systematic way to calculate (multidimensional) Gaussian integrals over polynomials. Those arise when expanding the remaining exponential factor.

[48]The difference between 'linked' and 'unlinked' diagrams is similar to those between moments and cumulants generated by $< e^L >$ or $\ln < e^L >$, respectively, in a high temperature expansion (See Section 5.3.2). The general relation between expanding a sum of exponentials and its logarithm is also known under the name 'linked cluster theorem'.

[49]On the other hand the $\hat{f}$–dependency is an interesting feature if $r(f^0,\hat{f})$ can be included , because then the saddle point approximation can be adapted to $\hat{f}$. It would lead to coupled maximization and minimization problems. (In contrast we will discuss below a maximization problem which is independent of the subsequent minimization problem.) This dependency of the maximization problem on the corresponding minimization problem, i.e. its adaption to $\hat{f}$, suggests for example an iterative procedure where both steps are performed alternately.

$\frac{1}{n} \sum_i \ln p(y_i|q_i,f^0)$ of the variable $z(y,q) = \ln p(y|q,f^0)$ becomes for large $n$ a nearly $n$–independent function. In this case a large $n$ (many data) correspond to a low temperature $1/\beta$.[50] Then the second order Taylor expansion of the log-posterior $\ln p(y|q,f^0)$ (i.e. a Gaussian approximation for the posterior) can be a good approximation. For example, under some regularity conditions, (e.g. the number of parameters included in $f^0$ is not chosen to increase with the sample size, the limit is not at the edge of the parameter space $F^0$, or one uses a model specification where different parameter values $f^0$ correspond to identical probabilities $p(y|q,f^0)$ at the maximum) this will be the case for i.i.d. random variables $z(y,q)$ with finite variance, according to the general asymptotic Gaussian limit theorem (Le Cam, 1953, 1986, Le Cam & Yang, 1990, see also the discussion and references in Chapter 4 and Appendix B of Gelman, Carlin, Stern, Rubin, 1995). Then the posterior will have a variance $(nJ(f^0))^{-1}$ where $J$ is the expectation under the true state of (the matrix of) the second derivatives of the log-posterior (Fisher information). For dependent variables this is not necessarily true, but can also be the case. Then, if $r(f^0,\hat{f})$ varies only weakly with $f^0$ compared with $p(f^0|D)$ (Gaussian $p(f^0|D)$ alone is not enough in this case) it does not strongly influence the location of the maximum. Then we can identify $\beta$ with $n$, and $E(f^0)$ with $-\frac{1}{n}\left(\sum_i \ln p(y_i|q_i,f^0) + \ln p(f^0)\right)$. (The factor $n$ disappears in the stationarity conditions for $h(x)$ if those are multiplied by $n$.) We do however not restrict to cases where we interpret $\beta = n$. Having nonlocal prior terms (interactions) in the exponent the saddle point approximation can be a good approximation even for a small $n$ of local data, if the dependencies induced by the prior (interaction) terms restrict the number of function with high probability strongly enough. Indeed, from physics it is known that mean field theories (saddle point approximations) can become exact when the correlations are strong, like for long–range forces or for local forces in high dimensional spaces. (For many physical models with local interactions, like the Ising model $d > 4$ is the dimension above which the mean field theory is valid.) In the case where only part of the exponent is multiplied by $\beta$, i.e. the integrand has the form $e^{\beta h(x) + \ln g(x)}$, we can apply the slightly more general formula

$$\int_{-\infty}^{\infty} df^0 \, r(f^0) e^{-\beta E(f^0)} \approx r(f^{0,*}) \left( \frac{2\pi}{\beta E^{(2)}(f^{0,*})} \right)^{\frac{1}{2}} e^{-\beta E(f^{0,*})}.$$

with $r(f^0)$ corresponding to $r(f^0,\hat{f})$, $-\beta E(f^0)$ to $L(f^0,\hat{f})$, and $f^{0,*}$ is the location of the minimum of $E(f^0)$ (maximum of $L(f^0,\hat{f})$), i.e. independent of $r(f^0)$.

This is called *maximum posterior approximation* (MaP). Especially for high dimensional spaces $F^0$ the

[50]Similarly, in field theories (e.g. quantum theory) in a Euclidean (imaginary time) formulation the system size is related to $n$, while in the corresponding interpretation as classical statistical system the parameter is an inverse temperatur $\beta$, and the evolution operator for imaginary times appears as "transfer matrix". In particular, the large $\beta$ limit corresponds to the limit of large system size. (See for example Zinn-Justin, 1989.)

off–peak contributions can be large requiring a large amount of data to allow a MaP approximation. Note also, that in the context of decision theory we only need to evaluate the Bayesian risk to select a optimal $\hat{f}$ and we do not have to require a good approximation of the risk itself, and $\hat{f}$–independent factors are therefore not important, like the ones including the second derivative $h^{(2)}$ or the factor $p(y^D|q^D)$.

Hence, to apply a $q-$, $y-$independent MaP approximation with respect to $f^0$ we write the $f^0$–dependent but $q-$, $y-$independent probabilities in exponential form

$$p(y_i^D|q_i^D, f^0) = e^{L^D(y_i^D|q_i^D, f^0)},$$

and

$$p(f^0|f) = e^{L^0(f^0)},$$

defining the log-probabilities $L^D$ (log-likelihood), $L^0$ (log-prior) and get

$$r(f, \hat{f}) \propto \int df^0 \, e^{\left(\sum_i L^D(y_i^D|q_i^D, f^0) + L^0(f^0)\right)} r(f^0, \hat{f}).$$

If we include the log-prior $\ln p(f^0)$ into $h(x)$, it is also included in the determination of the maximum $h(x^*)$. This allows to discuss situations where the $n$ data are not enough to yield sharp, especially non–degenerate, maxima (and indeed, choosing the relevant questions as basis questions, $X = Q^l$, the prior is always essential for local data if $n < |X|$). In contrast we assumed $r(f^0, \hat{f})$ not to be important for the location of the maximum. An extreme example for the opposite would be if the risk for the most probable state $f^{0,*}$ is for all $\hat{f}$ infinite, meaning that $f^{0,*}$ can be excluded from $F^0$. Less extreme, a risk can be strongly peaked at an $f^0$ with low posterior for some $\hat{f}$. On the other hand, it only matters which $\hat{f}$ is finally selected as the optimal one. This implies robustness against all errors, made somewhere on the way, which do not change this final decision.

We summarize that the MaP approximation includes probability aspects, but not aspects of relevance related to the loss function.

Maxima may be degenerated or weakly peaked within a subspace of parameters of $F^0$. Then one may perform a partial saddle point approximation for the subspace where the necessary conditions are fulfilled. For example, the risk may measure aspects of (i.e. depend on parameters describing) $f^0$ which are not measured by the data, so $p(f^0|D)$ does not dependent on them. Then the locations of the maxima depend on the risk $r(f^0, \hat{f})$ and the maxima are necessarily degenerated in direction of those relevant but not measured parameters (e.g. in a model with $X = Q^l$, uniform prior, and only local data). Then the MaP step returns not a single point but a subspace of important possibilities. More general, the maximum can, after incorporating data and prior, still be rather flat in some of the relevant dimensions. Then a MaP approximation should only replace a part of the $f^0$–integration while another part, i.e an integration over a subspace of parameters of $f^0$, remains. Performing this integration gives a new effective risk sensitive to the available data.

The MaP approximation consists in finding the most probable state $f^{0,*}$ to approximately calculate the $f^0$ integral. Then the factor $\sum_i L^D(y_i^D|q_i^D, f^{0,*}) + L^0(f^{0,*})$, being independent of $\hat{f}$, can be skipped and the optimal action state $\hat{f}^*$ can be found by minimizing $r(f^{0,*}, \hat{f})$. Often $\min_{\hat{f}} r(f^0, \hat{f})$ is a constant over all $f^0$, like in the usual regression case with mean square error, deterministic unrestricted $\hat{f}$, and Gaussian $f^0$ with $f^0$–independent variance. Thus, the full approximation procedure consists of two steps (MaP–MiR)

1. Maximization of the posterior (MaP):

$$f^{0,*} = \text{argmax}_{f^0} \left( \sum_i^N L^D(y_i^D|, q_i^D, f^0) + L^0(f^0) \right),$$

2. Minimization of risk (MiR):
   (for the state $f^{0,*}$ with maximal posterior)

$$\hat{f}^* = \text{argmin}_{\hat{f}} r(f^{0,*}, \hat{f}).$$

Note that in this approximation the MaP step is performed independent of the aspects important for the MiR step.[51] The first step can be interpreted as finding an approximation independent from its application. The second step uses the best found approximation for a specific application situation defined by the loss function. This is the usual implicit setting when looking for approximations without specifying applications for which they will be used. It has the advantage that the independence allows the same approximation to be used for several different applications. Thus, in the every day use of this procedure a statistician performs the first approximation step and potential users of the approximation the second, adapted to their problem. An example for a MaP–MiR related algorithm can be found in (Lemm, Beiu, Taylor, 1995). There the MaP step is implemented as a density approximation. A subsequent constructive algorithm tries to find a solution easy to implement in hardware, not working directly with the data but with the results of the density estimation of the first step. This is an attempt to minimize also aspects of the loss not related to approximation and corresponds to the MiR step.

The loss function depends on the action state $\hat{f}$ which may include aspects like complexity of $\hat{f}$. But note that the loss measures no aspects like complexity of the algorithm used to find the optimal (or a good) $\hat{f}$. So to say, action loss is included but no algorithmic loss. As a two-step procedure can often be expected to be more complex than a one–step procedure, the MaP–MiR procedure seems to be more appropriate for situations where

---

[51] In statistical practice or biological reality where on–line learning is required (and the model spaces $F^0$ and $\hat{F}$ may be adapted) both steps can of course be performed interlaced. See for example the Helmholtz machine (Dayan, Hinton, Neal, & Zemel, 1995; Hinton, Dayan, Frey, Neal, 1995), where in the 'learning phase' (MaP step) a 'generative model' (state $f^0$) is adapted and in the 'dreaming' phase (MiR step) the 'recognition model' (action $\hat{f}$) is optimized for given generative model $f^{(0,*)}$.

loss related with $\hat{f}$, for example its approximation ability and complexity, are more important than aspects of the loss related to the requirement of resources of the algorithm. But as stated in Section 6 one can also consider algorithms as part of $\hat{f}$ in a higher level problem. Then one can include specific algorithmic aspects of the loss and look for an optimal algorithm for a certain distribution of application (learning) situations. Again, to solve for the best algorithm one has to use a meta–algorithm and the same kind of problem appears on this higher level as here meta–algorithmic loss aspects are not included. Going further, meta–algorithms could be included into $\hat{f}$ using meta–meta–algorithms and so on, but usually the complexity from one level to the next increases so much that such applications are not expected to be feasible in most practical cases.

There are cases where $p(y|q, f^0)$ depends on a huge number of internal ('hidden') integration variables variables $z$. Then an approximated log-posterior must be maximized. Exchanging a nonlinear function with the integration is called annealed approximation (Seung, 1995). The replica approach is a special adaption of a saddle point method when the logarithm of a sum $g_i = \ln \sum_j^n e^{L_j}$ has to be averaged with weights $p_i$ (Mezard, Parisi, & Virasoro, 1987). This situation can occur, for example, if algorithms are compared with respect to a large ($n \rightarrow \infty$) number of application situations where the average is over the sampled data.

The second MiR step uses $p(y, q|f^0)$ to find the best alternative $\hat{f}$, i.e. for 'training' of the action model. A numerical evaluation of the integral can be seen as an empirical risk minimization (see Section 8) on a set of *($f^0$–generated) virtual examples* sampled according to $p(y, q|f^0)$, like a numerical evaluation of a full Bayesian approach can be based on $f$–based virtual examples generated according to $p(y, q|f)$.

Note that the distinction between a parallel and inverse decision problem, corresponding to choosing $p(\hat{y}|q, \hat{f})$ or $p(\hat{q}|y, \hat{f})$ to define the action model, only matters for the MiR and not the MaP step. Also, the model needs not necessarily to be reduced to only one remaining state. Thus, MaP–MiR can be generalized with a first step reducing the space $F^0$ to a smaller $\tilde{F}^0$ and a second step minimizing the risk within that $\tilde{F}^0$. But using more than one state for the minimization step, like taking for example the $n$ most probable states or skipping only very unlikely ones, requires calculation of the relative weights of the states related to second derivatives with respect to the parameters of $f^0$. Every decision problem defines an optimality mapping by $\hat{f}(f^0) = \operatorname{argmin}_{\hat{f} \in \hat{F}} r(f^0, \hat{f})$ for all $f^0 \in F^0$ and analogously for $f \in F$. In principle, one can restrict the search space $\hat{F}$ by eliminating $\hat{f}$ being never optimal. In the case available data have equal probability within $f^0 \in [f^0]_{r^*}$ i.e. they do not distinguish between states $f^0$ leading to the same decision $\hat{f}$ identifying those makes optimality mapping one–to–one. With respect to such a construction the space $F^0$, and therefore also the MaP step as maximization over $f^0 \in F^0$, is not independent of the loss function.

### 7.2.1   Perturbation theory beyond MaP

The MaP approximation can be extended by expanding the exponential to higher orders around the Gaussian reference point. We already mentioned that higher order contributions can be obtained by including all linked diagrams according to Wick's theorem. In general, a perturbation theory can be built upon any reference point based upon the general formula (See the Section about Heim's perturbation theory in Jaynes, 1996)

$$e^{A + \epsilon B} = e^A \left[ 1 + \sum_{n=1}^{\infty} \epsilon^n S_n \right],$$

with

$$S_n = \int_0^1 dx_1 \int_0^{x_1} dx_2 \cdots \int_0^{x_{n-1}} dx_n \prod_1^n B(x_n),$$

and

$$B(x_n) = e^{-x_n A} B e^{x_n A}.$$

Here $A$, $B$ stand for matrices, $x$ for a real number and $xA$ means multiplication of each entry of $A$ by $x$. In matrix notation we have

$$r(f, \hat{f}) = < \tilde{l}(q, y, \hat{f}) >$$

$$= \operatorname{Tr}(\rho(q, y, f^0, f)\tilde{l}(q, y, \hat{f})) = \operatorname{Tr}(e^L(q, y, f^0, f)\tilde{l}(q, y, \hat{f})),$$

with $L(q, y, f^0, f) = \ln(p(f^0|f)p(q)p(y|q, f^0))$, the trace denoting the integrals $\operatorname{Tr} = \int df^0 \int dq \int dy$ for diagonal matrices $L(f)_{i,j}$, $\tilde{l}(\hat{f})_{i,j}$ with indices $i = (q, y, f^0)$, $j = (q', y', f'^0)$, and $\rho(q, y, f^0, f) = e^L(q, y, f^0, f)$. Because the trace is invariant under similarity transformations $S$ this formulation allows, if convenient, to work with nondiagonal $L' = SLS^{-1}$, $\tilde{l}' = S\tilde{l}'S^{-1}$.[52]

If we now write $L = A + \epsilon B$, the expected risk can be expressed completely by unperturbed expectations $< \cdots >_A$ with respect to the reference $e^A$

$$< \tilde{l} > - < \tilde{l} >_A =$$

$$\sum_{n=1}^{\infty} \epsilon^n (< Q_n \tilde{l} >_A - < Q_n >_A < \tilde{l} >_A),$$

with

$$Q_1 = 1; Q_n + S_n - \sum_{k=1}^{n-1} S_k < Q_{n-k} >_A, \quad n > 1.$$

## 8   The Frequentist approach

### 8.1   Empirical risk minimization

The Frequentist paradigm is related to the general plug–in or bootstrap principle.[53] It is also called Monte Carlo estimate if applied to calculate expectations (Efron, B. & Tibshirani R.J. 1993). It is based on results of the

---

[52] See for example Derka, Bužek, Adam, & Knight (1996) for Bayesian inference with density operators used in quantum theory.

[53] Often the term bootstrap refers to the special case of estimating the standard error of some sample estimate which then requires resampling.

theory of uniform convergence. A functional of a population distribution is estimated by applying the same functional to an empirical sample drawn according to that distribution.[54] In our case the risk functional $r$ is replaced by an empirical risk $\hat{r}$: Training questions $q^D$ are generated according to $p(q^D|y_c^D, \hat{y}_c^D)$ (or usually $p(q^D)$) and the risk functional $r$ is applied to the empirical distribution of

$$l(q_i^D, y_i^D, \hat{y}_i).$$

In the case the risk functional is chosen as the expectation and $p(q|y_c, z_c) = p(q|y_c)$ we define an integrated loss function

$$\int dz\, p(z|q, y, \hat{f}) l(q, y, z) = \tilde{l}(q, y, \hat{f}),$$

Specifically, in parallel decision problems $z$ is equal to $\hat{y}$ and in inverse problems equal to $\hat{q}$. Note that in $\tilde{l}(q, y, \hat{f})$ $\hat{f}$ just denotes parameters. Yet, for a given loss function we can always introduce some effective deterministic function $\tilde{f}(q)$ (or $\tilde{f}(y)$) (not uniquely defined and possibly vector valued) containing the dependence of $\tilde{l}$ from the parameterization of $\hat{f}$ and (part of the dependence from) $q$ (or $y$) according to $\tilde{l}(q, y, \hat{f}) = \tilde{l}(q, y, \tilde{f}(q))$ or $\tilde{l}(q, y, \hat{f}) = \tilde{l}(q, y, \tilde{f}(y))$. That means every decision problem with $p(q|y_c, z_c) = p(q|y_c)$ can be seen as equivalent to a decision problem with deterministic function $\tilde{f}$. For a parallel decision problem with a deterministic $\hat{y}$–producing device we choose the functional dependency of $\tilde{l}$ from its first argument $q$ to be equal to that of $l$ from its first argument $q$, which gives $l = \tilde{l}$ and $\hat{f}(q) = \tilde{f}(q)$. Alternatively, we can simplify the notation by absorbing the first $q$ into the definition of $\tilde{f}(q)$ and write in such cases $\tilde{l}(y, \tilde{f}(q))$. Analogously, in an inverse setting we define $\tilde{f}(y)$ and can choose in the deterministic case $\tilde{f}(y) = \hat{f}(y)$.

The expected risk

$$r(f^0, \hat{f}) = \int dq \int dy\, p(q, y|f^0) \tilde{l}(q, y, \hat{f}),$$

with $p(q, y|f^0) = p(q|y_c) p(y|q, f^0)$ is approximated by

$$\hat{r} = \frac{1}{n} \sum_i^n \tilde{l}(q_i^D, y_i^D, \hat{f}) = \frac{1}{n} \sum_i \tilde{l}_i(\hat{f}),$$

with the data sampled according to $p(q, y|f^0)$.

The plug–in principle assumes the distribution of training data $D$ including the distribution of questions to be identical to that of the test data $\hat{D}^l$ including $p(q^l)$ for the relevant $q^l$ for which we want to calculate the functional $r$. This holds for example in a setting where both training and test set are generated by the same (stationary) device. Then an explicit knowledge of the generating distribution is not necessary.

Because the empirical risk values are used to make the decision, the chosen $\hat{f}$ depends on the training data and

---

[54]As empirical distributions of finite samples are a quite restricted class of distributions, functionals equal for them can differ on general distributions.

the empirical risk of the chosen $\hat{f}$ does not approximate its expected risk. To estimate the expected risk one has to reevaluate the risk for the chosen $\hat{f}$ with new, independent sample data, the empirical test set, not involved in the decision. (Not to be confused with the (true) test set $D^l$ in the definition of the decision problem.) Bounds for the difference between the expected risk of the true optimal $\hat{f}$ and the chosen $\hat{f}$ are given by the theory of uniform convergence (Vapnik, 1982; Dudley, 1984; Pollard, 1984; Haussler 1995, for an introduction e.g. Kearns & Vazirani 1994). These bounds are based on what we called structural information. Their local part, like bounds for absolute values or the local variance, allows locally the application of probability theoretic inequalities like Hoeffding's or Chebyshev's inequality. Their nonlocal part, formulated for example as finite $\epsilon$-entropy, or finite pseudo or VC dimension of a set $\tilde{l}(q, y, \hat{f})$ for $\hat{f} \in \hat{F}$, allows generalization. Then, when the training data are i.i.d. sampled according to the test distribution, the theory of uniform convergence gives bounds on the probability of deviations of the empirical risk from the expected risk in the true state $f^0$. For example, (Vapnik, 1995) gives bounds $\delta(f, N)$ in terms of the VC dimension for

$$p^{VC} = \sup_{f^0} p(\sup_{\hat{f}} |\hat{r}(D(N, f^0), \hat{f}) - r(f^0, \hat{f})| > \epsilon | f^0),$$

for bounded risk. The supremum over true states $f^0$ (with $p(f^0|f) \neq 0$) is implicit in the results using the definition of the VC dimension and can be replaced by $\forall f^0$. Using worst case considerations there is no need in this theory to explicitly calculate posterior probabilities.

There are some recent studies how specific nonlocal information affects the bounds of uniform convergence (Abu–Mostafa, 1990, 1993a, 1993b; Ratsaby, Maiorov, 1996). But for general nonlocal information a reformulation in terms of VC dimension or $\epsilon$-entropy is often difficult or even practically impossible, and one has to use an upper bound for them, giving the results of the theory of uniform convergence another worst case interpretation. Then, especially when only few local data are available, the uniform bounds can be trivial or weak.

## 8.2 Vocabulary and framework

We formulate the Frequentist setting in a decision theoretic language. Assume the availability of a stationary sampling process $S$ to generate training questions $q \in Q^D$ according to some $p^S(q)$ for which answers are available and the ratio with the distribution of relevant questions $p(q)/p^S(q)$ is known. We call the set $Q^S$ of questions with $p^S(q) \neq 0$ the sampling population (for questions) and specify for the present context $Q^D$ to be the set of sampled or training questions with answers used for the plug–in estimate. Including previously defined sets of questions we have the following listing

1. $Q^S$ the *sampling population* with $q \in Q^S$ the sampling questions,

2. $Q^D$ the set of sampled or *training questions*,

3. $Q^0$ the set of *prior questions* being the questions with data available but with $q \in Q^0$ not sampled according to $S$ or not used for the plug–in estimate,

4. $Q^l$ the set of *relevant questions*,

5. $Q^c = Q^l \setminus Q^S$ the set of *cost questions*,

Data $(q_i, y_i)$ are obtained by using measurement devices for $q_i$ to find results $y_i$. With respect to a sampling process $S$ we separate the available data into the two groups:

1. *Sampled data* or *training examples* $D^S$, we will write more shortly just $D$ skipping the superscript $S$, which are obtained using the available stationary sampling process $S$ and used to calculate the empirical risk via the plug–in principle. Nonlocal questions can be included in $Q^S$.

2. *Prior data* $D^0$, being all other data with questions generated from other processes $S' \neq S$. Such processes can be unknown processes from the past possibly different from $S$, they can be non–stationary like for active queries, they can use devices measuring other questions or represent active control. All priors $p(f^0|f)$ can be related to a factorial prior $p_{fact}(f^0) = \prod_{x \in X} p(f_x^0)$ by data $D^0$ with $p(f^0|f(D^0)) \propto p_{fact}(f^0)p(D^0|f^0)$. The data $D^0$ are not uniquely defined and the corresponding questions need not necessarily to be in the sampling population $Q^S$. Data which are sampled according to $S$ but not used for the plug–in principle are not sampled data but prior data.

For a Bayesian treatment the distinction is not important, but in a Frequentist approach only sampled data are used for the plug–in principle while prior data only enter in form of restrictions of the space $\hat{F}$. Note that this use of the term prior does not refer to the temporal aspect meaning information collected previously to the data $D$. We made things simpler by not trying to distinguish non–sampled non–prior data from priors. Thus, identification of the reference factorial prior as well as assumed and not measured data are understood as being part of the prior data $D^0$. In another context it might well be convenient, but not necessary, to distinguish non–sampled data from prior data with the latter having a more temporal connotation or referring to assumed and not measured data.

According to this distinction of data we can also split the log-posterior into a sampled and prior part

$$L(D, D^0, f^0) = L^D(D, f^0) + L^0(D^0, f^0)$$
$$= \sum_i L^D(y_i^D | q_i^D, f^0) + L^0(f^0).$$

Introducing a dummy index one can always achieve $Q^0 \cap Q^S = \emptyset$ even if for a question sampled and prior data are available at the same time. Non–relevant sampled questions $q \notin Q^l$ can be excluded from $Q^S$ because for them $p(q) = 0$. If a sampling process nevertheless produces them they can instead of simply being eliminated technically be treated like prior data influencing $p(f^0)$. The property of being sampled is for them not important and we choose in the following $Q^S \subseteq Q^l$. Including nonlocal questions like smoothness into $Q^S$ instead of nonlocal prior data could also enable generalization. But this requires nonlocal questions also to be relevant,

i.e. in $Q^l$, what is usually not the case. Also we prefer to treat the $q \in Q^l$ as being independent without priors.

On the other hand we do not assume all relevant questions necessarily to be available for sampling. That allows $Q^l \supset Q^S$ which represents situations where only part of the expected risk can be estimated by sampling. We call the remaining part *costs*. Complexity costs are a typical example. It must be determined by other informations which could come from another sampling process, a Bayesian calculation, or if $Q^c$ is finite and deterministic from a complete set of answers. Fig.8 shows the relations graphically.

We define prior data for question $q$ to be in the set of

a. *Hints* $D_h^0$ if questions from $Q^S$ depend on them, that is if $X^q \cap X^{Q^S} = X^S \neq \emptyset$, (Here $X^q$ denotes the basis of prior question $q$.)

b. *Cost priors* $D_c^0$ if cost questions from $Q^c$ depend on them, that is if $X^q \cap X^{Q^c} = X^c \neq \emptyset$. Here $X^q$ denotes the basis of prior question $q$.

A specific prior can be both and therefore the two sets need not to be disjunct.

Splitting the $q$–integrations of the expected risk into a *sampled risk* $r^D$ with $q \in Q^S$, determined by answers to $q \in Q^D \subset Q^S$ sampled by $S$, and an additional *(non–sampled) cost term* $r^0$ for $q \in Q^c$, determined by $D^0$, gives

$$r(f, \hat{f}) = r^D(f(D, D^0), \hat{f}) + r^0(f(D, D^0), \hat{f}),$$

with $r^D(f, \hat{f})$ being

$$\int_{Q^S} dq \int dy \int df^0 p(f^0|D, D^0)p(q)p(y|q, f^0)\tilde{l}(q, y, \hat{f}),$$

and analogous for $r^0(f, \hat{f})$ with $\int_{Q^c} dq$. We write $\tilde{l}^D$ for the part of the loss function depending on $q \in Q^S$, that means $\tilde{l}^D(q, y, \hat{f}) = \tilde{l}(q, y, \hat{f})$ for $q \in Q^S$ and zero otherwise, i.e. if $q \in Q^c$. Analogously, we write $\tilde{l}^0$ for the part depending on $q \in Q^c$. The parts are only defined up to $q$–independent terms, because such terms can be shifted between questions without changing the risk. We can use this freedom for a convenient choice.

Cases in which available data for questions in the loss function are partly sampled as well as not sampled are included in this definition by duplicating questions in $Q^S$ with a dummy index,

$$\tilde{l}(q, y, \hat{f}) \rightarrow \tilde{l}^D(q, y, \hat{f}) + \tilde{l}^0(q', y', \hat{f}),$$

if we define $q' \in Q^c$ for $q \in Q^S$.

The cost term can depend in general on all data including the training examples $D$. In such cases one may call $r^0(f, \hat{f})$ *posterior (non–sampled) costs* or *data dependent (non–sampled) costs* as they have to be determined after having seen all the data. If $p(f^0|D, D^0) = p(f_{Q^S}^0|D, D_h^0)p(f_{Q^c}^0|D_c^0)$ with $D_h^0 \cap D_c^0 = \emptyset$ then the $f_{Q^S}^0$–integration vanishes. Understanding $f^0 = f_{Q^S}^0$ to be the restriction on $Q^S$ we can write

$$r(f, \hat{f}) = r^D(f, \hat{f}) + r^0(\hat{f}),$$

and call $r^0(\hat{f}) = \tilde{l}^0(\hat{f})$ *prior costs*. Being $q$–independent, $\tilde{l}^0$ can also be seen as part of $\tilde{l}^D$.

Examples for possible (prior) costs include, storage requirements for parameters of $\hat{f}$ like the number of weights of a neural network or number of nodes in a decision tree, penalties for on–line evaluation times, criteria related to understandability for human experts, or to an effective and cheap hardware implementation in VLSI technology.

## 8.3 Sampling generalized questions

### Reweighting

Here we consider the case of relevant questions $q \in Q^l$ not directly sampled, i.e. $q \notin Q^S$ or equivalently $q \in Q^c$. We show that for such questions which have a basis $X^q$ of $q$ completely within $Q^S$ theoretically, but often not practically, the sampling process $S$ can be extended to an $S'$ so that $q \in Q^{S'}$. The basic fact used in reweighting for evaluating an integral by the plug–in principle is that the factorization of the integrand into function and probability is not unique,

$$\int dz\, p(z) g(z) = \int dz\, p'(z) \frac{p(z)}{p'(z)} g(z) = \int dz\, p'(z) g'(z),$$

where $g'(z)$ is equal to $g(z)$ multiplied by the reweighting factor $p(z)/p'(z)$. That means, instead of sampling $z$ according to $p(z)$ and summing up $g(z_i)$ for each sample point $z_i$ we can, assuming $p(z)$ and $p'(z)$ are known, alternatively sample according to $p'(z)$ summing up $g'(z_i) = (p(z_i)/p'(z_i)) g(z_i)$. For example, the method of importance sampling (Montvay & Münster, 1994) uses reweighting to reduce the plug–in error by choosing the reweighting factor so that the reweighted function $g'(z)$ is as constant as possible. But note that the factor $p'(z)$ must be a probability and can especially never be negative. In the situation we are discussing here, $p(z)$ corresponds to the relevant distribution of test questions $q \in Q^l$ with $p(q,y|f^0) = p(q)p(y|q,f^0)$ while $p'(z)$ is the distribution available to generate training data.

The simplest case of reweighting is when the set of test and the sampling population of potential training questions are the same , i.e. $Q^S = Q^l$, but have different probability distributions $p^S(q) \neq p(q)$. Then, if the ratio $p(q)/p^S(q)$ is known (not necessarily the test and available training distributions $p(q)$, $p^S(q)$ itself) one can use $p(q)/p^S(q)$ as reweighting factor for the loss function. Thus, the sampling data term of the risk can be related to a sampling process $S$ with $p^S(q)$ by

$$r(f^0, \hat{f}) = \int dq \int dy\, p^S(q) p(y|q,f^0) \tilde{l}^W(q,y,\hat{f}),$$

where the definition of the reweighted loss

$$\tilde{l}^W(q,y,\hat{f}) = \frac{p(q)}{p^S(q)} \tilde{l}(q,y,\hat{f})$$

compensates for deviations between $p(q)$ and $p^S(q)$ and $p^S(q) \neq 0$ according to the definition of $Q^S$. If not stated otherwise we understand in this paper implicitly $\tilde{l}$ to mean $\tilde{l}^W$ or $p(q)/p^S(q) = 1$.
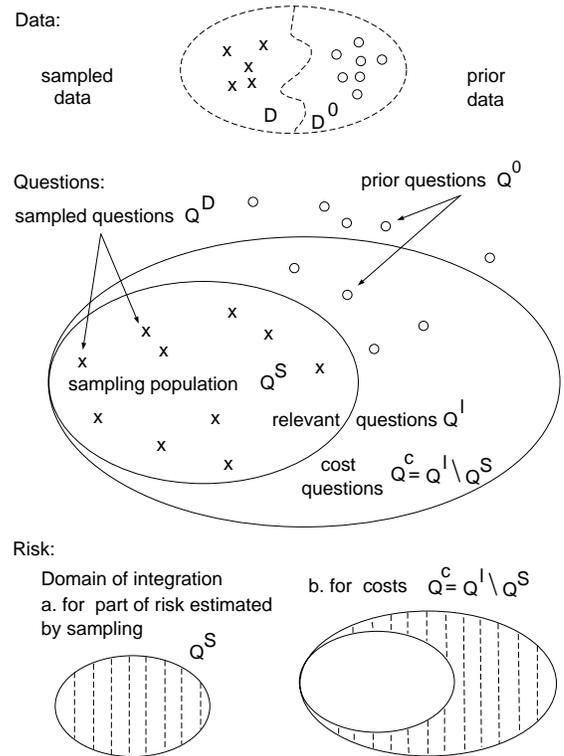


Figure 8: Shown is the distinction for data into the two subgroups of a. sampled or training data being sampled from a stationary process $S$ and b. prior data from any other sources. For questions there are 1. the set of relevant questions $Q^l$, 2. the sampling population $Q^S$ from which training questions are drawn according to a stationary $p^S(q)$, 3. the training questions ($\times$), and 4. prior or non–sampled questions ($\circ$). The sets of training and prior questions are finite. The part of the loss integration depending on $Q^S$ can be determined by sampling training data. The remaining part of the expected risk is called costs and must be determined for example by a Bayesian risk calculation, another sampling process or using exhaustive queries if $Q^c = Q^l \setminus Q^S$ consists of a finite number of deterministic questions.

Only a deterministic function must be included for a $q \in Q^l$ depending deterministically on one $q \in Q^S$. But for a $q \in Q^l$ depending deterministically on two or more $q_i \in Q^S$ the probability being sampled decreases with the product of the probabilities for (independent) $q_i$.

More general is a probabilistic dependence on answers to questions from $Q^S$. Consider relevant questions $q \in Q^l$ given in their dependence on the sample questions $q^S$

$$p(y|q, f^0) = \int dq^S \int dy^S \int dz^S \, p(q^S | y_c^S, z_c^S, q)$$

$$\times p(y^S | q^S, f^0) p(z^S | q^S, y^S, q) \delta(q(q^S, y^S, z^S) - y),$$

with the two $f^0$–independent and $q$–dependent probability factors and the defining function $q(q^S, y^S, z^S)$ assumed to be known. The probability $p(y^S | q^S, f^0)$ depends on the state of nature and is unknown. Then eliminating the $\delta$–function by performing the $y$–integration shows that the empirical risk is a sum of

$$l(q_i, y_i = q(q_i^S, y_i^S, z_i^S), z_i) = l(q_i, q_i^S, y_i^S, z_i^S, z_i),$$

when a sampling procedure for the variables is available according to

$$p(q, z, q^S, y^S, z^S) = p(q|y_c, z_c) p(q^S | y_c^S, z_c^S, q)$$

$$\times p(y^S | q^S, f^0) p(z^S | q^S, y^S, q) p(z | q, y(q^S, y^S, z^S), \hat{f}).$$

Using the corresponding devices the sampling steps within the single components are the following:

$$q \to q^S \to y^S \to z^S \to z$$

Assuming probability distributions $p^S(q^S | y_c^S, z_c^S, q) =$ to be available to generate training data gives a $q$–, $y_c^S$–, $z_c^S$– dependent reweighting factor for the loss function

$$\frac{p(q^S | y_c^S, z_c^S, q)}{p^S(q^S | y_c^S, z_c^S)}.$$

As already discussed for the deterministic case it is a principal problem for nonlocal questions that the $q^S$ can form a vector. Take as example a question measuring answer differences between points $x_1$ and $x_2 = x_1 + \Delta$

$$p(y|q, f^0) = \int dx_1 \int dx_2 \, p(x_1 | q) \delta(x_2 - \Delta - x_1)$$

$$\times p(y_1 | x_1, f^0) p(y_2 | x_2, f^0) \delta(y_1 - y_2 - y),$$

and assume that the $x_i$ are i.i.d. sampled. Here the probability to have a complete pair of two values $y_1(x_1)$ and $y_2(x_2)$ to insert as $y = (y_1, y_2)$ into the loss function has measure zero in standard continuous cases (even for $\Delta = 0$) if $\Delta$ is fixed and not integrated over. If the distribution is not as badly behaved as the $\delta$–function in the example there is still the high dimensionality of $x$.[55] That means the amount of available data is usually too small to sample such nonlocal questions.

---

[55]High dimensionality itself is not the problem if there are enough restrictions for the function. The situation is like in one–dimensional problems as a theorem from Kolmogorov shows that every continuous function of several variables can be expressed as superposition of functions of one variable for closed and bounded input domain (Kolmogorov, 1957). But the nonlocal priors used in the theory of uniform convergence, corresponding for example to a specific VC dimension or $\epsilon$–entropy, and related to a restricted $\hat{F}$, might not be considered as weak and harmless in those cases without empirical foundation.

## Missing values

One could think in remedy the situation by using guesses for the missing values necessary to answer a generalized question. Indeed, the Monte Carlo method or plug–in principle for expectations can be interpreted as replacing missing values by the data mean. This implicitly assumes the loss function to be constant (or compensating) for missing values. For example, the empirical risk estimate can be seen as minimization of a mean square error of the empirical loss with a prior (for the loss, not for $f^0$) $\lambda(\tilde{l}(q, y, \hat{f}) - E(\tilde{l}))^2$, with $E(\tilde{l})$ the expectation of $\tilde{l}$, in the $\lambda \to 0$ limit.

While the plug–in principle replaces missing values $y_x$ by the global sample mean $\sum_x y_x$, one might wish better locally varying approximations of $p(y|x, f^0)$. This can be done by methods of density estimation, including parametric approximations and nonparametric methods like splines, kernel and nearest neighbor methods, (see for example: Silverman, 1986; Härdle, 1990) or neural networks.

Optimal in a Bayesian sense would be sampling according to the posterior probability $p(y|q, f(D))$, but without reference to a specific model of $f^0$ all methods are somewhat ad–hoc.[56] But all those approximation methods used to replace missing values can be interpreted as an approximation of the posterior $p(f^0|D)$ with respect to some (implicit) model of $p(y|q, f^0)$. Specifically, Girosi, Jones, Poggio (1995) relate common interpolation methods (Radial Basis Functions, splines) to regularization terms, i.e. from a Bayesian point of view to priors within a maximum posterior approximation. Also note, that all these approximation methods which one can use to replace missing values are itself decision algorithms and have therefore, within a Frequentist approach, the same problems with general structural information as the Frequentist approach for our original decision problem. Thus, all the aspects we discuss for evaluating the empirical risk for generalized questions appear here again for a specific approximation problem.

### 8.4 Prior data

**Extended loss, indirect priors, and virtual examples**

Common is a situation where prior information is available in addition to training data sampled from the test (relevant) questions. This information can always be seen as corresponding to questions not included in the relevant set $Q^l$. An example are priors like an approximate symmetry and smoothness, being normally not sampled and not included in the set of test questions.[57]

---

[56]Such methods can also be used directly for the loss function. This corresponds to a $\hat{f}$–dependent prior on the loss function, usually difficult to relate to priors on $f^0$. The optimal solution $p(l|f(D), \hat{f})$ of this problem in the Bayesian framework is already nearly equivalent to the solution of the whole decision problem or even more complicated as not the whole loss distribution might be relevant for the risk functional $r$.

[57]We could define an induction problem as a situation where the generating distribution for test and training data

In those cases the loss function can be extended by additional terms to be defined also for data not in the test set. This means, the plug–in principle uses an extended expected risk $\bar{r}$ with respect to an extended set $Q^{\bar{S}}$

$$\int dq^{\bar{S}} \int dy \int d\hat{y}\, p(q^{\bar{S}}|y_c, \hat{y}_c) p(y|q^{\bar{S}}, f^0) p(\hat{y}|q^{\bar{S}}, \hat{f})\bar{l}(\bar{q}^{\bar{S}}, y, \hat{y}),$$

with an extended loss function $\bar{l}$ to approximate the (true) expected risk under the state of nature $f^0$. For example a smoothness property can be included as additional term. The optimal weight of the additional questions, i.e. $p(q^{\bar{S}})$, is usually determined by cross–validation or similar methods (see next Subsection). The extended loss should be chosen so that the relevant features according to the risk functional $r$ of $p(\bar{l}|f, \hat{f})$ are similar to $p(l|f, \hat{f})$ but this requires in general a model of $f^0$ to be determined. In usual approximation problems one chooses an extended loss which enforces $\hat{f}$ to answer similar as $f^0$ also to the not relevant questions as $f$ hoping that this results in similar answers to relevant questions, too. We will relate this ad–hoc method for approximation problems to the maximum posterior approximation within the Bayesian approach. We now show how extra terms can be interpreted as arising from a Lagrange implementation of indirect priors.

An indirect way to include nonlocal information is trying to transform our knowledge $f$ into knowledge about $\hat{f}$. More precisely, if we call the probability of choosing $\hat{f}$, i.e. $p(\hat{f}|f, \hat{a}, q_r)$ an *indirect prior*, one is interested in excluding alternatives $\hat{f}$ with zero (or small[58]) indirect prior probability. As the probability of selecting $\hat{f}$ depends in general not only on $f(D)$ but also on the decision problem, i.e. the $q_r$, including the risk functional, loss function, test data distribution, and the algorithm $\hat{a}$, its complete determination is a much more complicated than the decision problem itself requiring the representation of $f$ in a model with certain $f^0$. Ideally, we are interested in the zero part of the indirect prior resulting from the optimal algorithm defined by $q_r$.

coincide, a transduction problem as a situation where some data do not belong to the test set (Vapnik, 1982). In this formulation, most practical problems are not induction but transduction problems as the nonlocal (e.g. smoothness) information is normally not part of the test set.

[58]Knowledge about the form of the nonzero part of $p(\hat{f}|f, \hat{a}, q_r)$ is, in principle, of no help in searching for the absolute minimum, as the minimization is over all $\hat{f}$ even the unlikely ones, and only the impossible ones can be excluded. Exceptions are cases are where (also subsequent) knowledge of risk values for some $\hat{f}$ can be used to exclude others, i.e. if nonlocal information about the risk functional can be used to exclude certain possibilities. For example in the case of a decision problem with known minimal value, or if one is not looking for the absolute minimum but only for an acceptable minimum, checking $\hat{f}$ with high probability first is of course a good idea. Also in other cases one may ignore $\hat{f}$ with small probability so that the chance of missing the optimum is small. This is a problem on the level of comparing approximations to a decision problem and its analysis requires a corresponding risk and loss to be defined.

For example, in approximation problems with a quadratic loss function $\tilde{l}(q, y, \hat{f}) = (y - \hat{y}(\hat{f}))^2$ we know that the optimal solution is the true regression function $E(y(f^0, x))$ if contained in the search space $\hat{F}$. Then we can simply implement deterministic information about the true regression function of $f^0$ by the corresponding restrictions on $\hat{f}$, that is if we know $p(y|q, f^0) = \delta(q(\{E(y(f^0, x)), x \in X\})) - y)$ we only use $\hat{f}$ with $q(\hat{f}) = \hat{y} = y$.

Sometimes, assuming the existence of a state producing process with stationary distribution corresponding to a possibly unknown but reproducible state of knowledge one may also empirically estimate indirect priors by counting the results of learning algorithms.[59]

While some restrictions for $\hat{f}$ like the range of possible output values or specific symmetries are easily implementable, others, e.g. smoothness, are best taken into account by using the method of Lagrange multipliers. (For the exact conditions under which this is possible see for example Bertsekas, 1995.) Formulating the restriction in the form $q_a(\hat{f}) = a$ this means constructing an extended risk $\bar{r}$ by adding the following extra term to the risk

$$\lambda(a)(q_a(\hat{f}) - a).$$

Here $\lambda$ is the $a$–dependent Lagrange multiplier, the term $-\lambda a$ can be skipped because $\hat{f}$–independent, and $a$, $q_a$ and therefore $\lambda$ can be vectors. Note that for a given problem, including data, $a$ determines $\lambda$. For unknown $a$ determination of $a$ by cross–validation (see next Subsection) can be seen similar to imposing a prior on $a$.

Now we shortly discuss how sometimes additional nonlocal terms in the risk can be approximated by using virtual examples. If $q_a = \int dx\, p(x|q_a)q_a(x, \hat{f})$ is a sum or integral it might sometimes be easier in practice to use only part of the sum if the generalization ability is ensured in another way, for example by restriction of $\hat{F}$. Consider a deterministic $\hat{f}(x)$ with a symmetry or smoothness log-prior

$$\lambda(a)q_a = \lambda(a) \int dx\, w_{q_a}(x)(\hat{f}(x) - \hat{f}(sx))^2.$$

Then one can sample $x$–values according to $w_{q_a}(x)$ if those are positive and can be normalized. This sampling can be done independently of the sampling of the training examples according to $p^S(x)$ (for $x = q^D \in Q^D$). Thus, the term can be approximated by sampling $x$ and calculating $\hat{f}(x)$ and $\hat{f}(sx)$.

For a quadratic loss $\tilde{l}^D(x, y, \hat{f}) = (y - \hat{f}(x))^2$ and some $b$ constant with respect to $x$ with $b(a)p^S(x) = \lambda(a)w_{q_a}(x)$ the sampling for the data and symmetry terms can be combined. As $\lambda$ depends on $a$ so also does $b$. With $\tilde{l}^0 = \lambda(a)q_a$ and splitting the data terms in the empirical risk

$$\hat{r} = (1 - b(a))\hat{r}^D(f^0, \hat{f}) + b(a)\hat{r}^D(f^0, \hat{f}) + \tilde{l}^0.$$

[59]Referring to practical experience or literature about the use of specific priors approximates this method.

The integrations over $x$ and $y$ do not affect $\tilde{l}^0$ being independent of those variables, and we get for the last two parts

$$b(a) \int dx \int dy \, p^S(x) p(y|x, f^0) \left( \tilde{l}^D + \tilde{l}^0 \right)$$

$$= b(a) \int dx \int dy \, p^S(x) p(y|x, f^0)$$

$$\times \left( (y - \hat{f}(sx))^2 - 2(y - \hat{f}(x))(\hat{f}(x) - \hat{f}(sx)) \right).$$

When the second term in the last line vanishes, the whole integral can therefore be calculated approximately by using *virtual (input) examples*[60] $(sx, y(x))$ Abu-Mostafa, 1990, 1993a, 1993b ('hints'), Pomerleau, 1991 (ALVINN), Sietsma & Dow, 1991 (training with noise in practice), Vetter, Poggio, & Bülthoff, 1992 (virtual views of an object), Girosi & Chan, 1995 (for RBF), and more theoretical Webb, 1994, Leen, 1995, and Bishop (1995ab) who gives the explicit form of regularization terms, for quadratic and for cross–entropy error functions, for infinitesimal translations, by expanding $\hat{f}(sx)$ in a Taylor series around $\hat{f}(x)$. When $\hat{f}(x) \neq \hat{f}(sx)$ then all non-optimal $\hat{f}(x)$ are not equal to the regression function. For those $\hat{f}(x)$ the second term does not vanish and has to be considered in the actual minimization procedure. But when $(\hat{f}(x) - \hat{f}(sx))$ is zero or the optimal regression function is in $\hat{F}$ the second term vanishes at the minimum.

## Stratification: (cross–)validation and structural risk minimization

A trivial toy example may clarify the basic idea: Assume two deterministic questions $x_1$ and $x_2$ and a set of possible answers $Y = \{0, 1\}$ corresponding to the four possible functions $\hat{f}_i$ characterized by their answers $(\hat{f}_i(x_1), \hat{f}_i(x_2))$: $\hat{f}_1 : (0, 0)$, $\hat{f}_2 : (1, 1)$, $\hat{f}_3 : (1, 0)$, $\hat{f}_4 : (0, 1)$. Let us sample data $D$ until we have answers to both questions, for example $x_1 = 0$ and $x_2 = 1$. Clearly, $\hat{f}_4$ would minimize the mean square error. Now we perform the same minimization in a *hierarchical* way. We form the two groups (strata) of smooth functions $S_1 = \{\hat{f}_1, \hat{f}_2\}$ and non–smooth functions $S_2 = \{\hat{f}_3, \hat{f}_4\}$ and use the first data point, for example $x_1 = 0$ to minimize within the two groups finding $\hat{f}_1$ and $\hat{f}_4$. To decide between the optima of the two groups we sample more data until we get $x_2 = 1$. Again, we choose $\hat{f}_4$. Note that also forming non–disjunct, overlapping groups, e.g. $S_i \subset S_{i+1}$, would lead to the same result as long as they include all four functions. However, a stratified search is not always equivalent, to a full search. The difference is, that the new data are only used to decide between

'winners' of the strata, and 'loosers' are not reconsidered again. They may however better fit the complete data.

In the case of too large sets $\hat{F}$ the difference between the minima of the expected (training) risk for the empirically chosen and the optimal $\hat{f}$ can become too large (Vapnik, 1982) and solutions can depend too strongly (non–continuously) from the data (Tikhonov, 1963). Than it is necessary to restrict the minimization to a simpler subtask (regularization).

Subsets or strata are defined by some deterministic question $q_a$ requiring $q_a(\hat{f}) = s$. Then we search for the minimum first within the strata and compare the best solutions of different strata in a second step. We want to select $q_a$ so that minimization within each stratum $s$ is possible[61]. Selecting a specific stratum is equivalent to implementing an indirect prior, discussed in the last paragraph. Practically, it can be done by direct restriction of function values like using hardwired symmetries, or, more indirectly, by restricting parameter sets. Examples include choosing the number of nodes or the initial values of the weights of a neural network or the learning algorithm used. This can create overlapping strata, but this is no principal problem and only means that some solutions are considered more than once in the search process. (This is indeed the normal case for nonlocal $q^a$ like smoothness, where the value of $q^a$ can be increased without changing any function value at a data point.) All these constraints are easily kept constant during the minimization (learning) process.

Some constraints like smoothness might be difficult to implement directly. Then again, this minimization with constraints can be done by the method of Lagrange multipliers producing an extra term to be added to the empirical risk

$$\lambda(s)(q_a(\hat{f}) - s).$$

The only difference to the preceding paragraph is that the optimal value of $\lambda$ is not yet determined. The comparison between strata to find $\lambda$ (and therefore $s$) has to be done with an independent data set. That means we have to separate the available training data in a training set to be used within the strata and a training set to decide between the strata also called validation set[62]. Repeating the procedure several times with the same total data set but with different splittings into training and validation set is called *cross–validation* (Stone, 1974, 1977, Allen, 1974). A stratification method for strata where the VC dimension can be calculated is the method of structural risk minimization where the validation step is replaced by minimizing the worst case empirical risk.

---

[60] In this special case the input $sx$ is newly generated, while the same target $y$ for $x$ is used again for $sx$. We also used the term virtual examples in the Bayesian approach for sampling from the posterior probability $p(y|q, f)$ where for risk minimization both, $q$ (e.g. $x$) and $y$, are generated according to the posterior.

[61] Low VC dimension or $\epsilon$–entropy in the theory of uniform convergence or compact for a positive real $s$ and continuous one-to-one extremal equation in regularization theory. See (Verri & Poggio, 1986) for examples for $q_a$. The existence of a stable, unique solution is the result of practical interest in regularization theory. For finite, noisy cases the convergence theorems are not of so much practical value.

[62] The estimation of the true expected risk after the cross–validation procedure would require a third set of available data, the empirical test set.

Fig.9 illustrates that the selection of proper strata depends on indirect priors for the example of structural risk minimization. (See also Wolpert, 1994b, Ripley, 1996)

In general the minimization within some specific strata (e.g. for those strata including functions where a given smoothness functional is not even defined) or also between strata can be very difficult or even impossible. To make those problems solvable an indirect prior for $\hat{f}$ must be available to restrict the number of strata and also to enable minimization between strata. Practically, also the maximal number of different $\lambda(s)$ which can be considered is normally at least restricted by the available computational resources. Therefore, there is more what choosing a good stratification variable $q_s$ can do. It can make the minimization problem between the strata easy (for one or more algorithms under consideration). This is an algorithmic specific aspect not directly related to prior information about $f^0$. It means that $q_a(\hat{f})$ approximates already relevant $f$–independent aspects of $r(f, \hat{f})$. And indeed the most figures found in the literature plotting the empirical error against, for instance, some smoothness related $s$ or $\lambda(s)$ show a very simple one–minimum structure. In principle, even for restricted range of values for $s$ such a function could look arbitrarily wild having for example a multiminima or even random–like structure.

If information is available which restricts the range of $s$ this is a conceptual clear case of an indirect prior. However what one usually does is more like the following: begin with a starting value $s_0$ and explore every of its components in a given direction until the first common minimum is found and stop there. That is, one assumes a certain form of $s$–dependence of the risk $r(f, \hat{f}^*(s))$ for the $\hat{f}^*$ optimal for $s$.

### 8.5 Bayesian interpretation of the Frequentist approach

#### 8.5.1 Approximation and non–approximation loss

We defined an integrated loss function $\tilde{l}(q, y, \hat{f}) = \int d\hat{y} p(\hat{y}|q, \hat{f}) l(q, y, \hat{y})$ for test questions with $p(q|y_c, z_c) = p(q|y_c)$ or simply $p(q)$ as we will choose for the Frequentist setting. This integrated loss function has the same arguments as a log-posterior, except that $\hat{f}$ and $f^0$ are exchanged.

To interpret the Frequentist approach from a Bayesian point of view we choose the same parameter space for $\hat{F}$ and $F^0$. This corresponds to a one–to–one mapping, which we call *parameter mapping*, between $\hat{F}$ and $F^0$ which we can use to identify

$$\hat{F} = F^0 .$$

Specifically, $f^0 = \hat{f} \Leftrightarrow \forall q : f^0(q) = \hat{f}(q)$ when $f^0$ can be parameterized by a deterministic function $f^0(q)$. The *optimality mapping* for the risk functional $\hat{f}(f^0) = \mathrm{argmin}_{\hat{f}} r(f^0, \hat{f})$ defines another mapping between $F^0$ and $\hat{F}$. We already discussed that a one–to–one optimality mapping is obtained if $f^0$ which lead to the same optimal $\hat{f}$ can be identified, and by excluding
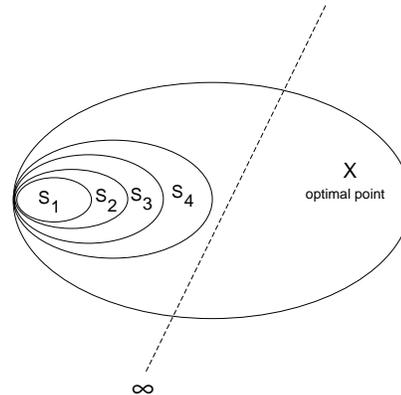


Figure 9: The VC horizon: Consider a set $\hat{F}$ with a subset $S$ of infinite VC dimension. If the sequence of sets $S_1 \subset S_2 \subset S_3 \subset S_4 \subset \cdots S_n \subset S$ does not contain the optimal $\hat{f}^*$ (marked by $\times$ in the figure) then the VC dimension of $S_n$ becomes infinite for $n \to \infty$ before the optimal $\hat{f}^*$ is found. On the other hand, choosing for example $S_1 = \{\hat{f}^*\}$ to be (or to include) $\hat{f}^*$, the optimal solution is already found in the first step. The final uniform bound is lower, if $\hat{f}^*$ is found earlier in a smaller $S_i$. The VC bound of Structural Risk Minimization depends therefore on the choice of the sequence $S_i$. The probabilistic aspects of this choice, however, do not enter the VC bound. Strictly deterministic, i.e. uniform, prior information only defines the set $S$ and not a chain of sets $S_i$. Here prior information $p(\hat{f})$ has the form of an indirect prior, i.e. it is about optimal actions and not states of nature. (Only in saddle point approximation for approximation problems, i.e. $l \propto -\ln P + c$ where we can identify $f^0$ with $\hat{f}$, this is equivalent to a direct prior $p(f^0)$ for state $f^0$.) If $\hat{F}$ is large, a prior free construction of the chain $S_i$ by uniform sampling will with high probability not contain good candidates. To yield reasonable results the sequence of $S_i$ has to be chosen depending on the probability that $S_i$ contains the optimal solution, i.e. depending on probabilistic prior information. Indeed, if we can attribute to most $\hat{f}$ very low probabilities to be the optimal solution, this increases the chance of testing good candidates.

$\hat{f}$ which are for no $f^0$ optimal. For such a construction the $f^0$ related to the $\hat{f}$ by the optimality mapping can be seen as effective states of nature for the decision problem involving $\hat{f}$.

We called a loss function approximation loss if for all $\hat{f} = f^0$ (See Section 5.3.4)

$$\tilde{l}(q, y, \hat{f}) = -c_1 L^D(y|q, f^0) - c_0,$$

with $\hat{f}$–independent constants $c_0$, $c_1$ adjusting the value of the minimum of $\tilde{l}$ or the normalization constant of $L^D$. We remark that an $f^0$–independent $c_0$ corresponds to a fixed normalisation of $\tilde{l}^A$ over $y$ for all $\hat{f}$ and all $q$. Equivalently we may say in this case

$$\tilde{l}(q, y, \hat{f}) = \tilde{l}^A(q, y, \hat{f}),$$

if we define the approximation part of the loss as

$$\tilde{l}^A(q, y, \hat{f}) = -c_1 L^D(y|q, f^0) - c_0.$$

In general this gives a ($c_1, c_0$–dependent) decomposition of the loss function into an *approximation* and *non–approximation part of the loss*

$$\tilde{l}(q, y, \hat{f}) = \tilde{l}^A(q, y, \hat{f}) + \tilde{l}^{NA}(q, y, \hat{f}).$$

Analogously, we define the (loss function dependent) set $Q^A$ of *approximation questions* to consist of all questions $q$ for which we can choose for all $\hat{f} = f^0$ and all $y$

$$\tilde{l}(q, y, \hat{f}) = -c_1 L^D(y|q, f^0) - c_0. \qquad (20)$$

To achieve a parallel notation for loss and log–probabilities we will use in this case the convention $l^A(q, y, \hat{f}) = l^D(q, y, \hat{f})$ for $q \in Q^S$, i.e. if possible we choose the sampled loss equal to the approximation part.

For those questions a large log-posterior for $f^0$ is equivalent to a small loss for $\hat{f}(f^0)$ which characterizes the situation for $q$ in a parallel decision problem as an approximation problem. In an inverse setting this is not necessarily true as for example a loss $l^D = (q - \tilde{f}(y))^2$ measuring the reconstruction quality of $\hat{q} = \tilde{f}(y)$ would correspond to a log-posterior $L^D = (q - f^0(y))^2$ instead of $L^D = (y - f^0(q))^2$.

The decomposition of the loss function into an approximation and non–approximation part of the loss induces the same decomposition for the risk

$$r(f, \hat{f}) = r^A(f, \hat{f}) + r^{NA}(f, \hat{f}).$$

Especially common is the case with prior costs and $Q^S \subseteq Q^A$ for which we have (using our convention for $\tilde{l}^D$ in such cases)

$$\begin{aligned}
\tilde{l}(q, y, \hat{f}) &= \tilde{l}^A(q, y, \hat{f}) + \tilde{l}^{NA}(\hat{f}) \\
&= \tilde{l}^A(q, y, \hat{f}) + \tilde{l}^0(\hat{f}) \\
&= \tilde{l}^D(q, y, \hat{f}) + \tilde{l}^0(\hat{f}).
\end{aligned}$$

Prior costs have the form of a fixed additional term implementing restrictions within $\hat{F}$ according to the method of Lagrange. Table 3 summarizes decompositions of loss, risk, and log-posterior:

**Decompositions**

| | |
|---|---|
| Sampled vs. Non–sampled | |
| Log-posterior | $L(D, D^0, f^0) = L^D(D, f^0) + L^0(D^0, f^0)$ |
| | $= \sum_i L^D(y_i|q_i, f^0) + L^0(D^0, f^0)$ |
| risk | $r(f, \hat{f}) = r^D(f, \hat{f}) + r^0(f, \hat{f})$ |
| Approximation vs. Non–approximation | |
| risk | $r(f, \hat{f}) = r^A(f, \hat{f}) + r^{NA}(f, f)$ |
| loss | $\tilde{l}(q, y, \hat{f}) = \tilde{l}^A(q, y, \hat{f}) + \tilde{l}^{NA}(q, y, \hat{f})$ |
| Special case: $Q^S \subseteq Q^A$ with prior costs | |
| loss | $\tilde{l}(q, y, \hat{f}) = \tilde{l}^A(q, y, \hat{f}) + \tilde{l}^0(\hat{f})$ |
| | $= \tilde{l}^D(q, y, \hat{f}) + \tilde{l}^0(\hat{f})$ |

Table 3: Decompositions of log-posterior, loss function, and risk

We now analyze the MaP–MiR procedure. To find the optimal $\hat{f}^*$ in a MaP–MiR approximation the first MaP step is followed by a second MiR step to find $\hat{f}^* = \mathrm{argmin}_{\hat{f}} r(f^0, \hat{f})$ for the most probable state $f^0 = f^{0,*}$. Despite the one–to–one mapping between $F^0$ and $\hat{F}$ and the same functional dependence of $L^D$ and $-\tilde{l}^D$ for $q \in Q^A$ there is no perfect symmetry between the two problems as the maximization is done with respect only to a finite sum over the set of data $D$ and the minimization with respect to the full $y$–, and $q$–integrals conditioned on the pure state $f^{0,*}$ with maximal posterior probability. Only if the mapping between $f^{0,*}$ and $\hat{f}(f^{0,*})$ already corresponds to the optimality mapping

$$\hat{f}(f^{0,*}) = \hat{f}^* = \mathrm{argmin}_{\hat{f}} r(f^{0,*}, \hat{f})$$

$$= \mathrm{argmin}_{\hat{f}} \int dq \int dy \, p(q) p(y|q, f^{0,*}) \tilde{l}(q, y, \hat{f}),$$

with

$$p(y|q, f^{0,*}) = e^{L^D(y|q, f^{0,*})}$$

then finding the state with maximal posterior probability already corresponds to minimizing the loss for it.

We already saw in Section 5.3.4 that for approximation problems the relative risk

$$c_1^{-1}(r(f^0, f^0) - r(f^0, \hat{f})) = K(p(y|q, f^0), p(y|q, \hat{f})),$$

is a Kullback–Leibler entropy and using Jensen's inequality that the optimal solution for the minimization step is

$$\tilde{l}^A(q, y, \hat{f}^*) = -c_1 L^D(y|q, f^{0,*}) - c_0.$$

We derive this well known result here again using an explicit calculation. For continuous parameter spaces $\hat{f}$, a necessary optimality condition for minima not on the boundary is a vanishing gradient with respect to the parameter vector $\hat{f}$

$$\left. \frac{d\, r(f^{0,*}, \hat{f})}{d\hat{f}} \right|_{\hat{f}^*} = 0.$$

We study the stationarity condition[63] for that part of the

---

[63]In case the parametrization of $f^0$ (or $\hat{f}$, respectively) does not ensure normalization, one has to add the normalization condition $\lambda(x)(1 - \int dy p(y|x, f^0))$ for all $x$, where $\lambda(x)$ is a $x$–dependent Lagrange multiplier.

loss integration which depends on approximation questions, i.e. with $q^l \in Q^A$,

$$\frac{d}{d\hat{f}} \int_{Q^A \subseteq Q^l} dq \int dy\, p(q) e^{L^D(y|q,f^{0,*})} \tilde{l}^D(q,y,\hat{f}) \Big|_{\hat{f}^*} = 0.$$
(21)

Interchanging integration and differentiation and using $-c_1 L^D(y|q,f^0) - c_0 = \tilde{l}^D(q,y,\hat{f})$ at $f^{0,*} = \hat{f}^*$ gives

$$\int dq\, p(q) \int dy\, e^{L^D(y|q,f^{0,*})} \frac{d\tilde{l}^D(q,y,\hat{f})}{d\hat{f}} \Big|_{\hat{f}^*}$$

$$= -c_1 \int dq\, p(q) \int dy\, e^{L^D(y|q,f^{0,*})} \frac{dL^D(y|q,f^0)}{df^0} \Big|_{f^{0,*}}$$

$$= -c_1 \frac{d}{df^0} \int dq\, p(q) \int dy\, e^{L^D(y|q,f^0)} \Big|_{f^{0,*}}$$

$$= -c_1 \frac{d}{df^0} \int dq\, p(q) \int dy\, p(y|q,f^0) \Big|_{f^{0,*}}$$

$$= -c_1 \frac{d}{df^0} 1 = 0.$$

That means that for $Q^l = Q^A$ the stationarity conditions of the MiR step are automatically fulfilled if $\hat{F} = F^0$.[64] (Here we discussed the saddle point approximation, however the same result holds in a full Bayesian approach with respect to $f$ (not $f^0$). if $\ln p(y|q,f) = -c_1 \tilde{l}(q,y,\hat{f}) + c_0$ and $F = \hat{F}$.) For $Q^l = Q^A$ the MaP–MiR solution is not altered if $p(q)$ is changed (for fixed data). In particular, terms like $r^A(f^0, \hat{f})$) $= \int dq\, p(q) \int dy\, p(y|q,f^0) l^D(q,y,\hat{f})$, with $l^A(q,y,\hat{f}) = l^D(q,y,\hat{f}) = -c_1 L^D(y|q,f^0) - c_0$ for every $\hat{f} = f^0$, can be added or dropped from the risk (with p(q) readjusted).[65] As long that term is $f^0$–dependent it cannot be seen as part of the ($f^0$–independent) loss, i.e. it cannot be written as $\tilde{l}(\hat{f})$. In this situation approximating one point is equivalent to approximating a whole function. This could also be done for the prior part $L^0$ writing it in its data dependence, giving rise to cost terms $l^0(f^0, \hat{f})$. The non–approximation part of the loss depending on questions $q \in Q^l \setminus Q^A$ can cause a deviation between optimality mapping and parameter mapping. Examples of typical non–approximation loss include

1. time and storage requirements for calculating $\hat{f}$,

2. costs producing a hardware (VLSI) implementatio a general $\tilde{l}(q,y,\hat{f})$ (energy function)n of $\hat{f}$,

---

3. understandability of the structure of $\hat{f}$.

We remark that the normalization condition $\int dy p(y|q,f^0)$ does not allow to choose for a general $\tilde{l}(q,y,\hat{f})$ $L = c_1 \tilde{l} + c_0$ with $f^0$ independent $c_0$. As discussed earlier non–approximation loss related to the algorithm instead of $\hat{f}$ define a higher level decision problem.

The typical example for approximation questions combines Gaussian noise with mean square error. Also, a question is an approximation question if a uniform $L$ is combined with a uniform $l$ finite on the same domain. Here log-probabilities $L$ (and analogously costs $\tilde{l}^0$) are called uniform on $F^0$ (or $\hat{F}$) if they are equal to a constant, i.e. independent of $\hat{f}$ (or $f^0$), on the domain where they are finite. An example are priors implemented by regularization terms to be determined by cross–validation with restricted interval for the regularization constant. Uniform priors can be skipped from the formalism by restricting the parameter spaces $F^0$ to the domain on which the log-priors are finite. The corresponding $\hat{f}(f^0)$ related to $f^0$ with zero prior probability by the optimality mapping can also be skipped. This is equivalent to the introduction of a uniform cost term. If their is no non–approximation part of the loss the optimality mapping is equal to the parameter mapping and trivially implemented by using the same restrictions for $F^0$ and $\hat{F}$.

While we saw that risk (not necessarily loss function) terms corresponding to log-likelihoods can be added to the risk integral without changing the problem, (non–approximation) loss terms on the other hand cannot simply be implemented in the log-posterior.[66] For example, uniform prior costs, equivalent to a restriction of $\hat{F}$, do not lead directly to a restriction of $F^0$. In principle one could regroup the $F^0$ by forming equivalence classes with respect to the restricted $\hat{F}$, but this may destroy the relation (20).

In general, for the MaP step being sufficient

$$\frac{d\tilde{r}^{NA}(f^{0,*}, \hat{f})}{d\hat{f}} \Big|_{\hat{f}=f^{0,*}} = 0$$

must hold. For prior costs $l^0(\hat{f})$ being $f^0$–independent this reads

$$\frac{d\tilde{l}^0(\hat{f})}{d\hat{f}} \Big|_{\hat{f}=f^{0,*}} = 0$$

for every $f^0$ if $\tilde{l}^D$ is approximation loss. This allows special cases where $\tilde{l}^0(\hat{f})$ has a minimum $\hat{f} = \hat{f}^*$ at the maximum $f^0 = f^{0,*}$ of the posterior $L^D + L^0$, i.e. $f^{0,*} = \hat{f}^*$, and $\tilde{l}^D(\hat{f})$ coincides with an approximation risk term $r^A(f^{0,*}, \hat{f})$. Then, even nonuniform costs do not change the optimal $\hat{f}$. However, this can only happen for specific $f^0$, as $\tilde{l}^0(\hat{f})$ is $f^0$–independent and we assume the MaP estimate to be data dependent. Otherwise calculating the MaP approximation would not be interesting. Thus, requiring that, depending on the possible data, every $f^0$

---

[64]For example minimizing a $L_1$ error (sum of unsquared distances) for Gaussian probabilities the MiR condition is not automatically fulfilled. An example is the support vector machine (Vapnik, 1995) for classification which uses a $L_1$ type of error in the non–separable case (which might e.g. result from noise). There the function to minimize consists of two terms, the norm of the weights of the optimal canonical hyperplane $||w||^2$ and the $L_1$ error, whose relative importance $C$ has to be determined for example by cross–validation or prior knowledge. (Individual $C_i$ for each data point $i$ could account for locally differing variances.)

[65]One could paraphrase this observation by: "One always can require what is already there."

[66]In contrast to the previous footnote this could be paraphrased by: "Reality is not always like one wants it to be."

can be the most probable one, i.e. be a MaP estimate $f^{0,*}$, the derivative of the cost term must be zero at every $\hat{f} = f^{0,*}$

$$\frac{d\tilde{l}^0(\hat{f})}{d\hat{f}} = 0.$$

Then, $\tilde{l}^0$ is uniform. Because nonuniform, data or $f^0$–independent costs belong to the non–approximation part of the loss their presence require the full two step MaP–MiR procedure.[67] An example would be using complexity costs to enforce simplicity of a model $\hat{f}$ (e.g. smoothness, sparseness, integer values) independent of the data, maybe even when knowing that this is not true for nature $f^0$ (which might for example allow real values).

### 8.5.2 MaP–MiR and ERM: Priors and prior costs

Often, complexity related prior costs, are included in empirical risk minimization, either by explicit penalty terms, or by choosing a specific structure for hypotheses $\hat{f}$. If those complexity aspects are requirements not related to priors ERM cannot be interpreted as MaP-MiR procedure. For example, a tree classifier might be chosen because it can be obtained rather effectively and/or because the resulting rules are relatively easy to interpret. If, there is no prior knowledge about $F^0$ having tree structure, and at the same time there is a more appropriate parameterization of $F^0$ available, then this one should be used in the MaP step while a tree classifier could be fitted in the MiR step.

We use the results from the previous paragraph to discuss in more detail the relations between the Bayesian and Frequentist point of view of empirical risk minimization in the presence of priors and prior costs (in the following shortly called costs). Consider the three problems:

1. Maximization of the posterior probability (MaP) given data $D$

$$\mathrm{argmax}_{f^0}\left(\sum_i^n L^D(y_i^D|q_i^D f^0) + L^0(f^0)\right),$$

2. empirical risk minimization (ERM) given data $D$

$$\mathrm{argmin}_f\left(\frac{1}{n}\sum_i^n \tilde{l}^D(q_i, y_i, \hat{f}) + \tilde{l}^0(\hat{f})\right),$$

being a sample estimate of

3. a full Bayesian risk minimization (MiR) for fixed $f^0$

$$\mathrm{argmin}_f\left(\int_{Q^S} dq \int dy\, p(q)p(y|q, f^0)\tilde{l}^D(q, y, \hat{f}) + \tilde{l}^0(\hat{f})\right).$$

Within the MaP–MiR procedure the fixed $f^0$ for the MiR step is the result $f^{0,*}$ of the MaP step while for ERM

---

[67]This means as soon as one knows model $\hat{f}$ and nature $f^0$ are not the same one should use two different descriptions for both. Otherwise the knowledge about their difference cannot be used.

$f^0$ is thought to be the true state of nature. We assume the sampling set to consist of approximation questions i.e. $Q^S \subseteq Q^A$ or equivalently $-c_1 L^D(y|q, f^0) - c_0 = \tilde{l}^D(q, y, \hat{f})$ for $q \in Q^S$, choosing in the following $c_1 = 1$ for simplicity. Consider the following cases in which arbitrary constants $c_0'$, $c_0''$ exist, so that

A: ('uniform costs and uniform priors')

$$\tilde{l}^0(\hat{f}) + c_0' = -L^0(f^0) + c_0'' = 0,$$

B: ('uniform costs but nonuniform priors')

$$\tilde{l}^0(\hat{f}) = c_0', \qquad L^0(f^0) \neq c_0'',$$

C: ('nonuniform costs')

$$\tilde{l}^0(\hat{f}) \neq c_0',$$

D: ('costs $\propto$ − priors')

$$n\,\tilde{l}^0(\hat{f}) + L^0(f^0) = c_0',$$

In case A we could choose $\tilde{l}^0 = 0$ by including the constant into $\tilde{l}^A = \tilde{l}^D$. It is the case with all priors and related prior costs already implemented as restrictions of $F^0$ and $\hat{F}$. Here numerical realization of ERM and the MaP step are identical and their interpretation from the Frequentist and Bayesian point of view are fully compatible because a MiR step is not needed. $F^0$ and $\hat{F}$ can be fully identified and there is no need to use a distinct notation for $f^0$ and $\hat{f}$.

In case B we also can choose $\tilde{l}^0 = 0$ so no MiR step is necessary, but the MaP step uses log-prior terms if available. As those terms are not part of the loss function ERM should in principle ignore them. Namely, the nonuniform parts of priors do not enter the procedure of ERM or the (worst case) bounds for uniform convergence. Using the priors nevertheless as cost terms in ERM leads to complete numerical equivalence with the MaP step and therefore the whole MaP–MiR procedure. In this sense priors for Bayesians are related to costs for Frequentists. Note, however, that while here the numerical calculations coincide, they are interpreted as different models of nature, as costs cannot be identified with priors.

Case C requires the MiR step. Therefore, in this case the two step MaP–MiR and the one step ERM differ. Exceptions are MaP results $f^{0,*}$ for which the nonuniform costs have a minimum. (For uniform costs all $\hat{f} \in \hat{F}$ are minima.) MaP–MiR incorporates priors and takes into account the differences between $F^0$ and $\hat{F}$. This might be especially important if they differ strongly because priors and costs are related to different aspects. MaP–MiR is expected to improve ERM in situations where the prior had a strong influence compared to the sampled data in the MaP step and/or the cost term is substantial compared to the data term and related to aspects different from those of the prior. The same remarks apply when the sampled loss $\tilde{l}^D(y|q, f^0)$ itself is non–approximation loss (this means according our convention it cannot be chosen as approximation loss), i.e. $Q^S \notin Q^A$.

Case D seems to show a perfect symmetry between MaP and ERM. Indeed, under these conditions ERM and MaP are numerically identical for the specific $n$

$$\text{argmax}_{f^0} \left( \sum_i^n L^D(y_i|q_i, f^0) + L^0(f^0) \right).$$

$$= \text{argmin}_{\hat{f}} \left( \frac{1}{n} \sum_i^n \tilde{l}^D(q_i, y_i, \hat{f}) + \tilde{l}^0(\hat{f}) \right).$$

But in as far as nonuniform costs are present the MaP step is not sufficient from a Bayesian viewpoint and the MiR step is missing. Indeed, for case $D$ the MiR step takes into account the same function $L^0(f^0)$ again but in form of costs $\tilde{l}^0(\hat{f})$ so the related aspects have a stronger influence in MaP-MiR than in ERM. This means, that, if not already in the minimum, a function $\hat{f}^* \neq f^{0,*}$ can be chosen which a lower $\tilde{l}^0(\hat{f}^*)$ than $\tilde{l}^0(\hat{f} = f^{0,*})$ of the result $f^{0,*}$ of the ERM or MaP step. One may expect this effect usually to be small in practice, as the saddle point approximation assumes a strongly peaked maximum for $L^D + L^0$, usually arising in the limit $n \to \infty$ where case C becomes case B or, for uniform priors, case A.

The theory of uniform convergence does not require all the conditions necessary for a Bayesian interpretation of ERM as long as the training data are sampled according to the (arbitrary) relevant distribution. It applies for general $\tilde{l} \neq -c_1 L^D - c_0$ and also if $\hat{F}$ and $F^0$ are chosen different, like in many examples of computational learning theory. Costs restricting the search space $\hat{F}$, and changing for example its VC dimension, influence the bounds of the theory of uniform convergence. Its bounds do not depend on the form of the finite parts of log-priors $L^0$ as they are worst case considerations. Prior costs $\tilde{l}^0$ do not contribute to the difference between empirical and expected risk (if they are not sampled itself), as they are data independent and included in both. The bounds only depend on the infinite part of prior costs $\tilde{l}^0$ restricting the space $\hat{F}$ which can for example reduce the VC dimension. But in most of these cases ERM differs in method and results from a MaP-MiR approach.

We can summarize the results by saying that for $F^0 = \hat{F}$ uniform costs allow an interpretation of the numerical ERM procedure as two step MaP-MiR procedure for approximation loss so that the MiR step is automatically satisfied. On the other hand, training data being sampled according to the test data distribution allow application of the bounds of the theory of uniform convergence to the results of empirical risk minimization.

For example, take a prior $p(f^0)$ depending on a symmetry property like $S(x, f^0) = (f^0(x) - f^0(-x))^2$. If we choose a uniform prior which is zero for $f^0$ if $S(x, f^0)$ is above a bound $B$ and constant for those $f^0$ below that bound, we can implement the prior as restriction on $F^0$ by excluding $f^0$ with $S(x, f^0) > B$. Through the optimality mapping this corresponds to the same restriction for the $\hat{f} \in \hat{F}$. This is the usual case.

Let us briefly write down the common example of a

(one-dimensional) Gaussian probability

$$p(y|q, f^0) = \frac{1}{\sigma_q \sqrt{2\pi}} e^{-\frac{(y - f^0(q))^2}{2\sigma_q^2}},$$

namely a log-posterior $L(y, f^0(q)) = \ln \sigma_q \sqrt{2\pi} - (1/2)\sigma_q^{-2}(y - f^0(q))^2$ corresponds to a quadratic loss function $l(y, \hat{f}) = (1/2)\sigma_q^{-2}(y - \hat{f}(q))^2$ with $\hat{f}(q)$-independent $\sigma(\hat{f}(q)) = \sigma_q$ skipping the constant. Note, that $f^0(q)$ is no random variable but parameterizes the states $f^0$ and corresponds to the regression function. As the regression function minimizes the mean square error we have $\hat{f}(f^0) = \text{argmin}_{\hat{f}} r(f^0, \hat{f})$ or written explicitly the optimality condition (21) gives

$$\int dy \, p(y|q, f^0) y = \hat{f}(q).$$

because $d\sigma_q/d\hat{f}(q) = 0$ if all $\hat{f}$ have the same $\sigma(\hat{f}(q)) = \sigma_q$. Inserting the form of $p(y|q, f^0)$ and performing the Gaussian integration we find again that the optimality condition is fulfilled. As $\hat{f}(q)$ represents the regression function all deterministic information about the regression function can be incorporated as indirect prior, that is by restricting the search space $\hat{F}$. This holds in particular for a deterministic bound on the smoothness of the regression function. On the other hand for nonuniform costs, like for example higher costs for states with regression far from zero $\tilde{l}^0(f^0) \propto \sum_q f^0(q)^2$, the two step MaP-MiR approximation is not equivalent to an empirical risk minimization even if we choose an prior $L^0 = \tilde{l}^0$. Fig.10 visualizes some of the relations.

We summarize how the classical Frequentist approach of empirical risk minimization with additional (e.g. regularization or penalty) terms can be interpreted as a specific Bayesian model. This 'classical' Bayesian model has the following specifications:

1. Definition of an effective loss function $\tilde{l}$ for $z$-independent generation of test questions.

2. Identification of the (parameter) space of actions $\hat{f}$ with that of states $f^0$.

3. The same function, up to a factor and a constant, is chosen for the $y$-dependent parts of the effective loss $\tilde{l}^D(q, y, \hat{f})$ and log-likelihood $L^D(y|q, f^0)$ depending on the same variables after identifying $\hat{f}$ with $f^0$. (We use a formulation where sampled data correspond to approximation questions, i.e. $Q^S \subseteq Q^A$.)

4. There are no nonuniform costs $\tilde{l}^0(\hat{f})$.

5. The decision relevant risk functional is the expectation functional (Bayesian expected risk).

Under these conditions empirical risk minimization corresponds to an exact risk minimization (i.e. no plug-in estimate) for the state with maximal posterior probability, regardless of how the training questions have been sampled.

Remarks on point 3:

i. It characterizes the situation as an approximation problem in a parallel decision setting.

ii. Note that AND and OR are exchanged in the following sense: When $y_1$ AND $y_2$ has been observed as training data then we assume $y_1$ OR $y_2$ can appear in a test situation and both log-posterior and loss consist of a sum. If we only know $y_1$ OR $y_2$ could have been the training data this would require adding probabilities and not log-probabilities resulting in a non–additive structure for $L$. A non–additive loss function depending on more than one outcome at a time would have the different interpretation of an interaction of losses for repeated tests, i.e. for cases where $y_1$ AND $y_2$ happens.

In physics mean field approximations and classical approximations of field theories are related to saddle point approximations. The relation between the Frequentist approach as a maximum posterior approximation and the full Bayesian approach is for example similar to the relation between classical physics and the path integral formulation of quantum mechanics with a field being the analogon to a pure state $f^0$.

### 8.5.3  MiR perturbation theory

The MiR step requires minimization of the expectation $< \tilde{l} > = < \tilde{l}^A + \tilde{l}^{NA} >$ under the distribution given by the MaP $f^{0,*}$. For a small enough 'perturbation' $\tilde{l}^{NA}$ we may expand $\tilde{l}(q, y, \hat{f})$ around $\tilde{l}(q, y, \hat{f} = f^{0,*})$. To change the location of a minimum we have to go at least to second order

$$\tilde{l}(q, y, \hat{f}) \approx \tilde{l}(q, y, f^{0,*})$$

$$+(\hat{f} - f^{0,*})\frac{d}{d\hat{f}}\tilde{l}(q, y, \hat{f})\big|_{\hat{f}=f^{0,*}}$$

$$+\frac{1}{2}(\hat{f} - f^{0,*})^2 \frac{d^2}{d\hat{f}^2}\tilde{l}(q, y, \hat{f})\big|_{\hat{f}=f^{0,*}} .$$

Normalization is assumed to be ensured by the parameterization of $f^0$, otherwise for example a Lagrange multiplier can be added. The stationarity condition $\frac{d}{d\hat{f}} < \tilde{l}(q, y, \hat{f}) >= 0$ is linear in second order. It gives for the parameter vector $\hat{f}$ the solution, assuming $c_1 = 1$

$$\hat{f}^* - f^{0,*} = -\frac{< \frac{d}{d\hat{f}}\tilde{l}^{NA}\big|_{\hat{f}=f^{0,*}} >}{< \frac{d^2}{d\hat{f}^2}\tilde{l}^A\big|_{\hat{f}=f^{0,*}} + \frac{d^2}{d\hat{f}^2}\tilde{l}^{NA}\big|_{\hat{f}=f^{0,*}} >}$$

$$= -\frac{< \frac{d}{d\hat{f}}\tilde{l}^{NA}\big|_{\hat{f}=f^{0,*}} >}{< (\frac{d}{d\hat{f}}\tilde{l}^A\big|_{\hat{f}=f^{0,*}})^2 + \frac{d^2}{d\hat{f}^2}\tilde{l}^{NA}\big|_{\hat{f}=f^{0,*}} >},$$

where $< \cdots >$ stands for the $y$ and $q$ integrals and we used $\int dy\, e^L = 1$ to get the second line, and therefore for $c_1 = 1$ we have $\int dy\, e^L \frac{d}{d\hat{f}}\tilde{l}^A = 0$ and $\int dy\, e^L \left( \frac{d^2}{d\hat{f}^2}\tilde{l}^A - (\frac{d}{d\hat{f}}\tilde{l}^A)^2 \right) = 0$ at the location $\hat{f} = f^{0,*}$.
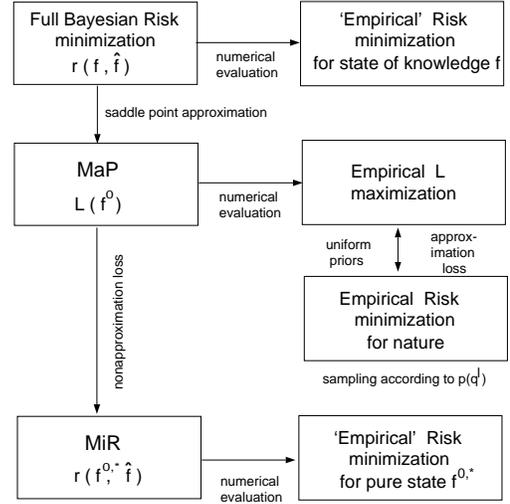


Figure 10: Numerical evaluation of the full Bayesian risk, for example by Monte Carlo methods, is technically the same as an empirical risk minimization (i.e. a use of the plug–in principle) for virtual data generated by $f$. (One might call this also a virtual (empirical) risk, however it is a real empirical risk if state $f$ is prepared as a mixture of $f^0$ according to $p(f^0|f)$.) A saddle point approximation of the full integral gives the MaP problem, which requires to find the state $f^{0,*}$ with maximal posterior probability. If the prior log-probability $L^0$ contains integrals (from nonlocal data) this may require a numerical evaluation of $L$, i.e. use of the plug–in principle. In the 'classical' model with uniform priors already implemented in $F^0 = \hat{F}$ and approximation loss, empirical risk minimization using the sampled data is equivalent to finding $f^0$ with maximal posterior probability. The bounds of the theory of uniform convergence require sampling according to the relevant distribution $p(q^l)$, which, however, can be arbitrary. For approximation loss, i.e. $\tilde{l} = -c_1 L^D - c_0$, the MiR step is not necessary.

58

## 8.6 Occam's Razor

Here is a good opportunity to discuss the celebrated principle of Occam's razor. Occam's razor states: If two theories explain the same phenomena equally well choose the simpler one. From a Bayesian point of view this simply corresponds to including complexity (prior) costs in the decision, but is sometimes also interpreted in the version: Simpler theories have higher prior probabilities.

This can be explained by the fact that empirical risk minimization has no possibility to include $\hat{f}$ dependent complexity in form of (prior) costs $\tilde{l}^0(\hat{f})$ independent of the priors $L(f^0)$. For uniform priors corresponding to uniform costs (case C) Occam's razor is automatically implemented in both versions. For nonuniform costs a second risk minimization step has to be included or if using a one step procedure Occam's razor has to be implemented via the priors. In contrary to the first case where the prior and cost versions of Occam's razor are equivalent, in the second case the Frequentist and Bayesian interpretations of the extra terms do not coincide. Then the prior version of the razor is the Bayesian interpretation of what appears for the Frequentists as the cost version. Relying on the cost version as the intended one, a Bayesian approach approximated by a two step MaP–MiR approximation differs from a one–step Frequentist approach.

In principle, within a Bayesian framework the concept of costs and priors are independent.[68] A MaP–MiR approximation is possible with arbitrary costs independent of priors, when such costs are included in the second risk minimization step and the MaP step is still justified. For example, a sparsity constraint can come from the implementation of complexity costs related to computational costs for nonzero numbers or from the prior information that the actual data are produced by a small number of prototypes. Both aspects cannot be modeled independently using ERM but could be taken into account using the two step MaP–MiR procedure.

## 9 Stationarity equations

### 9.1 Data for generalized questions

In this section we study the stationarity equations (or mean field equations) to find extrema of the log-posterior. There are a variety of methods to find an extremum of the posterior probability. If those are nonlinear they have usually be calculated by iteration. Gradient based methods, or EM (Expectation–maximization)

related algorithms (Dempster, Laird, Rubin, 1979, Tanner, 1993, Gelman, Carlin, Stern, Rubin, 1995) are special iteration schemes. In what we called a classical Bayesian model, with no additional nonuniform costs and $F^0 = \hat{F}$, log-posterior maximization (MaP) already includes risk minimization (MiR). But a Bayesian MaP–MiR approximation is not restricted to a classical model. Assuming that the MaP step yields a good approximation of the Bayesian integral, performing a second independent risk minimization step after maximizing the posterior allows

1. states $f^0$ independently modeled from actions $\hat{f}$, that is $F^0 \neq \hat{F}$,

2. arbitrary costs $\tilde{l}^0$ independent of priors $L^0$,

3. an inverse setting with states defined by $p(y|q, f^0)$ and inverse actions $p(\hat{q}|y, \hat{f})$.

Now we discuss in more detail the case of Gaussian basis questions and states $f^0$ parameterized by their regression functions.

Let us now have a sharper look at the MaP approximation in the case of nonlocal questions. We calculate the functional derivatives separately for different classes of questions. We choose one–dimensional Gaussian distributed basis questions $X$

$$p(y|x, f^0) \propto e^{-\frac{1}{2}\left(\frac{y - \bar{y}_x}{\sigma_x}\right)^2}.$$

Thus in this section we assume the $f^0$ to be parameterized by their regression functions $\bar{y}_x = \bar{y}_x(f^0)$.[69] We get

$$\frac{d}{df_x^0} L(y|x', f^0) = \frac{d}{d\bar{y}_x} L(y|x', f^0) \propto \delta(x - x') \sigma_x^{-2}(y - \bar{y}_x).$$

For general local Gaussian questions including differing variances

$$p(y|q_x, f^0) \propto e^{-\frac{1}{2}\left(\frac{y - \bar{y}_{q_x}(\bar{y}_x)}{\sigma_{q_x}}\right)^2},$$

with $\bar{y}_{q_x}(\bar{y}_x)$ being a deterministic function of $\bar{y}_x$, we have

$$\frac{d}{d\bar{y}_x} L(y|q_{x'}, f^0) \propto \delta(x - x') \sigma_{q_x}^{-2}(y - \bar{y}_{q_x}) \frac{d}{d\bar{y}_x} \bar{y}_{q_x}(\bar{y}_x),$$

---

[68]There are justifications of a relation between complexity and priors also for the Bayesian framework (MacKay, 1992a). The idea is that complex states appear usually in much more variations than simple ones and if their total probability is fixed their individual probabilities become small. Such a coupling between prior and complexity results from using uniform priors for specifically selected groups of states, and a grouping may be seen as more or less natural for specific situations. The problem of defining uniform priors in a such hierarchical situation is similar to the situation for continuous variables where uniform priors do not remain uniform under general transformations.

[69]We do here not discuss the mathematical difficulties related to the definition of functional integrals (see for example Glimm & Jaffe, 1987). Finally, in numerical calculations we discretize all functions so integrals are replaced by sums. In the language of field theory we use a sharp ultraviolet cutoff. See for example, Bialek, Callan, Strong, 1996, for a discussion of the continuum limit in density estimation. Wahba, 1983, shows the paradox that the expectation of $L^0 = -(\lambda/2) \int dx \, (d^2\bar{y}/dx^2)^2$ under $p(\bar{y}) = e^{L^0}$ is infinite (See also Green & Silverman, 1994). The functional integral $\int d\bar{y} e^{-(1/2)\langle \bar{y}|\mathcal{O}|\bar{y}\rangle}$ alone is according to the formula for Gaussian integrals formally $(\det(\mathcal{O})/(2\pi))^{1/2}$. For $\mathcal{O} = (-(d^2\bar{y}/dx^2) + m^2)$ which possesses a continuous spectrum this determinate cannot be defined (see e.g. Roepsdorff, 1991). In field theory renormalization group methods are used to find meaningful continuum limits (There is a huge literatur about renormalization. See for example Zinn–Justin, 1989, Itzykson & Drouffe, 1989 Le Bellac, 1991, Fernández, Fröhlich, & Sokal, 1992, Binney, Dowrick, Fisher, & Newman, 1992, and references therein.)

and for nonlocal Gaussian questions like the usual smoothness questions

$$p(y|q, f^0) \propto e^{-\frac{1}{2}\left(\frac{y - \bar{y}_q[\bar{y}_x]}{\sigma_q}\right)^2},$$

with $\bar{y}_q[\bar{y}_x]$ denoting a deterministic functional depending on the set $\{\bar{y}_x\}$,

$$\frac{d}{d\bar{y}_x} L(y|q, f^0) \propto \sigma_q^{-2}(y - \bar{y}_q)\frac{d}{d\bar{y}_x}\bar{y}_q[\bar{y}_x].$$

For general questions we write for the log-posterior $L = \ln p$ and find

$$\frac{d}{d\bar{y}_x} L(y|q, f^0) = \frac{d}{d\bar{y}_x} \ln p(y|q, f^0) = \frac{\frac{d}{d\bar{y}_x} p(y|q, f^0)}{p(y|q, f^0)}.$$

Remembering $p(y|x, f^0) = \prod_i p(y_i|x_i, f^0)$, $\int dx' = \sum_d \int \prod_i^d dx_i'^{(d)}$, and $\int dy = \sum_d \int \prod_i^d dy_i^{(d)}$ the derivative of the probability is found as

$$\frac{d}{d\bar{y}_x} p(y^q|q, f^0) = \frac{d}{d\bar{y}_x} \int dx' \int dy\, p(y|x', f^0)p(y^q, x'|y, q)$$

$$= \int dx' \int dy \left(\sum_i^d \sigma_{x_i'^{(d)}}^{-2} \delta(x_i'^{(d)} - x)(y_i - \bar{y}_x)\right)$$
$$\times p(y|x', f^0)p(y^q, x'|y, q)$$

for Gaussian $p(y|x, f^0)$. Not only for questions without input noise but also for $d = 1$, including possible input noise, the $x'$–integration vanishes. In the latter case this gives

$$\int dy\, (y - \bar{y}_x)p(y|x, f^0)p(y^q|x, y, q).$$

The general stationarity conditions are obtained by setting the functional derivatives of the total log-posterior with respect to $\bar{y}_x$ which parameterize the $f^0$ to zero

$$\forall x\; :\; 0 = \frac{d}{d\bar{y}_x} \sum_i^n L(y_i|q_i, f^0),$$

where $n$ is the number of training data $D = \{(q_i, y_i)\}$. The $q_i$ can be general nonlocal questions, and may for example be written in terms of distances to templates $T_x$. For the sake of simplicity we will in this case use the term 'data' for discrete templates, i.e. those defined only for a discrete set of $x$, and will call templates defined for a continuous set $X$ shortly templates. Now, we look to some examples.

## 9.2 Local quadratic templates

We study a situation where the log-posterior is a sum of quadratic terms each of them with templates $T_x$ only depending on one $x$. The standard case of training examples consisting only of local Gaussian basis questions has

$$L^D = \sum_i^n L(y_i|x_i, f^0)$$

The corresponding stationarity conditions are for $\sigma_x = 1$

$$0 = \sum_i^n \delta(x - x_i)(y_{i,x} - \bar{y}_x) = \sum_i^{n_x} (y_{i,x} - \bar{y}_x)$$

with $n_x$ the number of times $x$ is in $D_q$, $y_{i,x}$ an answer to question $x$, and $\sum_i^0 = 0$ $x$ not in the data. For $x \in D_q$ we find for unrestricted $F^0$ the well–known mean square solution

$$\bar{y}_x = \frac{1}{n_x} \sum_i^{n_x} y_{i,x}.$$

The $\bar{y}_x$ for $x \notin D_q$ are arbitrary.

Including other local Gaussian questions

$$p(y^{q^x}|q^x, f^0) \propto e^{-\frac{1}{2}\frac{(y_{q_x} - \bar{y}_{q_x})^2}{\sigma_{q_x}^2}},$$

with $\bar{y}_{q_x} = \bar{y}_{q_x}(\bar{y}_x)$ a function of one $\bar{y}_x$ only, we find

$$0 = \sum_i^{n_x} \sigma_x^{-2}(y_{i,x} - \bar{y}_x)$$

$$+ \sum_{q_x} \sum_i^{n_{q_x}} \sigma_{q_x}^{-2}(y_{i,q_x} - \bar{y}_{q_x})\frac{d}{d\bar{y}_x}\bar{y}_{q_x},$$

equivalent to

$$\bar{y}_x = \frac{\sum_i^{n_x} \sigma_x^{-2} y_{i,x} + \sum_{q_x} \sum_i^{n_{q_x}} \sigma_{q_x}^{-2} y_{i,q_x} \frac{d}{d\bar{y}_x}\bar{y}_{q_x}}{n_x \sigma_x^{-2} + \sum_{q_x} n_{q_x} \sigma_{q_x}^{-2} \frac{\bar{y}_{q_x}}{\bar{y}_x} \frac{d}{d\bar{y}_x}\bar{y}_{q_x}}.$$

For linear $\bar{y}_{q_x} = a_{q_x}\bar{y}_x + b_{q_x}$ this is a linear equation with solution

$$\bar{y}_x = \frac{\sum_{q_x} a_{q_x} \sigma_{q_x}^{-2} \sum_i^{n_{q_x}} (y_{i,q_x} - b_{q_x})}{\sum_{q_x} a_{q_x}^2 n_{q_x} \sigma_{q_x}^{-2}}$$

where we simplified the formula by including the $x$ into the $q_x$ with $a_x = 1$ and $b_x = 0$. This gives the reweighting to be done for varying variances and scaling and the correction for varying bias.

In general, sums of local quadratic terms

$$\frac{1}{2} \sum_i \int dx\, \Lambda_x^i ||\bar{y}_x - T_x^i||^2 \qquad (22)$$

$$= \frac{1}{2} \sum_i <\bar{y} - T^i\, |\Lambda^i|\, \bar{y} - T^i>$$

$$= \sum_i \left(\frac{1}{2} <\bar{y}\, |\Lambda^i|\, \bar{y}> - <\bar{y}\, |\Lambda^i|\, T^i>\right) + c,$$

give linear stationarity equations. Here $c$ is an $f^0$–independent constant, $\Lambda^i$ are nonnegative diagonal matrices with matrix elements $\Lambda_{x,x'}^i = \delta(x - x')\Lambda_x^i \geq 0$, and we assumed real scalar products. We define the projector $\mathcal{P}^i = (\mathcal{P}^i)^2$ into the $T^i$–space spanned by the $x$ with nonzero $\Lambda_x^i$. Then $\Lambda^i$ commutes with this projector $\Lambda = \mathcal{P}^i\Lambda = \Lambda\mathcal{P}^i$ and without loss of generality the template is understood to be restricted to $T^i$–space meaning $T^i = \mathcal{P}^i T^i$.

Sum of terms for several $T^i$ can be combined. To formulate this in general we define

$$\mathcal{N} = \sum_i \mathcal{P}^i,$$

60

with diagonal elements $\mathcal{N}(x,x)$ giving the number of templates active for $x$, and

$$\Lambda = \sum_i \Lambda^i.$$

The total projector for $\sum_i T^i$ has two contributions

$$\mathcal{P} = \mathcal{N} + \mathcal{M},$$

with $\mathcal{M}$ depending on the pairwise overlaps of spaces defined by the $\mathcal{P}^i$ (compare the AND for probabilities). For the example of two templates this reads

$$\mathcal{P} = \mathcal{P}^1 + \mathcal{P}^2 - \mathcal{P}^1 \mathcal{P}^2.$$

The operators $\mathcal{N}$ and $\Lambda$ can be inverted in the space where its diagonal matrix elements $n_x$ are nonzero and we write for the inverse in that subspace

$$\mathcal{N}_{\mathcal{P}}^{-1} = \mathcal{P}\,(\mathcal{P}\mathcal{N}\mathcal{P})^{-1}\,\mathcal{P},$$

$$\Lambda_{\mathcal{P}}^{-1} = \mathcal{P}\,(\mathcal{P}\Lambda\mathcal{P})^{-1}\,\mathcal{P}.$$

Introducing the ($\Lambda$–weighted) sum

$$T^\Lambda = \sum_i \Lambda^i T^i,$$

(with special case $T = \sum_i T^i$) the mean (over $i$ not $x$) of a set of templates $T^i$ can be written as

$$\overline{T}^\Lambda = \mathcal{P}\overline{T}^\Lambda = \Lambda_{\mathcal{P}}^{-1} T,$$

(special case $\overline{T} = \mathcal{N}_{\mathcal{P}}^{-1} \sum_i T^i$) and we have for the sum

$$\sum_i < \bar{y} - T^i|\Lambda^i|\bar{y} - T^i >$$

$$=< \bar{y} - \overline{T}^\Lambda|\Lambda|\bar{y} - \overline{T}^\Lambda > + \sum_i < T^i|\Lambda^i|T^i > - < \overline{T}^\Lambda|\Lambda|\overline{T}^\Lambda > .$$

The difference

$$\sum_i^{n_i} < T^i|\Lambda^i|T^i > - < \overline{T}^\Lambda|\Lambda|\overline{T}^\Lambda >= n_i\, \mathrm{VAR}_\Lambda(\{T^i\}),$$

proportional to a of a variance is $\bar{y}$–independent and therefore irrelevant for the derivative if it only appears as additive term in $L$.

### 9.3   Linear regularization

Many smoothness and symmetry functionals are examples for Gaussian nonlocal questions. Take as example a functional with $L^0$ of the form

$$L^0(f^0) - c = -\frac{1}{2}\int dx \int dx'\, \bar{y}_x \mathcal{O}_{x,x'} \bar{y}_{x'} = -\frac{1}{2} < \bar{y}|\mathcal{O}|\bar{y} >,$$

$$(23)$$

with a real symmetric positive (semi–)definite $\mathcal{O}$ representing a Gaussian probability[70] and a constant $c$ ensuring correct normalization. As we use angle brackets

$< \cdot\,|\,\cdot >$ for a scalar product in a Hilbert space of functions $\bar{y}$, we denote $\mathcal{O}_{x,x'} = \mathcal{O}(x,x') =< x\,|\,\mathcal{O}\,|\,x' >$ also by angle brackets. This notation can be used for any hermitian linear operator $\mathcal{O}$. Here the term $\lambda\frac{1}{2} < \bar{y}|\mathcal{O}|\bar{y} >$ is a quadratic regularization term in the sense of Tikhonov, and if the spaces $F_c^0 = \{f^0|\lambda < \bar{y}(f^0)|\mathcal{O}|\bar{y}(f^0) >\le c$ are compact for real $c$ then in the limit $\lambda \to 0$ asymptotic stability conditions hold (Tikhonov, 1963, Vapnik, 1982). Here we are not especially interested in asymptotic results (except if necessary to ensure the validity of the saddle point approximation), but for continuous $x$ we will refer to the case with quadratic regularization functional and therefore linear stationarity equation as *linear regularization*.

A matrix element

$$< \mathcal{D}\bar{y}|\mathcal{D}\bar{y} >= ||\mathcal{D}\bar{y}||^2 =< \bar{y}|\mathcal{D}^\dagger\mathcal{D}|\bar{y} >,$$

with $\mathcal{D}^\dagger$ denoting the adjoint of $\mathcal{D}$, gives an operator $\mathcal{O} = \mathcal{D}^\dagger\mathcal{D}$ on the domain where the operators are defined. To construct a smoothness functional for square integrable functions, $\mathcal{D}$ can be chosen as a hermitian linear differential operator. A first order example is[71]

$$\mathcal{D}_{x,x'} = \mathcal{D}(x,x') = i\delta'(x - x') = -i\delta(x - x')\frac{d}{dx},$$

giving

$$\mathcal{O}_{x,x'} = \mathcal{O}(x,x') = -\delta''(x - x') = -\delta(x - x')\frac{d^2}{dx^2}.$$

Instead of giving a template for $\bar{y}_x$ we could give also templates for other functions of $\bar{y}_x$. For example, a template $T'$ for $\mathcal{D}\bar{y}$ can be written

$$L^0(f^0) = -\frac{1}{2} < \mathcal{D}\bar{y} - T'|\Lambda_{T'}|\mathcal{D}\bar{y} - T' >$$

$$= -\left(\frac{1}{2} < \bar{y}\,|\mathcal{D}^\dagger \Lambda_{T'}\mathcal{D}\,|\,\bar{y} > - < \bar{y}\,|\,\mathcal{D}^\dagger \Lambda_{T'} T' >\right) + c,$$

for real scalar products and $c$ an $f^0$–independent constant. With $\mathcal{O} = D^\dagger \Lambda^T D$ and $\tilde{T} = \mathcal{D}^\dagger \Lambda_{T'} T'$ this reads

$$L^0(f^0) = -\left(\frac{1}{2} < \bar{y}\,|\mathcal{O}|\,\bar{y} > - < \bar{y}|\tilde{T} >\right) + c.$$

We may express $T'$ also by $\mathcal{D}$ and write $T' = \mathcal{D}T$ so that

$$L^0(f^0) = -\frac{1}{2} < \mathcal{D}\bar{y} - \mathcal{D}T|\Lambda_{T'}|\mathcal{D}\bar{y} - \mathcal{D}T >$$

$$= -\frac{1}{2} < \bar{y} - T|\mathcal{O}|\bar{y} - T >$$

$$= -\frac{1}{2}||\bar{y} - T||_{\mathcal{O}}^2.$$

---

[70]Thus, $\mathcal{O}$ is an inverse covariance operator $\mathcal{C} = \mathcal{O}^{-1}$. As matrix elements of an inverse operator are also called Green's functions, $G(x,x') = \mathcal{C}(x,x') = \mathcal{O}^{-1}(x,x')$ are the Green's functions fulfilling $\mathcal{O}\mathcal{C} = \mathcal{I}$. If $\mathcal{O}$ has the form $\lambda\mathcal{I} - \mathcal{O}$ $G$ are the matrix elements (kernel) of the resolvent operator. Usually the resolvent is seen as function a complex $\lambda$ with poles at the eigenvalues and a cut at the continuous spectrum of $\mathcal{O}$. See for example chapter 7 in Glimm & Jaffe, 1987 and any book on functional analysis.

[71]Notice that this formal notation does not mean the operators are diagonal in the $x$–representation. The $\delta$–function only restricts the derivatives to the location $x = x'$, but the derivatives itself depend also on the neighborhood of $x$. This is most easily seen by replacing the derivatives with a finite difference approximation. The operator $\mathcal{D}$ to be hermitian on a function space requires boundary terms to vanish, as can be checked using partial integration. This is fulfilled e.g. for periodic functions on its periodicity interval or for functions which vanish asymptotically.

Examples include distances in Sobolev spaces, where $\mathcal{O}$ consists of a sum over derivatives. A log-posterior of the form (23), having no term linear in $\bar{y}_x$, corresponds to a 'null' template, and no inhomogeneities appear in the corresponding derivative.

The normalization constant for a $d$–dimensional $y$ in the presence of linear terms is calculated according to

$$\int \left( \prod_{j=1}^{d} dy^{(j)} \right) e^{-<\bar{y}|\mathcal{O}|\bar{y}>+<\bar{y}|T>}$$

$$= (2\pi)^{\frac{d}{2}} 2^d (\det \mathcal{O})^{-\frac{1}{2}} e^{<T|\mathcal{O}^{-1}|T>}.$$

We define analogously to the previous Subsection

$$\mathcal{O} = \sum_i \mathcal{O}^i,$$

and, with invertability in the space defined by a projector $\mathcal{P}$,

$$\mathcal{O}_{\mathcal{P}}^{-1} = \mathcal{P} (\mathcal{P}\mathcal{O}\mathcal{P})^{-1} \mathcal{P},$$

$$T^{\mathcal{O}} = \sum_i \mathcal{O}^i T^i,$$

and

$$\overline{T}^{\mathcal{O}} = \mathcal{O}_{\mathcal{P}}^{-1} T^{\mathcal{O}}.$$

The projected equation may be solved for example with the pseudo inverse. Also, $\mathcal{O}$ may be extended so it its inverse exists in the whole space. One may, for example, adding the identity on the zero space, a mass term $m^2 \mathcal{I}$, or impose boundary conditions.

Then, like in the local case, we have also for nondiagonal $\mathcal{O}^i$, assumed to be real symmetric positive (semi–) definite, for a sum of quadratic terms

$$\Delta_{\mathcal{O},T} = \sum_i <\bar{y} - T^i|\mathcal{O}^i|\bar{y} - T^i>$$

$$= <\bar{y} - \overline{T}^{\mathcal{O}}|\mathcal{O}|\bar{y} - \overline{T}^{\mathcal{O}}> + E_{\mathcal{O}}$$

$$+ \sum_i <T^i|\mathcal{O}^i|T^i> - <\overline{T}^{\mathcal{O}}|\mathcal{O}|\overline{T}^{\mathcal{O}}>,$$

with minimum ("ground state")

$$E_{\mathcal{O}} = \sum_i <T^i|\mathcal{O}^i|T^i> - <\overline{T}^{\mathcal{O}}|\mathcal{O}|\overline{T}^{\mathcal{O}}>$$

at

$$\bar{y}^* = \mathrm{argmin}_{\bar{y}} \Delta_{\mathcal{O},T} = \overline{T}^{\mathcal{O}},$$

in the space on which $\mathcal{P}$ projects. Thus, we can call $\overline{T}^{\mathcal{O}}$ the *template average* of the set of $T_i$ with respect to the norm induced by the $\mathcal{O}^i$. The standard average is a special case. Like for a standard mean also a template average of two templates is always in the 'middle' between the two. That means, $\overline{T}^{\mathcal{O}}$ has has equal $\mathcal{O}$–distance from both templates

$$<\overline{T}^{\mathcal{O}} - T_1|\mathcal{O}|\overline{T}^{\mathcal{O}} - T_1> = <\overline{T}^{\mathcal{O}} - T_2|\mathcal{O}|\overline{T}^{\mathcal{O}} - T_2>.$$

This is easily seen because for $\mathcal{O} = \mathcal{O}^1 + \mathcal{O}^2$ one has

$$\overline{T}^{\mathcal{O}} - T_1 = \mathcal{O}_{\mathcal{P}}^{-1} \mathcal{O}^2 (T_1 - T_2) = -\left( \overline{T}^{\mathcal{O}} - T_2 \right),$$

with the minus sign disappearing in a quadratic form.

The functional derivative for the sum of a nonlocal quadratic $L^0$ as in (23) and a local quadratic term is

$$\frac{dL(\bar{y})}{d\bar{y}_x} = -\int dx' \mathcal{O}_{x,x'} \bar{y}_{x'} - \Lambda_x^T (\bar{y}_x - T_x)$$

$$= -<x|\mathcal{O}|\bar{y}> - <x|\Lambda^T|\bar{y}> + <x|\Lambda^T|T>.$$

Mean square error terms are special examples of such terms and therefore encompassed by this formulation. The variable $x$ can be multi–dimensional, assuming the vector is written in a basis where $\Lambda^T$ is diagonal. (The situation where for one $x$ several different $T_x$ are available will be discussed below.) For invertible $\Lambda^T + \mathcal{O}$ the stationarity condition reads

$$\bar{y} = (\Lambda^T + \mathcal{O})^{-1} \Lambda^T T, \tag{24}$$

which is for a linear operator $\mathcal{O}$ a linear inhomogeneous equation. If not invertible in the full space components of the null space, i.e. solutions of $(\Lambda^T + \mathcal{O})\bar{y} = 0$ can be added to a special solution. In cases $\mathcal{O}^{-1}$ can be calculated, it can be useful to rewrite Eq.(24) by separating the parts invariant under projection with $\mathcal{P}^T$ from those which are not

$$\mathcal{O}\bar{y} = \Lambda^T (T - \bar{y}).$$

Then the vector

$$a = \Lambda^T (T - \bar{y}) = \mathcal{P}^T a$$

being invariant under projection $\mathcal{P}^T$ can be calculated purely within the $T$–space. If $\mathcal{O}$ is invertible, the unknown $\bar{y}$ in the definition of $a$ can be eliminated using $\bar{y} = \mathcal{O}^{-1} a$, giving for $a$ the equation

$$\left( \mathcal{I} + \Lambda^T \mathcal{O}^{-1} \right) a = \Lambda^T T, \tag{25}$$

with $\mathcal{I}$ denoting the identity. Using $\mathcal{P}^T \Lambda^T = \Lambda^T \mathcal{P}^T$, $a = \mathcal{P}^T a$ and $T = \mathcal{P}^T T$ to insert the projector $\mathcal{P}^T$ gives

$$\left( \left( \mathcal{P}^T \Lambda^T \mathcal{P}^T \right)^{-1} + \left( \mathcal{P}^T \mathcal{O}^{-1} \mathcal{P}^T \right) \right) a = \mathcal{P}^T T = T.$$

Hence, only matrix elements of $(\mathcal{P}^T \mathcal{O}^{-1} \mathcal{P}^T)$ within the $T$–space are required to solve for $a$. We now consider in more detail the cases zero and nonzero quadratic templates.

### 9.3.1 Homogeneous linear regularization

Here we consider the special, but most common case of discrete Gaussian data with equal variance and quadratic nonlocal terms, e.g. smoothness, which can be seen as corresponding to a a zero–template $T_x = 0$, $\forall x$. Gaussian data terms with $\sigma_x^2 = \lambda$ give

$$\frac{d(L^D + L^0)}{d\bar{y}_x} \propto \sum_i^{n_x} (y_{i,x} - \bar{y}_x) - \lambda \int dx' \mathcal{O}_{x,x'} \bar{y}_{x'}.$$

The stationarity condition reads

$$(\mathcal{N}^D + \lambda \mathcal{O})\bar{y} = D, \tag{26}$$

Vector $D$ has components

$$D_x = \sum_i^n \delta(x - x_i) y_{i,x} = \sum_j^{n_x} y_{j,x},$$

and the operator $\mathcal{N}^D$ has matrix elements

$$\mathcal{N}^D_{x,x'} = \mathcal{N}^D(x,x')$$

$$= \delta(x-x')\sum_i^n \delta(x-x_i) = \delta(x-x')n_x,$$

giving the number of how often a specific $x$ appears in the data. We define the projector $\mathcal{P}^D$, projecting into the space spanned by the $x$ included in the data, by its matrix elements

$$\mathcal{P}^D_{x,x'} = \mathcal{P}^D(x,x') = \delta(x-x')\Theta\left(\sum_i^n \delta(x-x_i)\right),$$

with the step function $\Theta(x)$ restricting the matrix elements to zero or one. Its number of nonzero diagonal elements, i.e. $\tilde{n} = \mathrm{Tr}\mathcal{P}^D = n - \sum_x \sum_{i=2}^{n_x} 1$ is the number of different $x = x_i$ in the data. Then $D = \mathcal{P}^D D$ and $\mathcal{N}^D = \mathcal{P}^D \mathcal{N}^D \mathcal{P}^D$, with the operator $\mathcal{N}^D$ being equal to $\mathcal{P}^D$, and therefore an identity in that subspace, if $n_x = 1$ for all $x$. This is the usual case when i.i.d. sampling for continuous $x$. In a space where $(\mathcal{N}^D + \lambda\mathcal{O})$ is invertible the linear, inhomogeneous (e.g. integro–differential) equation (26) has the solution

$$\bar{y} = (\mathcal{N}^D + \lambda\mathcal{O})^{-1}D.$$

being a special case of Eq.(24) with $D/\lambda = \Lambda^T T$ and $\mathcal{N}^D/\lambda = \Lambda^T$. Components of a null space may be added. The matrix elements $\mathcal{O}^{-1}_{x,x'} = G(x,x')$ (or Green's function) satisfy by definition

$$\mathcal{O}G(x,x') = \delta(x-x').$$

For some $\mathcal{O}$ the Green's function can be calculated analytically. Then the solution of the resulting equation

$$\bar{y} = \frac{1}{\lambda}\mathcal{O}^{-1}\left(D - \mathcal{N}^D\bar{y}\right)$$

$$= \frac{1}{\lambda}\mathcal{O}^{-1}\mathcal{N}^D\left(\bar{D} - \bar{y}\right), \qquad (27)$$

(with $\bar{D} = \left(\mathcal{N}^D_{\mathcal{P}^D}\right)^{-1}D$) or in components

$$\bar{y}_x = \sum_i^n G(x,x_i)\frac{y_{i,x_i} - \bar{y}_{x_i}}{\lambda} = \sum_i^{\tilde{n}} a_i G(x,x_i),$$

is for fixed $x$ in a $\tilde{n}$–dimensional space spanned by known $G(x,x_i)$ with different $x_i$. In the vector

$$a = \frac{(D - \mathcal{N}^D\bar{y})}{\lambda}$$

with components

$$a_x = \sum_j^{n_x}\frac{y_{j,x} - \bar{y}_x}{\lambda}$$

only the components $a_{x_i} = a_i$ for $x_i$ belonging to the data are not equal to zero. Inserting $\bar{y} = \mathcal{O}^{-1}a$ into the definition of $a$ gives a $\tilde{n}$–dimensional matrix equation

$$\sum_i^{n_{x_i}} y_{i,x_i} = n_{x_i}\sum_j^d G(x_i,x_j)a_j + \lambda a_i,$$

or in operator notation

$$D = (\mathcal{N}^D(\mathcal{P}^D\mathcal{O}^{-1}\mathcal{P}^D) + \lambda\mathcal{I})a. \qquad (28)$$

The identity $\mathcal{I}$ commutes with $\mathcal{P}^D$ and $\mathcal{P}^D\mathcal{O}^{-1}\mathcal{P}^D$ has the matrix elements $G(x_i,x_j)$. This is the equivalent of Eq.(25).

If $\mathcal{O}$ has zero modes, then Eq.(27) becomes

$$\bar{y} = \mathcal{O}_1^{-1}(D - \mathcal{N}^D\bar{y}) + \sum_k^m b_k u_k$$

where $u_k$ represents an orthonormal basis of the zero space and $\mathcal{O}_1$ denotes the restriction of $\mathcal{O}$ to the sub–space where its inverse exists. In addition one has in the space of zero modes

$$0 = \mathcal{P}^0(D - \mathcal{N}^D\bar{y}),$$

where $\mathcal{P}^0$ denotes the projector into the space of zero modes, i.e.

$$\mathcal{P}^0(x,x') = \sum_j^m u_j(x)u_j(x').$$

This yields the two data space equations

$$\sum_i^{n_{x_i}} y_{i,x_i} = n_{x_i}\sum_j^d G(x_i,x_j)a_j + \lambda a_i, + \sum_k^m b_k u_k(x_i)$$

$$\sum_i u_k(x_i)a_i = 0, \quad \forall k.$$

The (pseudo-differential) operator

$$\mathcal{O} = \sum_{m=0}^{\infty} \mathcal{O}_m = \sum_{m=0}^{\infty} (-1)^m\frac{\sigma^{2m}}{m!2^m}\nabla^{2m}, \qquad (29)$$

with $\nabla^{2m}$ denoting the m–iterated Laplacian, results in a Gaussian $G(x,x')$. Diagonal $G(x,x')$ correspond to local questions, radially symmetric (Gaussian) Green's functions are called Radial Basis Functions. This and more examples, including various forms of splines, as well as the relation to conditionally positive definite and completely monotonic functions can be found in Poggio & Girosi, 1990 Wahba, 1990, and Girosi, Jones, & Poggio, 1995.

Restricting the terms in the summation to a number smaller than $\tilde{n}$, corresponding to an additional prior or cost term, can be combined with an algorithm to determine the optimal selection of $x_j$ (e.g. the centers of Gaussians). Including for example the variance $\sigma_x$ for the regularizer as parameter leads to nonlinear equations.

### 9.3.2  Inhomogeneous linear regularization

Choosing Gaussians data terms with equal variance in addition to a template $T_x$ we have to minimize

$$\sum(y_{i,x} - \bar{y}_{x_i})^2 + \lambda_T\int dx\,(\bar{y}_x - T_x)^2,$$

or for $h_x = \bar{y}_x - T_x$

$$\sum(y_i - h_{x_i} - T_{x_i})^2 + \lambda_T\int dx\,(h_x)^2,$$

with a shifted local error term and a penalty for deviation from zero. For the unshifted parameterization the stationarity condition reads

$$\bar{y}_x = \frac{\sum_i^{n_x} y_{i,x} + \lambda_T T_x}{n_x + \lambda_T}$$

which means $\bar{y}_x = T_x$ for non–data points $x \neq x_i, i = 1, \cdots, n$.

The last example becomes more interesting if combined with a nonlocal term, for example a differential operator implementing a smoothness prior. Then one has to minimize

$$\sum (y_i - \bar{y}_{x_i})^2 + \lambda_T \int dx\,(\bar{y}_x - T_x)^2$$

$$+\lambda_S \int dx \int dx'\,(\bar{y}_x \mathcal{O}_{x,x'} \bar{y}_{x'})^2,$$

with stationarity equation

$$\int dx'\left((n_x + \lambda_T) + \mathcal{O}_{x,x'}\right)\bar{y}_{x'} = \sum_i^{n_x} y_i + \lambda_T T_x.$$

For the example $\mathcal{O}_{x,x'} = \delta(x-x')\frac{d^2}{dx^2}$ this gives the linear inhomogeneous differential equation

$$\left(\frac{d^2}{dx^2} + n_x + \lambda_T\right)\bar{y}_x = \sum_i^{n_x} y_i + \lambda_T T_x.$$

We see, that in cases of nonzero templates the stationarity equations become besides the always present $\delta$–like data terms also continuous inhomogeneities. For continuous $x$ we will call the case where the regularization functional is a sum of a term quadratic in $\bar{y}$ and a term linear in $\bar{y}$ *inhomogeneous linear regularization*.[72]

Expressing for two templates that functions $f^0$ should be similar to $T_1$ OR $T_2$ leads to nonlinear equations, which we discuss in the next paragraph.

## 9.4 Nonlinear regularization

Priors which are constructed by combining quadratic subproperties $C_i$ using a real valued extension of logic do not need to be quadratic in the $\bar{y}_x$, for example, if an OR is implemented in a soft and not in a hard version. Also one could use a parameterization of $F^0$ so $\bar{y}_x$ is a nonlinear function of the parameters,[73] or use a template for nonlinear questions, like a correlation template for $\bar{y}_x \bar{y}_{x'}$ in terms like

$$||\bar{y}_x \bar{y}_{x'} - T_{x,x'}||^2.$$

Let us consider a case with two templates combined by a soft OR. The two templates could have been obtained

---

[72] Also for homogeneous linear regularization the stationarity equations are inhomogeneous, but the regularization functional adds nothing to the data inhomogeneities of $\delta$–form.

[73] Many methods, like for example sigmoidal neural networks use a nonlinear parameterization, but not independent for each $\bar{y}_x(f^0)$. One may contrast a genuine nonlinearity corresponding to a model of nature (i.e. $F^0$, $L$) and a non-linearity induced by choosing a nonlinear action model ($\hat{F}$, $l$).

by using the human interface discussed above and represent two typical functions $T^1$ and $T^2$ to which the actual $\bar{y}_x(f^0)$ is expected to be similar to at least one of them. For example, two such templates can be prototypical structures for electrocardiograms, or patterns in financial time series. This situation might, for example, be naively approximated according to Eqs.(7) by including terms like

$$\lambda_T \left( \int dx\,(\bar{y}_x - T_x^1)^2 \right) \left( \int dx'\,(\bar{y}_{x'} - T_{x'}^2)^2 \right),$$

assuming a possible bound $m$ incorporated into $\lambda_T$. This non–quadratic regularization term has the functional derivative with respect to $\bar{y}_x(f^0)$

$$\lambda_T \left( (\bar{y}_x - T_x^1) \int dx\,(\bar{y}_x - T_x^2)^2 \right.$$

$$\left. +(\bar{y}_x - T_x^2) \int dx\,(\bar{y}_x - T_x^1)^2 \right),$$

yielding a nonlinear stationarity equation. We may therefore call this case *nonlinear regularization*, which in case of nonzero templates is also inhomogeneous.

### 9.4.1 Finite temperature or mixture regularization

We can choose the realization (8) for OR, if we assume that the properties we combine are log-probabilities,

$$L = L^M = \ln p^M = \ln \sum_i^{N_i} \frac{Z_i}{Z}$$

$$\ln\left( \sum_i^{N_i} e^{-\frac{\beta}{2}\left(\sum_{j_i}^{N_{j_i}} <\bar{y}-T^{ij_i}|\mathcal{O}^{ij_i}|\bar{y}-T^{ij_i}>\right)+c_i-\ln Z} \right),$$

$$(30)$$

where the constants $c_i$ can include the logarithm of inverse normalization factors and weights for component $i$ and

$$Z_i = e^{-\frac{\beta}{2}\left(\sum_{j_i}^{N_{j_i}} <\bar{y}-T^{ij_i}|\mathcal{O}^{ij_i}|\bar{y}-T^{ij_i}>\right)+c_i} \quad (31)$$

and $Z = \sum_i Z_i$. This form is general enough to include the product terms of OR for non–disjunct events if $c_i$ is allowed to be imaginary, giving a negative factor for the exponential. For disjunct events the $c_i$ are real, thus all contributions to the sum are positive and we speak of a *mixture model* or for continuous $x$ of *mixture regularization*. For mixture models see Everitt & Hand, 1981, Titterington, Smith, & Makov, 1985, Kontkanen, Myllymäki, & Tirri, 1997. For every fixed $\bar{y}$ the log–probability (30) has the structure of a free energy of a system at temperature $1/\beta$. Thus, emphasizing the temperature–dependence we will also speak of a thermic realization of the OR and, accordingly, a *finite temperature regularization*. In Section 5.3.4 we defined energies as ($\beta$–scaled) shifted log-probabilities of elementary events. Thus, with respect to the set of (disjunct) elementary events $i = \omega_i \in \Omega$, the exponents define energies,

$$\beta E_i = -\frac{\beta}{2}\left( \sum_{j_i}^{N_{j_i}} <\bar{y}-T^{ij_i}|\mathcal{O}^{ij_i}|\bar{y}-T^{ij_i}> \right) + c_i$$

which are unique up to a factor $\beta$ (inverse temperature) and a constant. Here $i$ represent possible, disjunct 'states', and the function $\bar{y}$ plays the role of system parameters which we can adapt to minimize the free energy. Thus, the system is in state $i_1$ OR $i_2$ OR $\cdots$. The variables $(j_i, x)$ label the subsystems (e.g. internal or microscopic degrees of freedom, like single particle coordinates and momenta in a many particle system). Every state or elementary event $i$ is a complete collection of states for all subsystems, labeled by $(j_i, x)$, i.e. subsystem $x_1$ is in state $a_1$ (e.g. $\bar{y}_{x_1} = a$) AND subsystem $x_2$ is in state $a_2$ AND $\cdots$. The special form 30 has quadratic energies ('generalized oscillators' or 'generalized free fields')[74] and is therefore for finite $|X|$ a Gaussian mixture model, or, for continuous $X$, a mixture of Gaussian processes. $\mathcal{O}^{iji}$ defines the inverse covariance matrices of the processes.

We may remark that in (30) $p(f^0|f)$ is written as mixture while in applications in density estimation a mixture model is often used for $p(y|f^0)$. In regularization every training data point $(y_i, x_i)$ gives only one $x$ component of a whole vector $\bar{y}(x)$. In a density estimation problem a complete data vector with all its $x$ components are given and $d = |X|$ is discrete and usually relatively small, and $x$ is denoted by a discrete index like for example $i$ or $k$. In regularization a ('data') vector corresponds to a function $\bar{y}(x)$ given $\forall x$. For continuous $x$ this is a realization or complete sample path of a stochastic process. In addition to the finite number of training data, the mixture components are determined by "continuous data", i.e. templates $T^{iji}(x)$, and corresponding distances. To enable generalization at least one template has to be present, for example constant and equal to zero (zero template), together with a smoothness related $\mathcal{O}$. The problem of constructing the prior is, for example, discussed in Section 5.2. At this step we assume the mixture model for the prior fixed, except possibly for a few remaining parameters, like $\beta$, which can be adapted by cross-validation or after defining a corresponding prior by explicit Bayesian integration. Also, our interest is not restricted in identifying the maximal activated mixture component, like for example in deterministic clustering, which can be seen as low temperature approximation. Our problem is finding the maximal probable $f^{0,*}$, which does not necessarily have to coincide with the center of one of the single mixture components, given only (a small, finite) part of its vector components (not mixture components) $y_i(x_i)$.

Especially interesting is the limiting case of very large (or infinite) dimension $|X|$ of the space $X$ with only part of the components $x$ given. To allow a useful degree of generalization in this situation we must enforce strong correlations between different dimensions $x$. This can be done, like in the example of smoothness, by choosing special (e.g. metric, differentiable) structures on the set $X$ of dimensions $x$. An example of a finite version

of smoothness for the regression function is $\left(\frac{\bar{y}_i - \bar{y}_{i+1}}{i - i + 1}\right)^2$, writing $i$ for the dimension index $x$. This restricts the centers (means) to the neighborhood of the diagonals of adjacent dimensions.

Similar to a mixture model for function approximation is the clustering algorithm of (Rose, Gurewitz, & Fox, 1990). In clustering one tries to find for a given set of points $x_i$ ($y$ in our notation) a corresponding cluster centroid $j$ ($f^0$ in our notation) with respect to an error function $E(x, j)$ ($L(y, f^0)$ in our notation, or more precisely, an approximation problem with $\tilde{l}(y, \hat{f})$ and identification of $\hat{F}$ with $F^0$). Like in our case the temperatur parametrizes the convexity/concavity of the error surface, which in the case of clustering regulates the number of distinct centroids. The $\beta$ parameter is the Lagrange parameter in a maximal entropy approach determining the average error. However, we do not assume the average error to be fixed in advance and like often in statistical physics, the temperature or $\beta$ itself is the more natural parameter, even if uniquely related to the average error. For an application in pairwise clustering see Hofmann and Buhmann, 1997, and references therein. In contrast to the mixture model we are considering here, the error or log–probability of this problem is not a "free" Gaussian model and the method is combined with an additional mean–field approximation. For the use of a temperature parameter in optimization and matching problems see Yuille, 1990, Yuille, Kosowsky,1994, Yuille, Stolorz, Utans, 1994, and for simulated annealing, for example, Aarts-Korts-1989.

The sum over $i$ in the mixture model is analogous to the integration over $f^0$ in the full risk. Indeed, replacing to obtain a more symmetric notation $\int df^0$ by $\int df_1^0$ and $\sum_i$ by $\int df_2^0$, we can write for $r(f, \hat{f})$ for an approximation problem with $\hat{F} = F_1^0$

$$r(f, \hat{f}) = -\int df_1^0 \int df_2^0 \int dy\, p(f_1^0, f_2^0) p^D(y|f_1^0, f_2^0) \ln p^D(y|\hat{f})$$

$$= -\int df_1^0 \int dy\, p(f_1^0) \int df_2^0\, p(f_2^0|f_1^0) p^D(y|f_1^0, f_2^0) \ln p^D(y|\hat{f})$$

$$= -\int df_1^0 \int dy\, p(f_1^0) p^D(y|f_1^0) \ln p^D(y|\hat{f}).$$

If the $f_2^0$–integral (or summation over $i$) is performed exactly (or a saddle point approximation with multiple sadddle points is used) then

$$p(f_1^0) = \int df_2^0\, p(f_1^0, f_2^0),$$

is a mixture of components of the form of $p(f_1^0, f_2^0)$ i.e. for log–probabilities

$$L(f_1^0) = \ln \int df_2^0\, e^{L(f_1^0, f_2^0)},$$

An example of such an integration is given by the mixture–like elastic net energy (Durbin, Willshaw, 1987)

$$E_{mix}(y) = -\frac{1}{\beta} \sum_k \ln \sum_j^M e^{-\beta|x_k - y_j|^2} + \gamma \sum_i^N |y_i - y_{i+1}|^2,$$

---

[74]Which shall mean they are quadratic forms, which however may include $\delta$–like forces (data), potentially higher order derivatives or nonlocal (e.g. in time if $x$ corresponds to the time variable) terms, as well as linear terms in coordinates and derivatives (e.g. friction).

with given vectors $x_j$ and $N \geq M$ vectors $y_i$ to be optimized with $y_{N+1} = y_1$. This energy $E_{mix}$ can be obtained by summing over binary variables $V_{kj}$ as $-(1/\beta) \ln \sum_V e^{-\beta E}$ with

$$E(V, y) = \sum_{kj} V_{kj} + \gamma \sum_i |y_i - y_{i+1}|^2,$$

under the restriction $\sum_{kj} V_{kj} = 1$, $\forall k$ (Yuille, 1990, Yuille, Stolorz, Utans, 1994).

The form (30) uses one level of Gaussian mixtures. The one level structure is in principle no restriction, as every logical formula can be written in either conjunctive or disjunctive normal form. Those may however be very lengthy and contain negations, which one wants to avoid in continuous cases. Thus, a hierarchical model may be much more economical. Form (30) also uses with Gaussian processes the simplest possibility for a single mixture components. These can be seen as first term of a Taylor expansion of general more general mixture components. Correlation templates are higher order terms

$$< \bar{y} \otimes \bar{y} - T^{(2)} | \mathcal{O}^{(2)} | \bar{y} \otimes \bar{y} - T^{(2)} >,$$

Here, $\bar{y} \otimes \bar{y}$ are matrices with operator

$$\Delta^{(2)} = \mathcal{O}^{(2)} = \sum_{k,l} \mathcal{O}_k \otimes \mathcal{O}_l.$$

acting on the enlarged vector space of those matrices. Analogously one may consider higher order terms

$$\Delta^{(n)} = < \bigotimes_i^n \bar{y} - T^{(n)} | \mathcal{O}^{(n)} | \bigotimes_i^n \bar{y} - T^{(n)} > .$$

The natural choice for templates is

$$T^{(n)} = \bigotimes_i^n T.$$

For $T^{(n)} = \bigotimes_i^n T_i$ and $\mathcal{O}^{(n)} = \bigotimes_i^n \mathcal{O}_i$ the $\Delta^{(n)}$ factorize

$$\prod_i^n < \bar{y} - T_i | \mathcal{O}_i | \bar{y} - T_i > .$$

For linear $\mathcal{O}_i$ this non—quadratic interaction term (i.e. non–Gaussian probability) has minima at every $\bar{y}^{*,i} = T_i$. Thus, for different $T_i$ already one mixing component, or the corresponding energy, creates multi-modal, i.e. nonconvex or OR–like, functions. Those multi-model (interaction) terms can arise by integrating out (hidden or latent) variables, e.g. microscopic degrees of freedom. This integration or summation (over probabilities $p = e^{-\beta E_i}/Z$ and not over energies $E_i$) is a realization of OR and results in a mixture model (with maybe an infinite number of components). Such an 'effective' energy $E$, represents from the point of view of the finer system a free energy $F$.[75] Most often, but not always as we shall see, effective energies are used in the range where they are approximately $\beta$–independent.

---

[75]Instead of describing effective energies as the result of marginalizing, one can see the introduction of additional degree of freedoms as an improvement of an existing theory. which 'explain' a certain features of an older, i.e. then effective, energy function. Such additional degrees of freedom are introduced, for example, when finding new particles in physics, or better explanations for diseases in medicine.

## 9.5 Interactions and Landau–Ginzburg regularization

Interaction terms introduced at the end of the previous section are another possibility to write nonconvex OR–like log-probabilities and they provide finally the connection with our discussion a fuzzy implementation of prior knowledge in Section 5. There we pointed out, that one may use fuzzy properties not only for probabilities, but also possibly directly for log-probabilities. For such an *interaction regularization* we have, to specify an *interaction model*, like we have to specify the energy function (Hamiltonian) for a specific system in physics. Indeed, this specification of interactions for a physical model also provides an example that it is sometimes more natural to specify directly log-likelihoods (energies).

Usually, those are taken as polynomial functions

$$L = L^I = -g \left( \frac{1}{2} \sum_i \prod_{j_i} < \bar{y} - T^{ij_i} | \mathcal{O}^{ij_i} | \bar{y} - T^{ij_i} > \right).$$
(32)

In a mixture regularization with model (30) for the free energy the sum over $i$ is an implementation of OR for log-probabilities. For the example (32), the product over $j_i$ can be interpreted as some fuzzy OR according to Eqs.(7) applied to properties defining the log-probability. If the model consists of one term one could speak of a 'zero temperature regularization with interactions'. On the other hand the interactions might be effective, i.e. chosen to approximate a more fundamental mixture model. The interactions have with respect to the underlying model $\beta$, i.e. also temperature, dependent parameters. In this case it is better to speak of an 'effective interaction regularization'. Effective, i.e. $\beta$–dependent, interactions may be resulting from a Taylor expansion of the a 'true' underlying free energy or log-probability. In contrast to a high energy expansion, where $L$ is expanded around $1/\beta = \infty$ one may also expand around a not infinite $1/\beta = 1/\beta^*$. We will call the case where we introduce a temperature–like parameter $1/\beta$ (or reduced temperature $t = (1/\beta) - (1/\beta^*)$) in the energy in analogy to the celebrated phenomenological treatment of phase transitions in physics a *Landau–Ginzburg regularization*. (See Landau & Lifshitz, 1980, (§145) and for example Goldenfeld 1992; Safran 1994; Ivanchenko & Lisyansky, 1995.) The Landau–Ginzburg theory is used to model systems near a phase transition and many results at the critical temperature are independent from many details of the system (Universality classes with respect to critical exponents). These results, however, are not necessarily valid at the phase transition when fluctuations are important. This is , for example, the case in low dimensional systems with local, i.e. short range, forces. We are not only interested in an effective $L$ in the immediate neighborhood of phase transitions, but it is this neighborhood where most problems can arise and where a nonlinear regularization is most different from a linear approach.

Clearly, a mixture model (32) and interaction model (30) can be quantitatively quite different, they may however share common qualitative features. Probabilities

according to a mixture model can easily be generated in a two step process (first choose $i$ according to the mixture coefficients, then draw from $i$) if the probability processes for the components $i$ are available. In an interaction model a decomposition in simple Gaussian processes is not necessarily possible. The form (32) has however always a polynomial structure, which can be helpful for calculational purposes. For fourth order polynomials the solutions can always be given explicitly, so the two template case is analytically solvable.[76] Notice, that the usual mean square data terms are encompassed in both formulations, as operators $\mathcal{O}^{ij_i}$ diagonal in $x$-representation and with only a finite number of nonzero elements.

## 9.6 Mean field equations

Now we get to the problem of solving these nonlinear stationarity or mean field equations. We already discussed that a maximum does not change under strictly monotonically increasing functions $h(L)$, i.e. functions with $L(f^0) > L(f'^0) \Rightarrow h(L(f^0)) > h(L(f'^0))$ and $\frac{dh}{dL} > 0$. Strictly monotonically decreasing functions $h(L)$ only change a maximum into a minimum. Including such functions $h$ the stationarity conditions, obtained by setting the functional derivative to zero, read

$$0 = \frac{dh(L)}{dL}\frac{dL(f^0)}{df^0} = \frac{dh(L)}{dL}\frac{dL(\bar{y})}{d\bar{y}},$$

where in our case $\bar{y}$ represents the parameter vector $f^0$, and we will see that iteration procedures can be related to different $h$.

We will now give the mean field equations, for finite temperature and (effective) polynomial interaction regularization, in the form

$$\mathbf{O}\bar{y} = \mathbf{t}, \qquad (33)$$

with in general $\bar{y}$–dependent $\mathbf{O} = \mathbf{O}(\bar{y})$ and $\mathbf{t} = \mathbf{t}(\bar{y})$, so the equation is nonlinear.

### 9.6.1 Mean field equations for finite temperature regularization

For a log-probability of mixture form (30) we find

$$\mathbf{O}^M = \overline{\mathcal{O}}^Z \quad = \quad \frac{\sum_i Z_i \mathcal{O}^i}{Z} = \frac{\mathcal{O}^Z}{Z}, \qquad (34)$$

$$\mathbf{t}^M = \overline{T}^{Z,\mathcal{O}} \quad = \quad \frac{\sum_i Z_i T^i}{Z} = \frac{T^{Z,\mathcal{O}}}{Z},$$

with according to (31) $\bar{y}$–dependent

$$Z_i = Z_i(\bar{y}) = e^{-\frac{\beta}{2}\left(\sum_{k_i}<\bar{y}-T^{ik_i}|\mathcal{O}^{ik_i}|\bar{y}-T^{ik_i}>\right)+c_i},$$

$$Z = Z(\bar{y}) = \sum_i Z_i(\bar{y}),$$

and $\bar{y}$–independent

$$\mathcal{O}^i = \sum_{j_i} \mathcal{O}^{ij_i},$$

---

[76] Already 'solvable' nonlinear polynomial equations require iteration methods: most roots have to be calculated iteratively. There is no doubt, however, that this can be done quite efficiently.

$$T^i = T^{\mathcal{O}^i} = \sum_{j_i} \mathcal{O}^{ij_i} T^{ij_i} = \mathcal{O}^i \overline{T}^{\mathcal{O}^i}.$$

Thus, denoting by $< \cdot >_{(Z_i/Z)(\bar{y})}$ the expectation under the probability $Z_i/Z$ the stationarity equation can be written

$$0 = \frac{\sum_i Z_i \left(\mathcal{O}^i \bar{y} - T^i\right)}{Z} = <O^i \bar{y} - T^i>_{(Z_i/Z)(\bar{y})} . \quad (35)$$

Using the monotonic transformation $h(L) = e^L = p$, or, equivalently, multiplying the mean field equation, by $Z$, gives

$$\mathcal{O}^Z \bar{y} = T^{Z,\mathcal{O}}.$$

For only one template $T^{ij_i} = T$, i.e. $\mathbf{t}^M = \mathcal{O}T$ a trivial solution is always $\bar{y} = T$. The stationarity equation can be solved in the space where the inverse of $\mathcal{O}^Z$ exists, using for example the pseudo inverse. Alternatively, $\mathcal{O}^Z$ may be extended to be invertible, for example by adding the identity on the zero space, a mass term proportional to the identity operator, or by imposing boundary conditions.

In the high temperature limit $\beta \to 0$ and $c_i = c$, $\forall i$ all $Z_i$ become equal so that

$$\bar{y} = (\mathcal{O}_{\mathcal{P}}^Z)^{-1}T^{Z,\mathcal{O}} \Rightarrow \bar{y} = \mathcal{O}_{\mathcal{P}}^{-1}T^{\mathcal{O}} = \overline{T}^{\mathcal{O}}.$$

Hence, we find the template average $\overline{T}^{\mathcal{O}}$ as high temperature limit of the mean field solution for the mixture model.

In the low temperature limit $\beta \to \infty$ only the largest $Z_i$ survive(s), so that all $T^{\mathcal{O}^i}$ with ("positive error gap")

$$\left(\sum_{k_i} <T^{\mathcal{O}^i} - T^{ik_i}|\mathcal{O}^{ik_i}|T^{\mathcal{O}^i} - T^{ik_i}>\right) - \frac{2c_i}{\beta}$$

$$< \left(\sum_{k_{i'}} <T^{\mathcal{O}^i} - T^{i'k_{i'}}|\mathcal{O}^{i'k_{i'}}|T^{\mathcal{O}^i} - T^{i'k_{i'}}>\right) - \frac{2c_{i'}}{\beta},$$

for all $i' \neq i$ in the limit $\beta \to \infty$, become a low temperature solution.

Using the same (invertible) operator for all mixture components so that $\mathcal{O}^i = \mathcal{O}/N_i$ results in the equation

$$\bar{y} = \overline{T}^{\mathcal{O}^i,Z} = \sum_i \frac{Z_i \overline{T}^{\mathcal{O}^i}}{Z}. \qquad (36)$$

The equation is still nonlinear, because of $Z_i = Z_i(\bar{y})$, and the $\bar{y}$ are still nonlocally coupled. In this situation the space of possible solutions is the convex hull spanned by the low temperature limits $\overline{T}^{\mathcal{O}^i}$, which are $\bar{y}$–independent. The high temperature limit becomes $\overline{T}^{\mathcal{O}} = \frac{1}{N_i}\sum_i^{N_i} \overline{T}^{\mathcal{O}^i}$. We will call the $\mathcal{O}$–distance

$$d_{\mathcal{O}}(\bar{y}, \bar{y}') = \sqrt{<\bar{y} - \bar{y}'|\mathcal{O}|\bar{y} - \bar{y}'>}$$

$$= \sqrt{||\bar{y} - \bar{y}'||_{\mathcal{O}}^2}, \qquad (37)$$

the *canonical distance* of a finite temperature regularization problem with $\mathcal{O} = \mathcal{O}^i$, and

$$d_{\mathcal{O},T}(\bar{y}, \bar{y}') = \frac{d_{\mathcal{O}}(\bar{y}, \bar{y}')}{\max_{i,j} d_{\mathcal{O}}(T^i, T^j)}, \qquad (38)$$

the *normalized canonical distance*. Solutions $\bar{y}^*$ of Eq.(36) must have $0 \leq d_{\mathcal{O},T}(\bar{y}^* - T^i) \leq 1$, for all $i$. Notice, that the canonical distance of a regularization problem depends over $\mathcal{N}^D$ and the normalized canonical distance over $\mathcal{N}^D$ and $\bar{D}$ from the actual given data.

For equal $\mathcal{O}^i$ the exponents can be diagonalized simultaneously, so that

$$< \bar{y} - \overline{T}_i^{\mathcal{O}} |\mathcal{O}| \bar{y} - \overline{T}_i^{\mathcal{O}} > = < \bar{y} - \overline{T}_i^{\mathcal{O}} |\mathcal{U}\mathcal{D}\mathcal{U}^\dagger| \bar{y} - \overline{T}_i^{\mathcal{O}} >$$

with diagonal $\mathcal{D}$. Hence, in the corresponding eigenvector representation the operators in all exponents are local simultaneously. We saw already that $\mathcal{O}$ is data dependent, and so are therefore also its eigenvectors.

### 9.6.2 Mean field equations for regularization with polynomial interaction

For the example (32) one finds with $h(L) = g^{-1}(L)$

$$
\begin{aligned}
\mathbf{O}^I &= \sum_i \sum_{j_i} M^{ij_i} \mathcal{O}^{ij_i} \\[2mm]
\mathbf{t}^I &= \sum_i \sum_{j_i} M^{ij_i} \mathcal{O}^{ij_i} T^{ij_i},
\end{aligned}
\tag{39}
$$

with

$$M^{ij_i} = M^{ij_i}(\bar{y}) = \prod_{k_i \neq j_i} < \bar{y} - T^{ik_i} |\mathcal{O}^{ik_i}| \bar{y} - T^{ik_i} > .$$

In the case of only one template $T^{ij_i} = T$, i.e. $\mathbf{t}^I = \mathcal{O}T$ this has also as trivial solution $\bar{y} = T$. For $\mathcal{O}^i = \mathcal{O}/n_i$ the equation reduces to

$$\bar{y} = \frac{\sum_i \sum_{j_i} M^{ij_i} T^{ij_i}}{M},$$

with

$$M = M(\bar{y}) = \sum_i \sum_{j_i} M^{ij_i}.$$

Hence, the solutions are restricted to the space spanned by convex combinations of the $T^{ij_i}$.

In most cases nonlinear equations can only be solved numerically with the help of iteration procedures.

### 9.7 Iteration procedures = learning algorithms

Here we will discuss iteration procedures to solve for the inhomogeneous integro–differential equations.[77] Iteration procedures correspond to the actual learning algorithms. We consider here their application to MaP equations, but the methods also apply to MiR equations or a full Bayesian approach.

An iteration procedure is defined by a function $G$

$$f^{0,i+1} = G^i(f^{0,i}),$$

producing new guesses $f^{0,i+1}$ from a current guess $f^{0,i}$, and with fixed points being solutions of the stationarity conditions. We restrict in the following to solving the stationarity conditions, and assume the maximum

---

[77] For an introduction to numerical methods see for example Hackbusch, 1989, Press, Teukolsky, Vetterling, Flannery, 1992, and references therein.

(or minimum, saddle point) conditions, i.e. the second derivatives, to be checked separately. We can construct an iteration procedure by choosing a function $H^i$, which we also allow in general to be (also stochastically) $i$–dependent, and write for $\bar{y}^{i+1} = G^i(\bar{y}^i)$

$$f^{0,i+1} = G^i(f^{0,i}) = f^{0,i} + \tilde{H}^i\left(\frac{dL(f^{0,i})}{df^{0,i}}\right).$$

to ensure fixed points are solutions of the stationarity equations we require $\tilde{H}^i(0) = 0$, and not to create additional spurious solutions one must have $\tilde{H}^i(x) \neq 0$ for $x \neq 0$. We can fulfill those conditions by defining a (possibly $x$–dependent) nonsingular linear mapping (a matrix for vector $x$) $H^i(x)$ acting on $x$ with $\det(H^i(x)) \neq 0, \forall x \neq 0$, by $\tilde{H}^i(x) = H^i(x)x$. The matrix $H^i(x)$ is usually chosen positive definite (and symmetric) for a maximization problem. $H^i(x)$ may result from a (bijective) transformation of the independent variables $f^0 = T^i(f'^0)$ and/or from a strictly monotonic transformation $h^i(L)$ of the log-posterior. For everywhere differentiable transformations the chain rule gives

$$\frac{dh^i(L(T^i(f'^0)))}{df'^0} = \frac{dh^i(L)}{dL}\frac{dL(f^0)}{df^0}\frac{dT(f'^0)}{df'^0}.$$

Strictly monotonic and therefore invertible transformations $h^i$, $T^i$ do not lead to additional stationary points, i.e. spurious solutions. We write

$$f^{0,i+1} = f^{0,i} + H^i(f^{0,i})\frac{dL(f^{0,i})}{df^{0,i}}. \tag{40}$$

For example, gradient, Newton, or Quasi–Newton algorithms correspond to special $H^i$. The $i$–dependence allows also to include methods like conjugate gradient. In general an iteration procedure can be given by an implicit equation $\tilde{G}^i(f^{0,i}, f^{0,i+1}) = 0$. A solution $f^{0,i+1}$, however, has to be given in an explicit form. We can formally extend the parameter vector $f^0$ to a population vector (of parameter vectors) by introducing dummy variables $m$ (population indices), split $f^0$ into several $f_m^0$. Then $H^i$ can induce interactions between different $f_m^0$. Iteration procedures of order $m$ can for example be obtained by

$$f^{0,i+1} = f^{0,i} + \sum_{j=i-m+1}^{i} H^{i,j}(\{f^{0,k}\}_{i-m+1 \leq k \leq i})\frac{dL(f^{0,j})}{df^{0,j}}.$$

The matrix $H^i$ can be a stochastic function (e.g. stochastic annealing) deterministic only in some (zero temperature) limit or individual $H^i$ can be chosen nonsingular only in subspaces like in line search algorithms, as long as a higher iterated equation can be written with a nonsingular $H$ and the convergence check is made for this iterated equation. Similarly, for transformation methods (homotopy, EM–like algorithms) only at a fixed point of the iteration procedure the log-posterior $L = L^i$ corresponds to a strictly monotone transformation of the original problem. For example, $L^i$ can represent a smoothed version of $L$ (e.g. deterministic annealing).

If $G^i$ contains integrals over $\bar{y}_x$, like in the nonlinear terms, already much 'nonlocal information' may be contained in one iteration with $G^i$. Such integrals are for

example introduced by the EM *algorithm* which defines hidden variables $u$, i.e.

$$p(f^0) = e^{L(f^0)} = \int du\, p(f^0, u) = \int du\, e^{L(f^0, u)},$$

with also nonnegative $p(f^0, u) \geq 0$. Then we define the corresponding conditioned variables which have by construction equal norm if summed over $u$, independent of $f^0$

$$p(u|f^0) = e^{L(u|f^0)} = \frac{p(f^0, u)}{\int du\, p(f^0, u)}.$$

This allows to add to $L$ a cross–entropy term in the conditioned variables

$$L(f^0) + \int du\, e^{L(u|f^{0,ref})} L(u|f^0) \tag{41}$$

$$= \int du\, e^{L(u|f^{0,ref})} L(f^0, u) = Q(f^0, f^{0,ref}).$$

This transformation is strictly maximum sufficent relative to $L(f^{0,ref})$ (see Section 6.5) by construction, i.e. because conditioned variables have $f^0$–independent norm over $u$, the additional term is positive and maximal if $L(u|f^0) = L(u|f^{0,ref}), \forall u$, (See Section 5.3.4), which is the case for $f^0 = f^{0,ref}$. The reference point $f^{0,ref}$ has to be adapted during iteration. This has to be done at latest if a local maximum for fixed $f^{0,ref}$ is reached. However, maximizing $Q$ does not require necessarily to find a maximum for every fixed $f^{0,ref}$, it is enough to increase $Q$ with every iteration (Generalized EM). Summarizing, EM–like algorithms use a transformation $h^i(L)$ which is strictly monotonic only at a fixed point and during iteration only strictly maximality (or, respectively, minimality) sufficient relative to some previous guess (see Section 6.5).

While nonsingular $H^i$ ensure that the fixed points of the iteration procedure are zeros of the gradient, this does not guarantee convergence. In general iteration procedures can produce all varieties of features known from discrete dynamical systems, including limit cycles or chaotic behavior (See for example Devaney, 1986, Beck & Schlögl, 1993 and references therein).

We now discuss the example Eq.(33) in more detail. Firstly, we write Eq.(33) in a form

$$\bar{y} = G(\bar{y}),$$

by choosing some additive decomposition $\mathbf{O}(\bar{y}) = \mathcal{A}(\bar{y}) + \mathcal{B}(\bar{y})$, or a decomposition of some $\tilde{H}^i(\mathbf{O})$, with $\mathcal{A}$ positive definite and therefore invertible. This can also be done by directly selecting a convenient $\mathcal{A}$, for example $\bar{y}$–independent, which then defines a corresponding $\mathcal{B} = \mathbf{O} - \mathcal{A}$.

Here $\mathcal{A}$ is usually chosen to be a linear operator, but invertability and not linearity is the crucial property. (However, inverting an nonlinear operator has usually to be done again by iteration, requiring another linear operator to be inverted.) Then we obtain an iteration procedure by defining the left hand side to be $\bar{y}^{i+1}$ and $\bar{y}$ on the right hand side to be $\bar{y}^i$. For our examples this gives a $G(\bar{y})$ of the form

$$G(\bar{y}) = \mathcal{A}^{-1}(\mathbf{t} - \mathcal{B}\bar{y}),$$

$$= \bar{y} - \mathcal{A}^{-1}(\mathbf{O}\bar{y} - \mathbf{t}),$$

or in components

$$G_x(\bar{y}) = \int dx'\, \mathcal{A}^{-1}_{x,x'}\left(t_{x'} - \int dx''\, \mathcal{B}_{x',x''}\bar{y}_{x''}\right).$$

The equations, or their variants described below, are solved by choosing a representation (i.e. a linear basis) to write it in component form, for example, in $x$–representation:

$$\bar{y}^{i+1}_x = G_x(\bar{y}^i).$$

In other cases one may prefer to work in another basis, for example plane waves (or general coherent states, Blaizot & Ripka, 1986) to transform a differential equation into an algebraic equation.

To achieve convergence one usually has to include a step-size $\eta$, a method which is also called *relaxation*. Then a new guess $\bar{y}^{i+1}$ is generated from a previous guess $\bar{y}^i$ by mixing only part of the new solution to the old one, giving[78]

$$\bar{y}^{i+1} = G_\eta(\bar{y}^i) = (1 - \eta)\bar{y}^i + \eta G(\bar{y}^i)$$

$$= (1 - \eta)\bar{y}^i + \eta \mathcal{A}^{-1}(\bar{y}^i)(\mathbf{t}(\bar{y}^i) - \mathcal{B}(\bar{y}^i)\bar{y}^i)$$

$$= \bar{y}^i + \eta(G(\bar{y}^i) - \bar{y}^i) \tag{42}$$

$$= \bar{y}^i + \eta \mathcal{A}^{-1}(\bar{y}^i)(\mathbf{t}(\bar{y}^i) - \mathbf{O}(\bar{y}^i)\bar{y}^i)$$

$$= \bar{y}^i + \mathcal{A}^{-1}_\eta(\bar{y}^i)\frac{dL(\bar{y}^i)}{d\bar{y}^i},$$

where $\eta$ can be included in the definition of the operator $\mathcal{A}^{-1}_\eta = \eta \mathcal{A}^{-1}$. The expression $\mathbf{t} - \mathbf{O}\bar{y}$ is the gradient of the log–probability $L$ (or negative gradient of the energy $E$) or the (negative or positive, respectively) residual at point $\bar{y}$. For linear $\mathcal{A}$ we recognize an iteration procedure of the form (40) with

$$H(L) = \eta \mathcal{A}^{-1}.$$

When a linear approximation of $\frac{dL}{d\bar{y}}$ is possible in the neighborhood of some $\bar{y}^0$ the convergence depends on the spectral radius of $\mathcal{I} + \mathcal{A}^{-1}\mathcal{H}(\bar{y}^0)$, where $\mathcal{H}$ denotes the Hessian and $\mathcal{I}$ the identity. In the linear approximation the Newton algorithm is optimal. For a linear equation and fixed $\mathcal{A}$, choosing $\eta < 1$ is also called underrelaxation, and $\eta > 1$, which can improve convergence, is called overrelaxation. (See for example Press, Teukolsky, Vetterling, Flannery, 1992). For example, in the finite temperature model (30) the Hessian for the functional $p = \sum_i Z_i/Z$ reads

$$\mathcal{H}^M(p) = -\beta \sum_i Z_i \left(\mathcal{O}^i + \beta \mathcal{O}^i|\bar{y} - \overline{T}^{\mathcal{O}^i}><\bar{y} - \overline{T}^{\mathcal{O}^i}|\mathcal{O}^i\right).$$

---

[78]Multiplying with $(1/\eta)\mathcal{A}$ and projecting onto an infinitesimal $< d\bar{y}|$ the iteration procedure $\bar{y}^{i+1} = \bar{y}^i - \eta \mathcal{A}^{-1}(\bar{y}^i)(\mathbf{O}(\bar{y}^i)\bar{y}^i - \mathbf{t}(\bar{y}^i))$ can be written $1/\eta < d\bar{y}|\mathcal{A}|\Delta\bar{y}^i > = < d\bar{y}|(\delta L/\delta\bar{y})|_{\bar{y}=\bar{y}^i} >= dL$. For infinitesimal $|\Delta\bar{y}^i >= |\bar{y}^{i+1} - \bar{y}^i >$ approximately equal to $d\bar{y}$ this shows that for positive (semi) definite $\mathcal{A}$ the differential $dL$ is larger (or equal) to zero. Thus, the functional $L$ increases during iteration for $\eta$ small enough.

while one finds for $L = \ln \sum_i Z_i / Z$

$$\mathcal{H}^M = \sum_i \frac{Z_i}{Z} \left( -\beta \mathcal{O}^i + \beta^2 \, |\mathcal{O}^i \bar{y} - T^i{>}{<}\mathcal{O}^i \bar{y} - T^i| \right)$$

$$-\beta^2 \left( \sum_{ij} \frac{Z_i}{Z} |\mathcal{O}^i \bar{y} - T^i{>}{<}\mathcal{O}^j \bar{y} - T^j| \frac{Z_j}{Z} \right).$$

Notice, that in the high temperatur limit the terms proportional to $\beta^2$ vanish faster than $-\beta \sum_i \frac{Z_i}{Z}\mathcal{O}^i$. For only one mixture component the last two terms compensate for $L = \ln Z$ and the second term vanishes at the stationary point $\bar{y} = \mathcal{O}^{-1}T$ for $p = Z$. For large deviations $\mathcal{O}^i \bar{y} - T^i$ the Hessian does not need to be negative definite even for positive definite $\mathcal{O}$. Similarly, the second derivative $d\left( dg^{-1}(L)/d\bar{y}(x) \right)/d\bar{y}(x') = d\left( \mathbf{t}^I - \mathbf{O}^I \bar{y} \right)/d\bar{y}(x')$ gives for the Landau–Ginzburg model (39)

$$\mathcal{H}^I = \sum_i \sum_{j_i} M^{ij_i} \mathcal{O}^{ij_i}$$

$$+ \sum_i \sum_{j_i} \sum_{k_i} M^{ij_i k_i} \mathcal{O}^{ij_i} |\bar{y} - T^{ij_i}{>}{<}\bar{y} - T^{ij_i}| \mathcal{O}^{ij_i},$$

with

$$M^{ij_i k_i} = \prod_{l_i \neq k_i \neq j_i} <\bar{y} - T^{il_i}|\mathcal{O}^{ij_i}|\bar{y} - T^{il_i}> .$$

$G_\eta(\bar{y}^i) = G_{\eta(i,\bar{y}^i)}(i, \bar{y}^i)$ can be chosen $i$–dependent by varying $\eta = \eta^i$ or $\mathcal{A} = \mathcal{A}_{(i)}$, which can include dependence on past values $y^{0,k}, k \leq i$. The operator $\mathcal{A}$ should be chosen adapted to the problem, i.e. approximating $\mathcal{O}$ or, at least near a stationary point even better, the Hessian $\mathcal{H}$. A not exactly positive definit $\mathcal{A}$ might be helpful in the beginning of an iteration step, if easy to invert and leading 'mainly' in the right direction. Then a proper $\mathcal{A}$ (e.g. equal to the identity, see below) can be chosen in subsequent iterations, when the solution $\bar{y}$ is already approximately correct. Convergence is not necessarily guaranteed, but depends, besides on choosing a good initial guess, on adjusting the relaxation factor $\eta$. Choosing $|\eta|$ small enough the change $|\bar{y}^{i+1} - \bar{y}^i|$ becomes arbitrary small, increasing the 'resolution' of the search and also reducing oscillations. This usually allows methods which search in directions with always positive (or negative) projections on the gradient to reach at least a local maximum (minimum) if the step-size is small enough. (For convergence results see for example Bertsekas, 1995, Golden, 1996, and references therein.) The step-size can also be determined by a line search in the direction given by $G^i(\bar{y}^i) - \bar{y}^i$.

The usual *gradient algorithm* or method of steepest descent is a special iteration procedure of this type: If $dL/d\bar{y} = (\mathbf{O}\bar{y} - \mathbf{t})$ and if $\mathcal{A}$ is the identity operator, then the term $\bar{y} - G(\bar{y})$ is the gradient of $L$. Iteration schemes can also be related to the gradient of other surfaces but with the same stationary points. For example, for Eq.(33) the iteration (42) is for $\mathcal{A} = 1$ a gradient algorithm on the surface $p = e^L$ and for a general linear operator $\mathcal{A}$ a gradient on a surface parameterized by

variables $T(f^0) = \bar{z} = \sqrt{\mathcal{A}}\bar{y}$. For positive definite $\mathcal{A}$ the square root exists and we have:

$$\frac{dL(\bar{z})}{d\bar{z}} = \left( \sqrt{\mathcal{A}} \right)^{-1} \frac{dL(\bar{y})}{d\bar{y}}.$$

Here the square root is equal to its transpose as we understand positive definite to include symmetric. Thus, by multiplying with $\sqrt{\mathcal{A}}$ one finds

$$\bar{y}^{i+1} = \bar{y}^i + \eta \mathcal{A}^{-1} \frac{dL(\bar{y}^i)}{d\bar{y}}$$

$$\Leftrightarrow \bar{z}^{i+1} = \bar{z}^i + \eta \frac{dL(\bar{z}^i)}{d\bar{z}}.$$

To apply EM–like algorithms we can choose for a log-posterior of the form (30) with only positive terms ($c_i$ real) the summation index $i$ (which has nothing to do with the iteration $i$) as hidden variable

$$p(f^0) = \int du \, p(f^0, u) = \sum_i p(f^0, u_i),$$

with

$$p(f^0, u_i) = e^{L(f^0, u_i)} = e^{-\frac{\beta}{2} \left( \sum_{j_i}^{N_{j_i}} <\bar{y} - T^{ij_i}|\mathcal{O}^{ij_i}|\bar{y} - T^{ij_i}> \right) + c_i}$$

so that

$$p(u_i|f^0) = e^{L(u_i|f^0)} =$$

$$\frac{e^{-\frac{\beta}{2} \left( \sum_{j_i}^{N_{j_i}} <\bar{y} - T^{ij_i}|\mathcal{O}^{ij_i}|\bar{y} - T^{ij_i}> \right) + c_i}}{\sum_i e^{-\frac{\beta}{2} \left( \sum_{j_i}^{N_{j_i}} <\bar{y} - T^{ij_i}|\mathcal{O}^{ij_i}|\bar{y} - T^{ij_i}> \right) + c_i}} = \frac{Z_i}{Z},$$

and choosing a reference $f^{0,ref}$

$$Q(f^0, f^{0,ref}) = \frac{\sum_i Z_i^{ref} \ln Z_i}{Z^{ref}} = \overline{\ln Z_i}^{ref},$$

where

$$Z_i^{ref} = e^{-\frac{\beta}{2} \left( \sum_{j_i}^{N_{j_i}} <\bar{y}^{ref} - T^{ij_i}|\mathcal{O}^{ij_i}|\bar{y}^{ref} - T^{ij_i}> \right) + c_i},$$

$$Z^{ref} = \sum_i Z_i^{ref},$$

and

$$\ln Z_i = -\frac{\beta}{2} \left( \sum_{j_i}^{N_{j_i}} <\bar{y} - T^{ij_i}|\mathcal{O}^{ij_i}|\bar{y} - T^{ij_i}> \right) + c_i.$$

So the stationarity condition for fixed $f^{0,ref}$ reads

$$\overline{\mathcal{O}}^{ref} \bar{y} = \overline{T}^{ref},$$

with

$$\overline{\mathcal{O}}^{ref} = \frac{\sum_i Z_i^{ref} \mathcal{O}^i}{Z^{ref}} = Z^{ref} \mathcal{O}^{ref},$$

$$\overline{T}^{ref} = \frac{\sum_i Z_i^{ref} T^i}{Z^{ref}} = Z^{ref} T^{ref},$$

or, after multiplying by $Z^{ref}$,

$$\mathcal{O}^{ref} \bar{y} = T^{ref}.$$

With $\mathcal{O}^{i,j_i}$ being linear operators and therefore quadratic $\ln Z_i$, this is a linear equation, which can be solved in one step by inverting $\mathcal{O}^{ref}$. In case, a direct inversion is not feasible, this inversion can also be approximated by iterative procedures. In general, however, the $\ln Z_i$ can be a non–quadratic function. The model $F$ can for example allow varying mixture coefficients (included in the $c_i$) or different variances (included in the factors of $\mathcal{O}^{ij_i}$) which then have to be included in the optimization process. Such additional parameters are part of the description of $f^0$. If we implement them by a 'hard OR' with uniform prior on the allowed space this gives not rise to additional terms and means practically minimizing the resulting equations also with respect to the additional parameters. In general we can also add prior terms for the additional parameters. One must be careful however about the range of parameters consistent with prior knowledge. Allowing for example to optimize the relative weight of data and smoothness terms on the training set can end up in the so called '$\delta$–catastrophe', i.e. a solution having peaks at every data point and in case of a zero template being zero elsewhere, a situation most times not intended to be a very likely member of $F^0$. Thus, the EM algorithm for Gaussian mixtures can be seen as a method solving a reference equation linear in $\bar{y}$. The stationarity equations for fixed reference $f^0$ can become at least partly nonlinear if $\mathcal{O}^{ij_i}$, $c_i$ or the $\bar{y}$ itself are parameterized nonlinearly. One may also use cross–validation to determine those parameters.

The EM transformation (41) does not yet define the maximization procedure used to maximize $Q(f^0, f^{0,ref})$, i.e. a $\tilde{H}^i(\mathcal{O}^{ref})$ can be chosen and splitted in $\mathcal{A}$ and $\mathcal{B}$ in various ways. Every iteration procedure has to separate the occurrences of $\bar{y}$ in the stationarity conditions into old $\bar{y}^i$ and new $\bar{y}^{i+1}$. An $EM$–algorithm treats occurrences of $\bar{y}$ at two time scales: some are renewed during maximizing $Q$, others when changing $\bar{y}^{ref}$. Thus, $Q$ may be maximized by any method, including such based on a random search, gradient–like, or EM–like algorithms. See for example the Helmholtz machine (Dayan, Hinton, Neal, & Zemel, 1995; Hinton, Dayan, Frey, Neal, 1995) for an application of the EM algorithm to hierarchically defined $p(f^0)$.

Table 4 gives the, in general ieration dependent, matrices $\mathcal{A}_{(i)}^{-1}$ for some common iteration (or learning) procedures. They are special cases of relaxation techniques and most of them are local, i.e. they only depend on one previous guess $\bar{y}^i$ and derivatives at that location, The gradient corresponds to choosing $\mathcal{A}^{-1}$ equalt to the identity $\mathcal{I}$. Jacobi iteration uses a diagonal $\mathcal{A}$, e.g. the diagonal part of $\mathbf{O}$, the Gauss–Seidel method includes also the lower triangular part, e.g. of $\mathbf{O}$. Newton's method takes the negative Hessian $\mathcal{H}$ for minimization and maximization. More precisly, the Newton method uses the given formula at locations where the Hessian is negative definite (or positive definite for minimization), at other locations the method has to resort to any other minimization algorithm. Quasi–Newton methods try to approximate the Hessian $\mathcal{H}$. In the table the abbreviations $\Delta_{(0)}^i = \bar{y}^i - \bar{y}^{i-1}$ and $\Delta_{(1)}^i =$

**Local learning algorithms**

| | |
|---|---|
| Gradient | $\mathcal{A}^{-1}_{(i)} = \mathbf{I}$ |
| Jacobi | $\mathcal{A}^{-1}_{(i)}$ diagonal |
| Gauss–Seidel | $\mathcal{A}^{-1}_{(i)}$ triangular |
| Newton | $\mathcal{A}^{-1}_{(i)} = -\frac{d^2 L}{(d\bar{y}(x)d\bar{y}(x'))}\big|_{y^i} = -\mathcal{H}^{-1}$ |
| Quasi–Newton | $\mathcal{A}^{-1}_{(i)} = \mathcal{A}^{-1}_{(i-1)} + \frac{\Delta_{(0)}^i \Delta_{(0)}^i{}^T}{\Delta_{(0)}^i{}^T \Delta_{(1)}^i}$ |
| ( DFP ) | $-\frac{\mathcal{A}^{-1}_{(i-1)} \Delta_{(1)}^i \Delta_{(1)}^i{}^T \mathcal{A}^{-1}_{(i-1)}}{\Delta_{(1)}^i{}^T \mathcal{A}^{-1}_{(i-1)} \Delta_{(1)}^i}$ |
| CG | $\mathcal{A}^{-1}_{(i)}$ |
| | $= \left(1 - \frac{\mathcal{A}^{-1}_{(i-1)}}{\eta^{i-1}}\frac{dL}{dy}\big|_{\bar{y}^{i-1}}\frac{\Delta_{(1)}^i{}^T}{\left(\left(\frac{dL}{d\bar{y}}\right)^T\frac{dL}{d\bar{y}}\right)\big|_{\bar{y}^{i-1}}}\right)$ |
| EM | $\mathcal{A}^{-1}_{(i)} = \mathcal{A}'^{-1}_{(i)}\frac{dh^{ref}(L)}{dL}$ |
| | with $h^{ref}(L) =$ |
| | $L + \int du\, e^{L(u|f^{0,ref})} L(u|f^0)$ |

Table 4: Some local learning algorithms

$\frac{dL}{d\bar{y}}\big|_{\bar{y}^i} - \frac{dL}{d\bar{y}}\big|_{\bar{y}^{i-1}}$ are used. The given formula refers to the DFP (Davidon–Fletcher–Powell) method, in the BFGS (Broyden–Fletcher–Goldfarb–Shanno) method, for example, a term $a^i b^i b^i{}^T$ is added with $a^i = \Delta_{(1)}^i{}^T H^{i-1}\Delta_{(1)}^i$ and $b^i = \frac{\Delta_{(0)}^i}{\Delta_{(0)}^i{}^T \Delta_{(1)}^i} - \frac{H^{i-1}\Delta_{(1)}^i}{a^i}$ ($T$ denotes the transpose). Conjugate gradient methods (CG) determine the stepsize by a line search in conjugate directions, obtained by a Gram–Schmidt–procedure. Directions are called conjugate if they are orthogonal in $\mathbf{O}$–distance, assuming we are solving $\mathbf{O}\bar{y} = \mathbf{t}$. For non–quadratic problems they are usually combined with a heuristic to restart the Gram–Schmidt procedure. (See for example Bertsekas, 1995.) For the EM algorithm $H'^i$ defines the chosen iteration algorithms used for $Q$ with fixed reference state. EM algorithms are not restricted to local methods. (So for them the word local in the table caption does not necessarily apply.) Note that $dL/df^0$ in the nonlinear case usually contains integrals over $x$ and therefore even optimization methods which are local with respect to $f^0$ are nonlocal with respect to $x$.

Nonlinear equations do normally have multiple solutions corresponding to several extremal points. If several solutions have to be considered in the last risk minimization step of a MaP–MiR approximation this requires calculation of the relative weight factors of the solutions or widths of the maxima depending on the second derivatives, or a corresponding estimate or assumption (See for example Gelman, Carlin, Stern & Rubin 1995 and especially for neural networks: Buntine & Weigend, 1991; MacKay 1991, 1992b, 1992c; Neal, 1996).

Nonlinear inhomogeneous equations appear for example in scattering theory as approximation to higher dimensional linear inhomogeneous equations (See for example, Austern, 1970, Taylor, 1972, Newton, 1982). There the inhomogeneities (data) are related to the in and out channels representing the boundary conditions or asymptotic states of the wave functions. Numeri-

cal aspects, applications to scattering theory and related higher order approximations are, for example, discussed in (Giraud & Nagarajan, 1991, Lemm, Giraud, & Weiguny, 1994 and Lemm, 1995ab).[79]

Nonlocal templates, or technically the inhomogeneities, can be used in the following way: Instead of using a small space $F_1^0$ to represent possible states $f^0$ of nature and corresponding to hard implemented priors, one allows a larger space $F_2^0$ and implements $f_1^0 \in F_1^0$ within $F_2^0$ as priors with an soft OR by taking $f_1^0 \in F_1^0$ as templates for $F_2^0$. This allows to go beyond $F_1^0$ if the data require. A soft implemented template for $p(f^0|f)$ is not equivalent to using noisy answers for $f^0$: The state of knowledge about $f^0$, that is $p(f^0|f)$, is updated through data, while the noise levels of pure states are assumed to be stationary, i.e. clamped during learning. Templates can be seen as a method of transfer of knowledge between tasks.

# 10    An introductory example

## 10.1    The models

To exemplify the techniques we study a case with one-dimensional $x$ and two full templates, $T^1$, $T^2$, i.e. which are defined for all $x$, in addition to standard data $\bar{D}$. We will study as well the mixture as interaction regularizations.

## 10.2    Finite temperature regularization

To express ($\bar{D}$ AND $T^1$) OR ($\bar{D}$ AND $T^2$) we choose a probability of the form

$$P(\bar{y}) = e^L = \frac{Z}{Z_{F^0}} = \frac{Z_1 + Z_2}{Z_{F^0}}$$

with normalization constant

$$Z_{F^0} = \int_{F^0} df^0 \left( Z_1(f^0) + Z_2(f^0) \right).$$

and Gaussian components

$$P \propto Z = e^{-\beta\frac{1}{2}(\Delta_D + \Delta_1) + c_1} + e^{-\beta\frac{1}{2}(\Delta_D + \Delta_2) + c_2}$$

[79]For example in the Time Independent Mean Field Theory (TIMF) for quantum mechanical scattering one obtaines approximate variational solutions for matrix elements $< \chi'|\mathcal{O}^{-1}|\chi >$ (e.g. $\mathcal{O} = E - H$, with energy $E$ and Hamiltonian $H$ so $\mathcal{O}^{-1}$ is the resolvent of $H$) by choosing $\phi$, $\phi'$ from a space of possible trial functions for which $< \mathcal{O}^{-1}\chi' - \phi'|\mathcal{O}|\mathcal{O}^{-1}\chi - \phi >$ is stationary. Notice the similarity in the role of $\mathcal{O}^{-1}\chi$ or $\mathcal{O}^{-1}\chi'$ and that of a template average $\mathcal{O}^{-1}T^{\mathcal{O}} = \mathcal{O}^{-1}\sum \mathcal{O}^i T^i$. For a mean field approach one chooses $\phi$, $\phi'$ as product of single particle functions. Expanding the quadratic functional gives as variational solution $< \chi'|\mathcal{O}^{-1}|\chi > = < \chi'|\phi > + < \phi'|\chi > - < \phi'|\mathcal{O}|\phi >$. In contrast to the error minimization problems in scattering theory $E$ is in general a complex number and the wave functions $\chi$, $\chi'$, $\phi$, $\phi'$ are allowed to be complex functions. Thus, the stationary points are not maxima or minima but saddle points, and the variational solutions do not yield bounds for the exact solutions. At a saddle point, on the other hand, the effect on the numerical value of the matrix element of different directions of deviations of $\phi$, $\phi'$ from the true solution can have different signs. Deviations from the optimal solution can therefore partly compensate, improving the variational solution.

$$= e^{-\beta\frac{1}{2}(\tilde{\Delta}_1)} + e^{-\beta\frac{1}{2}(\tilde{\Delta}_2)}$$

$$= e^{-\frac{\beta}{2}\frac{\tilde{\Delta}_1 + \tilde{\Delta}_2}{2}} 2\cosh\left(\frac{\beta}{4}(\tilde{\Delta}_2 - \tilde{\Delta}_1)\right)$$

$$\propto e^{-\frac{\beta}{2}\frac{\tilde{\Delta}_1 + \tilde{\Delta}_2}{2}} 2\cosh\left(\frac{\beta}{4}(\tilde{\Delta}_2 - \tilde{\Delta}_1)\right),$$

with parameter $\beta > 0$, $c_i$ real,

$$\Delta_D = < \bar{y} - \bar{D}|\mathcal{O}^D|\bar{y} - \bar{D} >,$$

$$\Delta_i = < \bar{y} - T^i|\mathcal{O}^{s,i}|\bar{y} - T^i >,$$

and for combining data and template term in the same exponent

$$\tilde{\Delta}_i = \bar{\Delta}_i + \sum_j <T^{i,j}|\mathcal{O}^{i,j}|T^{i,j}> - <\overline{T}^{\mathcal{O}^i}|\mathcal{O}^i|\overline{T}^{\mathcal{O}^i}> -2c_i/\beta,$$

$$\bar{\Delta}_i = <\bar{y} - \overline{T}^{\mathcal{O}^i}|\mathcal{O}^i|\bar{y} - \overline{T}^{\mathcal{O}^i}>,$$

$$\overline{T}^{\mathcal{O}^i} = \left(\mathcal{O}_{\mathcal{P}_i}^i\right)^{-1} \sum_j \mathcal{O}^{i,j} T^{i,j},$$

$$\mathcal{O}^i = \mathcal{O}^D + \mathcal{O}^{s,i},$$

$$T^{i,1} = \bar{D}, \quad T^{1,2} = T_1, \qquad T^{2,2} = T_2,$$

and we will write

$$T^i = T^{i,1} + T^{i,2}.$$

The data operator is diagonal in $x$–representation

$$\mathcal{O}^D = \lambda_D \mathcal{N}^D,$$

and for the template operator $\mathcal{O}^s$ we choose the same for $T_1$ and $T_2$, i.e. $\mathcal{O}^{s,1} = \mathcal{O}^{s,2}$. This gives for the operator $\mathbf{O}^M = \overline{\mathcal{O}}^Z$

$$\mathbf{O}^M = \overline{\mathcal{O}}^Z = \frac{1}{Z}\sum_i e^{-\frac{\beta}{2}\tilde{\Delta}_i}\mathcal{O}^i = \left(\lambda_D\mathcal{N}^D + \mathcal{O}^s\right). \quad (43)$$

Thus, the two terms of Eq.34 corresponding to $i = 1$ and $i = 2$ coincide and the factor $Z = \sum_i Z_i$ cancels out.

We select an operator related to a smoothness measure,

$$\mathcal{O}^s = \lambda_0\mathcal{O}_1 + \lambda_2\mathcal{O}_2 + \lambda_4\mathcal{O}_4,$$

with

$$\mathcal{O}_1(x,x') = \mathcal{I}(x,x') = \delta(x - x'),$$

$$\mathcal{O}_2(x,x') = -\delta(x - x')\frac{d^2}{dx^2},$$

and

$$\mathcal{O}_4(x,x') = \delta(x - x')\frac{d^4}{dx^4}.$$

The $\lambda_i$ and $\lambda_D$ allow changing the weight of the four parts. The inhomogeneous side has the form

$$\overline{T}^Z = \frac{e^{-\beta\frac{1}{2}(\Delta_D + \Delta_1) + c_1}}{Z}\left(\lambda_D\mathcal{N}^D\bar{D} + \mathcal{O}^sT_1\right)$$

$$+ \frac{e^{-\beta\frac{1}{2}(\Delta_D + \Delta_2) + c_2}}{Z}\left(\lambda_D\mathcal{N}^D\bar{D} + \mathcal{O}^sT_2\right), \qquad (44)$$

with $Z \propto P$. As the operators are the same for both $T^i$ they have equal normalization constants and the $c_i$ are directly the logarithms of the mixture coefficients. The proportionality factor $Z_{F^0}$ cancels out as well as

the variance–like term arising from combining the data templates terms for equal $x$. For the model equation $\overline{\mathcal{O}}^Z \bar{y} = \overline{T}^Z$ we could, for example, write in the general case $\mathcal{O}^1 \neq \mathcal{O}^2$

$$\left( \frac{\mathcal{O}^1 + \mathcal{O}^2}{2} + \tanh\left( \frac{\beta}{4}(\tilde{\Delta}_2 - \tilde{\Delta}_1) \right) \frac{(\mathcal{O}^1 - \mathcal{O}^2)}{2} \right) \bar{y}$$

$$= \frac{\mathcal{O}^1 T^1 + \mathcal{O}^2 T^2}{2} + \tanh\left( \frac{\beta}{4}(\tilde{\Delta}_2 - \tilde{\Delta}_1) \right) \frac{\mathcal{O}^1 T^1 - \mathcal{O}^2 T^2}{2}.$$

However, we have already seen in Section 9.6.1 that in the case of equal operators $\mathcal{O}^i$ the equation

$$\left( \lambda_D \mathcal{N}^D + \mathcal{O}^s \right) \bar{y} = \mathcal{N}^D \bar{D}$$

$$+ \frac{e^{-\beta \frac{1}{2}(\Delta_D + \Delta_1) + c_1}}{Z} \mathcal{O}^s T_1 + \frac{e^{-\beta \frac{1}{2}(\Delta_D + \Delta_2) + c_2}}{Z} \mathcal{O}^s T_2,$$

can be simplified to

$$\bar{y} = \overline{T}^{\mathcal{O}^i, Z} = \sum_i \frac{Z_i \overline{T}^{\mathcal{O}^i}}{Z}. \tag{45}$$

For two templates this may also be written as

$$\bar{y} = \overline{T}^{\mathcal{O}} + \tanh\left( \frac{\beta}{4}(\tilde{\Delta}_2 - \tilde{\Delta}_1) \right) \frac{\overline{T}^{\mathcal{O}^1} - \overline{T}^{\mathcal{O}^2}}{2}.$$

Eq. (45) shows that in this simplest case of only two templates with equal operators $\mathcal{O}^i$ the space of solutions is effectively one–dimensional. It is the line spanned by convex combinations of the two $\overline{T}^{\mathcal{O}^i}$. (The superscript $\mathcal{O}^i$ is kept as reminder that an operator inversion $(\mathcal{O}^i)^{-1}$ is needed to obtain this template average and the $i$–dependence of the $T^i$ remains.) We also see, that for $\beta = 0$ where the tanh is also zero, this gives the correct high temperature solution $\bar{y}(\beta = 0) = \overline{T}^{\mathcal{O}}$. For $\beta \to \infty$ the tanh becomes $\pm 1$ depending on the sign of $\tilde{\Delta}_2 - \tilde{\Delta}_1$. Hence, we find correctly as self–consistent low temperature solutions the component templates $\overline{T}^{\mathcal{O}^1}$ and $\overline{T}^{\mathcal{O}^2}$.

For the symmetric case with $\tilde{\Delta}_2 - \tilde{\Delta}_1 = \bar{\Delta}_2 - \bar{\Delta}_1$ one stationary solution in $\mathcal{P}$–space is easily found. Then the equation $\bar{y} = \overline{T}^{\mathcal{O}}$ for $\bar{\Delta}_2 - \bar{\Delta}_1 = 0$ is consistent with the definition of the template average for which we found $\bar{\Delta}_1 = \bar{\Delta}_2$. We will see however, that this solution of the stationarity equation is only a minimum for small enough $\beta$, (the "high temperature phase").

Choosing only one $x$, i.e. $|X| = 1$, with $\mathcal{O}^s = 1$ and $T_1 = T^1 = 1$, $T_2 = T^2 = -1$ the model equation reduces to the celebrated mean field equation of a ferromagnet with uniform couplings

$$\bar{y} = \tanh\left( \beta \bar{y} \right). \tag{46}$$

Here the templates $T^i$ (representing prior knowledge) are the analogon of possible states[80] of a physical sys-

tem (which has also to be specified a-priori)[81] and the mean $\bar{y}$ represents in both cases the observable we are interested in. The data, which update our knowledge can be called local fields, changing (correcting) a-priori given templates $T_i$ to combined templates $T^i$, i.e. the final (posterior) mixture components.

At $\beta = 1$ Eq.(46) shows a bifurcation phenomena, as for $\beta \leq 1$ there is only one solution $\bar{y} = 0$, while for $\beta > 1$ two new solutions appear and $\bar{y} = 0$ becomes unstable. If the equation is seen as a phenomenological description of a large (or infinite) system this is also called a *phase transition*. Indeed, such bifurcations or phase transitions are typical for mixture distributions. For example, magnetic systems have many connections to neural networks, especially to Hopfield nets, see e.g. Hertz, Krogh, & Palmer 1991. A interesting clustering algorithm showing these phase transitions can be found in Rose, Gurewitz, & Fox, 1990. For a real magnetic implementation see also Blatt, Wiseman, & Domany, 1997.

### 10.2.1 Some remarks on local templates

Generalization requires correlations between different $x$, i.e. nonlocal dependencies. We implemented those dependencies by giving explicitly global templates. Alternatively, those dependencies can arise from coupled local templates $T_x^i$. Local templates means that we construct a model by adding mixture components for every possible local state $i_x$ for every single $x \in X$. Nonlocal dependencies are present if mixture terms for $x$ depend also on $x' \neq x$. For the whole system the number of states grows exponentially with $|X|$. This leads for large $|X|$ (especially in the limit of continuous $x$ ) to two obvious problems:

---

[81] One may consider templates for physical systems (states) to be more 'real' than templates originating from fuzzy implementation of prior knowledge. This means, however, only that for physical systems a lower temperature can be actually realized. With regard to 'fuzzy' templates 'nature' is usually in a state of higher temperature and corresponding low temperature states may not be preparable. However, in physical systems the energy function (i.e. states) is (depending on the level of description) not exactly known. Then different possibilities have to be combined by OR, giving mixture probabilities. Examples are spin glasses where the variables on which the log-likelihood depends are separated into states, and interactions. (They are for spin glasses, however, treated asymmetrically: marginalization over interactions not states). As long as the probability distribution over different interactions cannot be made deterministic and one does not restrict to non–fluctuating (self-averaging) observables (i.e. those for which averaging over interactions can be skipped) a spin glass cannot be prepared at ('knowledge') temperature zero, i.e. $\beta = \infty$, even if it is at 'thermic temperature' $\beta^\Phi$ (near to) zero. (We have seen in Section 5.3.4 that one may define many temperatures $\beta^j$, related to different parameters of the generating process for $f^0$, and if we distinguish a physical from a knowledge temperature this depends on what part of the process we label thermic. In principle we may substitute every process under the label 'physical' process.) Also for the 'thermic' temperature it is practically impossible to reach the absolute zero point. Hence, the difference between temperature ranges for 'fuzzy systems' and physical states is not of qualitative, but of quantitative nature.

---

[80] Not equal to the states $f^0$ of possible $\bar{y}$.

1. At low temperature the optimization can become extremely difficult or impossible when the number of local minima of $p(f^0)$ is too large to be considered completely.

2. At high temperature the generalization ability can be to small to be useful when the probability distribution $p(f^0)$ is too broad.

For a mixture model with only a few global templates $p(f^0)$ remains non–factorial in the high temperature limit. This means that some combinations of local states remains always excluded and generalization possibilities remain at least for finite $|X|$ in this limit. (While for infinite $|X|$ one must require a remaining finite factor dimension of $p(f^0)$ to allow generalization with respect to data and relevant questions depending on a finite number of $x$. See Section 2.)

The Hopfield model, for example, is a special mixture model with quadratic components (and usually combined with a special iteration dynamic). It is constructed with coupled local templates $T_x^i = \pm 1$. so its number $2^{|X|}$ of global templates $T_X^j$ grows exponentially with $|X|$. Typically the Hopfield net is used as associative memory. There one is interested in retrieving (a large number of) stored patterns by varying the starting point for the iteration procedure.[82] Even although optimally not done at zero temperature, where many unwanted mixture states are stable, its use as associative memory has the nature of a low temperature application, because the memories shall be retrieved as near as possible to the original stored pattern. Its use is limited by the on–set of the spin glass phase. (See for example Amit, Gutfreund, & Sompolinski, 1987; Amit, 1989.) In general for a system with local templates the generalization possibility can break down completely in the high temperature limit when all templates become equally likely and $p(f^0)$ gets factorial (See Section 2).

For nonlinear regularization we have mainly fuzzy logical applications in mind where the number of fuzzy logical alternatives (templates) is not extremely large but comparable to the number of alternatives in typical problems solved by logical methods. We are, however, especially interested in the interpolation between templates, i.e. in the deformed solutions under given data fields at finite temperature.

One can see the two–(or few–)template case as an effective model of an intermediate temperature range where two templates $T^1$ and $T^2$ are near their phase transition at $\beta^*$. Those two templates are then considered as high temperature averages of finer 'constituent' templates which at temperature $\beta^*$ cannot be distinguished while not considered templates are treated in their low temperature approximation and already excluded. Our Gaussian two template example then corresponds to a Gaussian approximation, ('oscillators' for discrete $x$, 'free field' or random phase approximation

for continuous $x$, however in a quite general form corresponding to the chosen $\mathcal{O}$) for the two effective templates $T^1$ and $T^2$. Thus, a mixture model with two global templates, defined for all $x$, represents a system capable of two Gaussian (process) states. This system is at zero temperature in one of those two possible global states, and at nonzero temperature in a mixture (weighted OR for disjunct events) of those two states.

For a model with global templates the local states $T_x^i$ for all $x$ are already combined into global states $T^i$ by AND i.e. by the sum inside the scalar product in the log-probability. The logarithm of the partition sum for $N$ global templates

$$\ln Z = \ln \sum_i^N e^{-\beta E_i(X)} = \ln \sum_i^N \prod_x^{|X|} Z_{x,i}(X)$$

is of the form $\mathrm{OR}_i\ \mathrm{AND}_x\ Z_{x,i}$. with $Z_{x,i}$ being the (effective) partition sum for a single $x$ in the global template (state) $i$ of the complete system. Notice, that if the system includes nonlocal interactions (e.g. smoothness) then $Z_{x,i} = Z_{x,i}(X)$ depends also on $x' \neq x$ and $Z$ does not factorize into local components depending only on single $x$.

Constructing a system instead out of global states (templates $T^j$) out of combinations of local states (templates $T_x^i$) for each $x$ leads to

$$\ln Z = \sum_x^{|X|} \ln \sum_i^n Z_{x,i}(X) = \ln \prod_x^{|X|} \sum_i^n Z_{x,i}(X)$$

$$= \ln \sum_{i_1}^n \cdots \sum_{i_{|X|}}^n \prod_x^{|X|} Z_{x,i_x}(X) = \ln \sum_{i'}^{n^{|X|}} \prod_x^{|X|} Z_{x,i'}(X),$$

with multi–index $i' = (i_1, \cdots, i_{|X|})$ and $Z_{x,i'}(X) = Z_{x,i_x}(X)$. This corresponds to $N = n^{|X|}$ global templates. The sum reduces if probabilities for certain combinations of local templates are zero. Otherwise one has $n$ different disjunct local states $i$ for every subsystem $x$ and correspondingly $n^{|X|}$ templates for the whole system with, depending on the considered interaction and dynamic as many potential candidates for (meta)stable states. Notice the similarity of this form of $\ln Z$ with the form obtained by averaging for spin glasses (See first footnote in Section 5.3.4). The sums over all configurations $x \in X$ of composite systems, can sometimes be reduced to a product over single component sums (that means 'exchanging' $\sum$ with $\prod$ under the logarithm) at the cost of introducing new variables. Such an embedding of the system in a larger space can make it easier solvable (similar to the idea of introducing Lagrange parameters), or a certain approximation scheme (e.g. saddle point approximation) can become applicable. Quadratic interactions for example are linearized by the Hubbard–Stratonovich transformation where the Gaussian integral formula is used 'backwards': $e^{\bar{y}_x^2/4a} = \sqrt{a/\pi} \int d\mu\, e^{-a\mu^2 \pm \bar{y}_x \mu}$. Then the total partition sum factorizes in $x$ and the sum can be performed over local components. The remaining integral over $\mu$ (called order parameter) can be performed in saddle point approximation. Similarly, restrictions like $\delta$–functions or

---

[82]Hence, the iteration procedure corresponds in this case to retrieval and not to learning. Learning in the Hopfield net, i.e. finding the weights $\mathcal{W}$ so that the correct patterns (templates) are stable corresponds on this level of comparison to the determination of priors $p(f^0)$.

step functions $\Theta(x)$ can be written in an integral representation, which creates new order parameters for a subsequent saddle point approximation.

For example, a ferromagnetic mean field equation with nonuniform coupling looks like $\bar{y} = \tanh\left(\beta\mathcal{W}\bar{y} + \beta h^{ext}\right)$, to be read component-wise, with an external field vector $h^{ext}$. The coupling matrix or operator $\mathcal{W}$, causes couplings between different $x$ values of a vector $\bar{y}$. To obtain such an equation one has to use local templates, independent for different $x$. A model log–probability would look like $\sum_x \ln \sum_{i=1}^2 e^{L_x + L_{x,i}(X)} = \sum_x (L_x + \ln\cosh(L_{x,i}(X)))$ for $L_{x,1} = -L_{x,2}$ with, for $L_x = 0$, a derivative of the form $\frac{dL}{d\bar{y}_x} = \sum_{x'} \tanh\left(\frac{dL_{x',i}}{\bar{y}_x}\right)$. Neglecting the non–diagonal terms give equations of the structure of mean field equations with nonuniform coupling operator (e.g. nearest neighbors) $\mathcal{W}$ and external field $h^{ext}$ contained in $L_{x,i}$. For quadratic interactions such equations are usually obtained using the Hubbard–Stratonovich transformation (see below).

## 10.3 Landau–Ginzburg regularization

For an interaction version we can approximate a structure ($\bar{D}$ AND ($T^1$ OR $T^2$)) in a naive (fuzzy) implementation as
$$L^{I_1} = -g(\gamma_D \Delta_D + \gamma \Delta_1 \Delta_2),$$
or, similarly, use the structure ($\bar{D}$ AND $T^1$) OR ($\bar{D}$ AND $T^2$)
$$L^{I_2} = -g((\gamma_D \Delta_D + \gamma\Delta_1)(\gamma_D \Delta_D + \gamma\Delta_2))$$
$$= -g(\gamma_D^2 \Delta_D^2 + \gamma_D \gamma \Delta_D (\Delta_1 + \Delta_2) + \gamma^2 \Delta_1 \Delta_2).$$

These $L$ have a polynomial structure. (The strictly monotonically increasing $g$ does not change the location of extrema, even if non–polynomial.) The parameter $\gamma$, $\gamma_D$ parameterize the relative weights of data and template terms in the energy (log-probability, error) function. Because extrema of $L^{I_i}$ are independent of a scaling factor, we can always choose $\gamma_D = 1 - \gamma$. To have a parameter with values between zero and infinity, like $\beta$ in a mixture or finite temperature regularization, we can use
$$\beta^I = \frac{\gamma}{1-\gamma}, \quad \gamma = \frac{\beta^I}{1+\beta^I},$$
so that for example after multiplying with $1/(1 + \beta^I)$, skipping from now on the superscript $I$ for $\beta^I$
$$L^{I_1} = -\tilde{g}(\Delta_D + \beta\Delta_1 \Delta_2),$$
where $\beta$ as in $L^M$ can be seen (after multiplying again with $\beta$) as a common scaling factor for $\lambda_D, \lambda_1, \lambda_2$. Analogously, we get
$$L^{I_2} = -\tilde{g}(\Delta_D^2 + \beta\Delta_D(\Delta_1 + \Delta_2) + \beta^2 \Delta_1 \Delta_2).$$

Conversely, using for the temperature $\beta$
$$\gamma^M = \frac{\beta}{1+\beta}, \quad \beta = \frac{\gamma^M}{1-\gamma^M},$$
has the advantage that the infinite interval $[0, \infty]$ is mapped into the finite interval $[0, 1]$.

Eqs. (39) give for model $L^{I_1}$,
$$\mathbf{O}^{I_1} = \mathcal{O}^D + \beta(\Delta_1 + \Delta_2)\mathcal{O}^s, \tag{47}$$
$$\mathbf{t}^{I_1} = \mathcal{O}^D \bar{D} + \beta\mathcal{O}^s(\Delta_2 T^1 + \Delta_1 T^2), \tag{48}$$
and for model $L^{I_2}$,
$$\mathbf{O}^{I_2} = \left(\mathcal{O}^D + \beta\mathcal{O}^s\right)\left(2\Delta_D + \Delta_1 + \Delta_2\right) \tag{49}$$
$$\mathbf{t}^{I_2} = \left(2\Delta_D + \Delta_1 + \Delta_2\right)\mathcal{O}^D \bar{D} \tag{50}$$
$$+ \beta\mathcal{O}^s\left((\Delta_D - \Delta_2)T_1 + (\Delta_D - \Delta_1)T_2\right).$$

Both models have in the no–data case, $\Delta_D = 0$, the solutions $\bar{y} = T_1$ and $\bar{y} = T_2$, while in the case of missing templates $L^{I_1}$ reduces to a Gaussian data model and the second model $L^{I_2}$ keeps a quadratic term in $\bar{y}$, which however is equivalent to a Gaussian using $\tilde{g}(x) = g(\sqrt{x})$. For $T_1 = T_2$, i.e. $\Delta_1 = \Delta_2 = \Delta_T$, $L^{I_1}$ is not equivalent to a Gaussian $L = \Delta_D + \Delta_T$ and treating data and templates not symmetrically. Using $\sqrt{\Delta_1 \Delta_2}$ in $L^{I_1}$ would restore the Gaussian in the limit, but destroy the polynomial structure.

A high temperature expansion of the mixture model $L^M$ according to (9) gives
$$L^{HT} = c_1 - \beta\frac{\tilde{\Delta}_1 + \tilde{\Delta}_2}{2} + \frac{\beta^2}{8}(\tilde{\Delta}_1 - \tilde{\Delta}_2)^2$$
$$= c_3 - \beta\frac{\bar{\Delta}_1 + \bar{\Delta}_2}{2} + \frac{\beta^2}{8}(\tilde{\Delta}_1 - \tilde{\Delta}_2)^2$$
$$= c_3 - \beta\left(\Delta_D - \frac{\Delta_1 + \Delta_2}{2}\right) + \frac{\beta^2}{8}(\tilde{\Delta}_1 - \tilde{\Delta}_2)^2.$$

Here the no-template case gives the Gaussian $L^{HT} = c - \Delta_D$. For $T_1 = T_2$ and $\mathcal{O}^1 = \mathcal{O}^2$ the difference $(\tilde{\Delta}_1 - \tilde{\Delta}_2)$ is zero, so that $L^{HT} = c - \Delta_D - \Delta_T$, with $\Delta_T = \Delta_1 = \Delta_2$, which is symmetric between data and templates and Gaussian. For $\mathcal{O}^1 = \mathcal{O}^2 = \mathcal{O}^s$ the terms quartic in $\bar{y}$ cancel in $(\Delta_1 - \Delta_2)$ and the $L^{HT}$ quadratic in $\bar{y}$ can only have one extremum.

For the high temperature approximation we find, after dividing by $\beta$
$$\mathbf{O}^{HT} = \frac{\mathcal{O}}{2} + \beta\frac{(\tilde{\Delta}_2 - \tilde{\Delta}_1)}{2}\frac{(\mathcal{O}^1 - \mathcal{O}^2)}{2}, \tag{51}$$
$$\mathbf{t}^{HT} = \frac{T^{\mathcal{O}}}{2} + \beta\frac{(\tilde{\Delta}_2 - \tilde{\Delta}_1)}{2}\frac{(T^1 - T^2)}{2}.$$

For $\mathcal{O}^1 = \mathcal{O}^2$ and $T^1 = T^2$ the high temperature equation becomes $\beta$–independent and, as we already saw, linear in $\bar{y}$, so only one solution can exist. More general, one sees that the temperature independent $\overline{T}^{\mathcal{O}} = \mathcal{O}_{\mathcal{P}}^{-1}\sum_{i=1}^2 \mathcal{O}^i T^i$, for which $\tilde{\Delta}_1 = \tilde{\Delta}_2$, is a self-consistent solution of the high temperature equation. We have already seen that the template average $\overline{T}^{\mathcal{O}}$ is also one mean field solution for finite temperature or mixture regularization.

We may think of a form similar to the high temperature expansion to obtain an effective Landau–Ginzburg log-likelihood which possesses both the high and low

temperature limits of the mixture model.[83] Instead of implementing the OR and using a parameter $\gamma$ weighting the data against the template influence a more temperature–like parameter should interpolate between an AND in the high temperature phase (corresponding to the first term in the high temperature expansion which is according to Section 5.3.4 the first moment or average with respect to the mixture coefficients $a_i/Z_a$) and an OR for the low temperature limit. Thus, we can choose

$$L^{PF} = -g_1\left(\Delta_D + \frac{\Delta_1 + \Delta_2}{2} + \beta\bar{\Delta}_1\bar{\Delta}_2\right)$$

$$= -g_2\left(\frac{\bar{\Delta}_1 + \bar{\Delta}_2}{2} + \beta\bar{\Delta}_1\bar{\Delta}_2\right)$$

$$= -g_3(\bar{\Delta} + \beta\bar{\Delta}_1\bar{\Delta}_2)$$

with $\bar{\Delta}$ resulting from the combination of $\bar{\Delta}_1 + \bar{\Delta}_2$. The superscript in $L^{PF}$ refers to an interpretation of $\bar{y} - \overline{T}^{\mathcal{O}}$ as (an self–interacting) prior field with $\bar{\Delta}$ describing the propagation of an average field and the term $\tilde{\Delta}_1\tilde{\Delta}_2$ the additional 'repulsive' self–interaction. Varying the interaction strength $\beta$ allows to go from the pure average field ($\beta = 0$ or high temperature case) to the purely interacting field ($\beta = \infty$ or low temperature).

Here we have

$$\mathbf{O}^{PF} = \mathcal{O}(1 + 2\beta(\bar{\Delta}_1 + \bar{\Delta}_2)), \tag{52}$$

$$\mathbf{t}^{PF} = T^{\mathcal{O}} + 2\beta(\bar{\Delta}_1 T^{\mathcal{O}^2} + \bar{\Delta}_2 T^{\mathcal{O}^1}), \tag{53}$$

with $\mathcal{O} = \sum_i \mathcal{O}^i$, $\mathcal{O}^i\overline{T}^{\mathcal{O}^i} = T^{\mathcal{O}^i}$, and $T^{\mathcal{O}} = \sum_i T^{\mathcal{O}^i}$. The resulting mean field equations

$$\bar{y} = \frac{\overline{T}^{\mathcal{O}} + 2\beta\mathcal{O}^{-1}(\bar{\Delta}_1 T^{\mathcal{O}^2} + \bar{\Delta}_2 T^{\mathcal{O}^1})}{1 + 2\beta(\bar{\Delta}_1 + \bar{\Delta}_2)},$$

show the correct high and low temperature behavior. For $\beta = 0 \Rightarrow \bar{y} = \overline{T}^{\mathcal{O}}$, and for $\beta \to \infty$ only the second terms with self–consistent solutions $\bar{y} = T^{\mathcal{O}^2}$ and $\bar{y} = T^{\mathcal{O}^2}$. For $\mathcal{O}^1 = \mathcal{O}^2$, this reads, similar to the case of the mixture model

$$\bar{y} = \frac{\sum_i z_i \overline{T}^{\mathcal{O}^i}}{z},$$

with

$$z_i = \frac{1}{2} + \beta\bar{\Delta}_i, \quad z = \sum_i z_i,$$

replacing $Z_i = e^{-\frac{\beta}{2}\tilde{\Delta}}$ and $Z$. However, this equation has a special usually not wanted feature. According to

its symmetry against exchanging $\bar{\Delta}_1 \leftrightarrow \bar{\Delta}_2$, the solution has either $\bar{\Delta}_1 = \bar{\Delta}_2$, i.e. $\bar{y} = \overline{T}^{\mathcal{O}}$, or for a solution $\bar{y}_1$ there exists another solution $\bar{y}_2$ with exchanged $\bar{\Delta}_i$.

Thus, as soon as the solution $\overline{T}^{\mathcal{O}}$ gets deformed, there exist always two of them with the same value of $L$, so there is no way to choose between them. In this sense the equation implements a model where for all given data (i.e. a posteriori) both components $\overline{T}^{\mathcal{O}^i}$ are equally likely. One can, for example, replace the term $\beta\bar{\Delta}_1\bar{\Delta}_2$ in $L^{PF}$ by $L^{I_1}$ or $L^{I_2}$ to 'enforce a decision'.

Because it has the simplest structure and at the same time shows the typical phase transition phenomenon we choose for the following numerical study the model $L^{I_1}$ for comparison with the mixture model, also with $\mathcal{O}^D = \lambda_D\mathcal{N}^D$ and $\mathcal{O}^s = \lambda_0\mathcal{O}_1 + \lambda_2\mathcal{O}_2 + \lambda_4\mathcal{O}_4$.

---

[83]In the context of regularization we want to fix the low temperature solutions and find possible parameterizations for all $\beta$, without being only interested in the neighborhood of the phase transition. Hence, we use a form for $L$ where the low temperature solutions, i.e. the templates (combined for every mixture component), can directly be read of, not however necessarily the corresponding critical $\beta^*$. Alternatively, we could also express $L$ in terms of the reduced temperature $t = 1/\beta - 1/\beta^*$, and choose polynomial terms in $\bar{y}$ to produce a phase transition at $t = 0$. This is more natural when studying phase transitions. For fourth order polynomials it is easy to solve for the extrema and therefore to relate the two formulations.

## 10.4 Bifurcations, phase transitions: one dimensional case

We have discussed the special two template case with $\mathcal{O}^1 = \mathcal{O}^2$. In that case the solutions of the mixture model $L^M$ are restricted to a one dimensional line in the function space $F^0$ of $\bar{y}_x$. Similarly, using a Landau–Ginzburg form for the log-likelihood, quartic effective interaction terms yields stationarity equations with at most two stable solutions. Hence, it may help in understanding the features of higher dimensional situations if we recall the well–known one–dimensional case. Therefore, we present some figures which gives a visually oriented summary over the bifurcation/phase transition behaviour for the models $L^M$ and $L^{I_1}$ in one dimension. This can also be seen as an illustration of the discussion of the maximum posterior approximation in Sections 7 and 8.5.

Figs.11 and 12 show that $L^M$ and $L^{I_1}$ indeed possess similar behavior, including a phase transition. There howevr are quantitative differences, especially the high and low temperature limits are not the same. At high temperature only one solution exists, decreasing the temperature a second solution can become stable. However, except for the case $a = 0$ the high temperature solution follow under decreasing temperature the more probable and thus better solution. This exemplifies the principle of *annealing techniques*.

In contrary, varying the value of $a$, corresponding to data, instead of temperature leads to typical hysteresis effects. Then a quite unlikely solution can remain stable for a long time. This may be seen as prototype for sequential updating or *on–line learning methods*.

Because $\gamma = \frac{\beta}{1+\beta}$ and not $\beta$ is the convex mixing coefficient between of data and template terms in $L^{I-1}$ the Figs.13 and 14 show the one dimensional case parameterized with $\gamma$. One sees that it can make a big difference to test different equal spaced values of either $\beta$ or $\gamma$ for example in cross–validation. The parameter $\gamma$ has the advantage of being completely in the interval $[0, 1]$.

Figs.15 and 16 compare the saddle point approximation with the full Bayesian approach depending on the distance of the data $a$ to the two templates $T_i = \pm 1$. The distances $|a - 1|/2$ and $|a + 1|/2$ correspond to the relative canonical distance $d_{\mathcal{O},T}(\bar{y}, \overline{T}^{\mathcal{O}'})$ and the value zero to the high temperature template average $\overline{T}^{\mathcal{O}}$. Clearly, for $a = 0$ the full Bayesian solution remains zero at all temperatures. For higher $a$ the mean field solution becomes quickly better. The, in contrast to the full risk, more pronounced structure shows the low temperature character of the MaP approximation, i.e. its tendency to single templates instead of their template average. Notice however, that the mean field solution despite resulting from an expansion in $1/\beta$ is not worst at high temperatures but in the neighborhood of the phase transition. The high temperature limit of the saddle–point approximation coincides, in the shown case of an approximation problem, with the exact solution again. This is the case because for a Gaussian distribution mean and mode coincide.
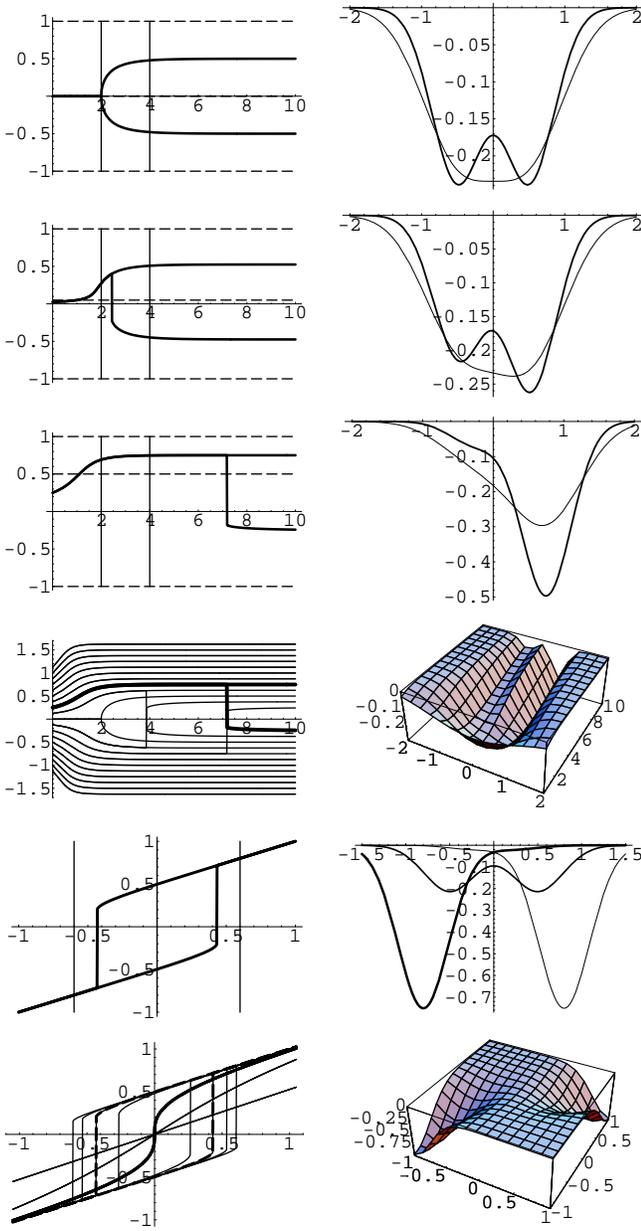
Figure 11: One–dimensional Gaussian mixture model: $N(a, 1/\sqrt{\beta})N(b, 1/\sqrt{\beta}) + N(a, 1/\sqrt{\beta})N(c, 1/\sqrt{\beta})$, at $b = -1$, $c = 1$, with $N(\mu, \sigma)$ denoting a Gaussian with mean $\mu$ and variance $\sigma^2$. We denote the independent variable by $\bar{y}$ to relate the example to the Bayesian framework. The variable $a$ is meant to represent data values, $b$ and $c$ templates. Rows 1– 3 (from top) left: $a = 0, 0.05, 0.5$; right: cut at $\beta = 2$ (10 times $f$) and $\beta = 4$; row 4, left: $\beta$ vs. $\bar{y}_{min}$ for $-2.25 \leq a \leq 2.25$ by 19 steps by 0.25 (thick line: $a = 0.5$); right: function at $a = 0.10$ for $0.5 \leq \beta \leq 10$; row 5, left: $\beta = 6$, $-1 \leq a \leq 1$; right: $\beta = 6$, $a = -0.6$(thick), 0, 0.6(thin); row 6, left: $a$ vs. $\bar{y}_{min}$ for $\beta = 0, 1, 2$(thick), 4, 6(dashed), 8, 10; right: function at $\beta = 6$ for $-1 \leq a \leq 1$.

Figure 12: One–dimensional Landau–Ginzburg regularization with a product term as effective interaction, representing a version of 'Fuzzy OR': $(\bar{y} - a)^2 + \beta(\bar{y} - b)^2(\bar{y} - c)^2$ at $b = -1$, $c = 1$. Rows 1– 3 (from top) left: $a = 0, 0.5, 1$; right: cut at 0.5 and 2.0; row 4 left: $\beta$ vs. $\bar{y}_{min}$ , $-1.5 \leq a \leq 1.5$ by 13 steps 0.25 (thick line $a = 0.5$); right: function at $a = 0.5$ for $0 \leq \beta \leq 3$; row 5: left: $-1 \leq a \leq 1$, $\beta = 1.2$; right: cuts at $\beta = 1.2$ for $a = -0.6$(thick) , 0.0, 0.6(thin); row 6: left: $a$ vs. $\bar{y}_{min}$ for $\beta = 0, 0.5$(thick), 1.0, 1.2(dashed), 1.5, 2.0, 2.5, 3.0; right: function at $\beta = 1.2$, $-8 \leq a \leq 8$.

Figure 13: One–dimensional Gaussian mixture model: parameterized by $\gamma = \frac{\beta}{1+\beta}$ taking values in $[0,1]$, corresponding to $\beta = \frac{\gamma}{1-\gamma}$. Rows 1–3 (from top) left: $a = 0, 0.5, 1$; right: cut at $\gamma = 2/3$ and 0.85; row 4: left: $\gamma$ vs. $\bar{y}_{min}$, $-1.5 \leq a \leq 1.5$ by 13 steps 0.25 (thick line $a = 0.5$); right: function at $a = 0.1$ for $0.12 \leq \gamma \leq 0.4$; row 5: left: $-1 \leq a \leq 1$, $\gamma = 0.15$; right: cuts at $\gamma = 0.15$ for $a = -0.5$(thick), 0.0, 0.5(thin); row 6: left: $a$ vs. $\bar{y}_{min}$ for $0 \leq \gamma \leq 1$ by steps of 0.2 (thick: $\gamma = 2/3$, $\gamma = 0.85$); right: function at $\gamma = 0.85$, $-1 \leq a \leq 1$.

Figure 14: One–dimensional Landau–Ginzburg regularization with a 'Fuzzy OR' in its 'natural' convex parameterization $\gamma = \frac{\beta}{1+\beta}$ taking values in $[0,1]$, so that $\beta = \frac{\gamma}{1-\gamma}$. Rows 1–3: (from top) left: $a = 0, 0.5, 1$; right: cut at $\gamma = 1/3$ and $2/3$; row 4: left: $\gamma$ vs. $\bar{y}_{min}$, $-1.5 \leq a \leq 1.5$ by 13 steps 0.25 (thick line $a = 0.5$); right: function at $a = 0.5$ for $0 \leq \gamma \leq 11$; row 5: left: $-1 \leq a \leq 1$, $\gamma = 0.5$; right: cuts at $\gamma = 0.5$ for $a = -0.5$(thick), -0.0, 0.5(thin); row 6: left: $a$ vs. $\bar{y}_{min}$ for $\gamma = 0, 0.25, 1/6$(thick), 0.5(dashed), 0.75, 1; right: function at $\gamma = 0.5$, $-1 \leq a \leq 1$.

Figure 15: Mean field (Maximum posterior approximation or empirical risk minimization) vs. full Bayesian approach for a Gaussian mixture model at $a = 0.1$. Row 1: The posterior probability $p(f^0|f)$ used for the MaP step (left), and the corresponding optimal MaP–solution $\bar{y}^* = f^{0,*} = \mathrm{argmax}_{f^0 \in F^0} p(f^0|f)$ (right). Row 2: The true effective probability $p(y|f) = \int df^0\, p(f^0|f) p(y|f^0)$ (left) and its maximal value $\mathrm{argmax}_y p(y|f)$ (right). Row 3: The full Bayesian risk for the corresponding approximation problem (so we can identify $f^0$ and $\hat{f}$) $r(\hat{f}, f) = -\int df^0 \int dy\, p(f^0|f) p(y|f^0) \ln p(y|\hat{f})$ (left) and its minimal value $\hat{f}^* = \mathrm{argmin}_{\hat{f} \in \hat{F}} r(\hat{f}, f)$ (right). Row 4: Show all three curves combined for comparison (right) and on the left the actual loss distribution $-p(y|f) \ln p(y|\hat{f})$ for the example $\hat{f} = \bar{y} = 0$. We may remark that even in this case where the mean field approximation is no good approximation for the true $\beta$–dependency of the risk because of the small $a$ (for $a=0$ the true optimal solution would always be zero), it is nevertheless possible to obtain the whole range by adapting the 'mean field temperature', for example by cross–validation. Notice that neither the linear high temperature regularization nor the two linear low temperature limits can access the whole range as they can never cross the value $a$ by changing their $\beta$.
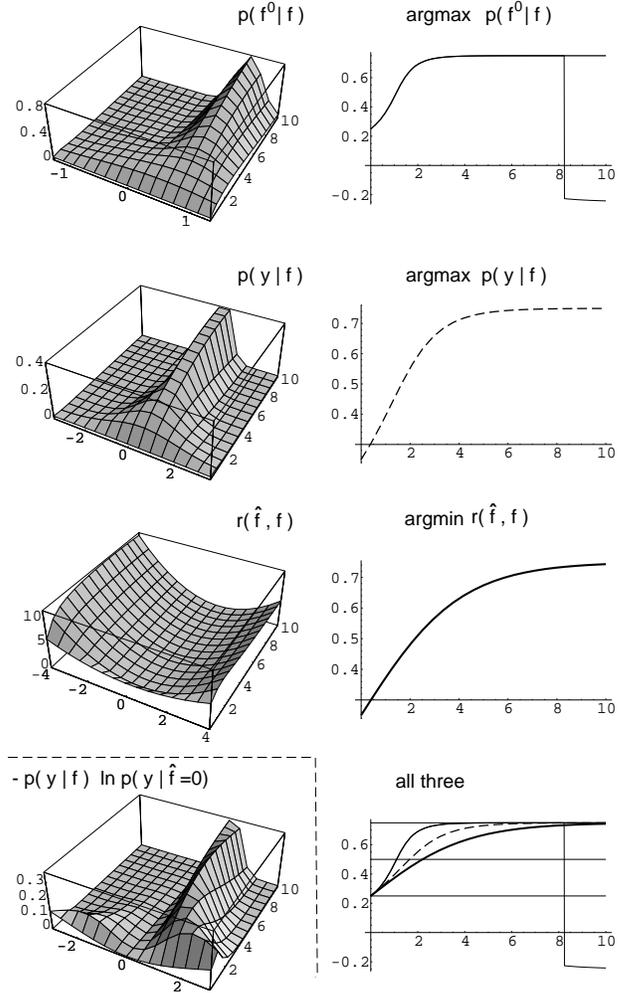
Figure 16: Mean field (Maximum posterior) vs. full Bayesian approach at a=0.5. The same situation as in Fig.15 with $a$ nearer to the 1–template, so the mean field approximation becomes better. (For still larger $a$ (not shown) the mean field approximation improves quickly.) In both figures the low temperature character of the maximum posterior method is nicely seen. The posterior probability is much sharper peaked than the true risk, amplifying therefore differences between alternative $f^0$. The true risk, containing two integrations, is much smoother. The fact that the maximum of the true $y$ distribution $p(y|f)$ in state $f$ does not coincide with the optimal $\hat{f}$ shows the asymmetry of this distribution. Obviously, the mean field approximation is much better for larger $a$. For a non–approximation problem the risk minimization under $f^{0,*}$ would have to be included in the MaP–MiR procedure. The results depend from the chosen non–approximation loss. One may remark here, that in situations where the template represents a prototypical situation for which actions are available and cheap, it is reasonable to add a loss term increasing with the distance from the nearest template. Including such a 'template-distance' loss favors a low temperature approximation and improves the validity of the mean field solution.
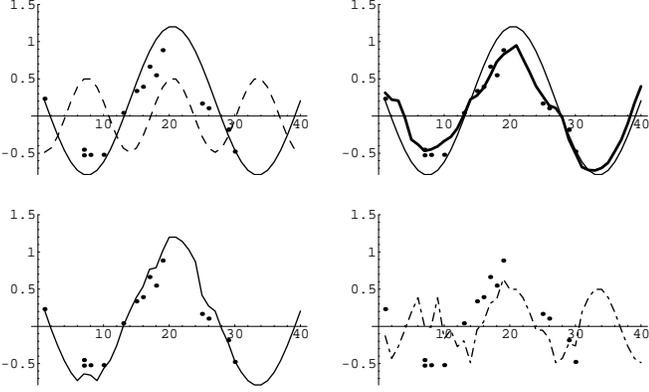
Figure 17: The two template example. The upper left diagram shows the two templates $T^1$ and $T^2$ and data (drawn from the interval $[1, 30]$). The upper right diagram shows the state of nature $f^0$ (thickly dashed) and $T^1$. The second row from above shows the two ($\lambda_D$-, $\lambda_0$-dependent) $\overline{T}^{\mathcal{O}^i}_{\lambda_2 = \lambda_4 = 0}$ which are the solutions for either $T_1$ or $T_2$ combined with the data $\bar{D}$ under vanishing smoothness coefficients $\lambda_2 = \lambda_4 = 0$. They are in the following figures given as reference to estimate the effect of smoothing. (Left: for $T_1$, right: for $T_2$.)

.

## 10.5 Numerical results

As examples of templates we choose (see the two dashed curves in the upper left picture in Fig.17)

$$T_1 = -\sin\left(\frac{3\pi(x-1)}{m-1}\right) - a_1,$$

$$T_2 = \sin^2\left(\frac{3\pi(x-1)}{m-1}\right) - a_2,$$

with $m = 40$ and $a_i$ adjusted so that both functions have mean zero on the interval $[1, 40]$. We consider the case that the learner expects the actual function to be similar, but not identical, to either $T_1$ OR $T_2$. Thus, the templates represent function prototypes. They may stand for two typical structures for a time series or, in case of an incomplete (here one–dimensional) image to be reconstructed for two expected spatial patterns.

A mixture model can be easily realized by a hierarchical sampling process. Then firstly a mixture component is chosen corresponding to one of the templates and representing disjunct events. In a second step the actual $f^0$ is generated from that mixture component. An interaction model may be sampled by Monte–Carlo methods. We do not intend to generate $f^0$ exactly according to the one of the learning models. Instead we generate the state of nature $f^0$ by a different hierachical process. Specifically, we use the following method to generate $f^0$: Firstly, we choose one $T^i$ ($T^1$ in the below examples which, however, is assumed not to be known by the learner) and add Gaussian noise (with $\sigma = 0.2$) for every $x$. In a second step this wiggly function is smoothed by
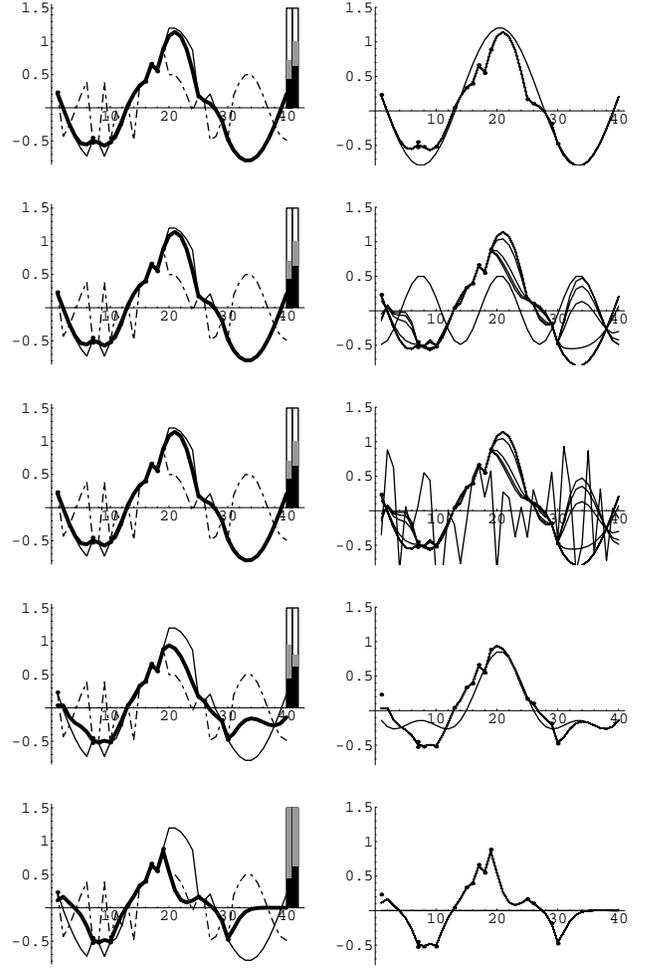


Figure 18: Mixture model at $\beta = 1 < \beta^*$. ($L^M$, relaxation, $\eta = 1.0$, at $\lambda_D = 100$, $\lambda_0 = 1$, $\lambda_2 = 1$, $\lambda_4 = 1$). The first three rows, here and in the following figures, shows the results of relaxation learning for the mixture model $L^M$ according to (43, 44) with $\mathcal{A} = \mathcal{O}$ for different starting configurations $\bar{y}^0$, which are (top to bottom) ($\bar{y}^0_1 = T_1$, $\bar{y}^0_2 = T_2$, and a random $\bar{y}^0_r = T_{random}$). The last two rows show for comparison two one-template models with the same choice of parameters and starting point $\bar{y}^0 = T$: Row 4: a mixture template $T = \overline{T} = (T_1 + T_2)/2$, row 5: the usual zero template $T = T_0 = 0$ of homogeneous linear regularization. Here and in the corresponding following figures, the diagrams on the right show the evolution of the solution $\hat{y}$ during iteration, and the diagrams on the left the final solutions (thick line). For comparison also shown are, data (points) and the two templates $\overline{T}^{\mathcal{O}^i}_{\lambda_2 = \lambda_4 = 0}$ (see Fig.17) for the given $\lambda_D$, $\lambda_0$. The bars show the (mean square) generalization error (gray) calculated for 1000 newly generated random points in the intervals $[1, 40]$ (left) or $[1, 30]$ (right), respectively, and, here and in the following, always normalized with respect to the largest of the errors for all five cases ($L^M$ with $\bar{y}^0 = T_1$, $T_2$, $T_{random}$, and linear regularization with $\overline{T}$, $T_0$) under the same parameter combinations. The black part represents the minimal possible generalization error, with absolute value always equal 0.04.
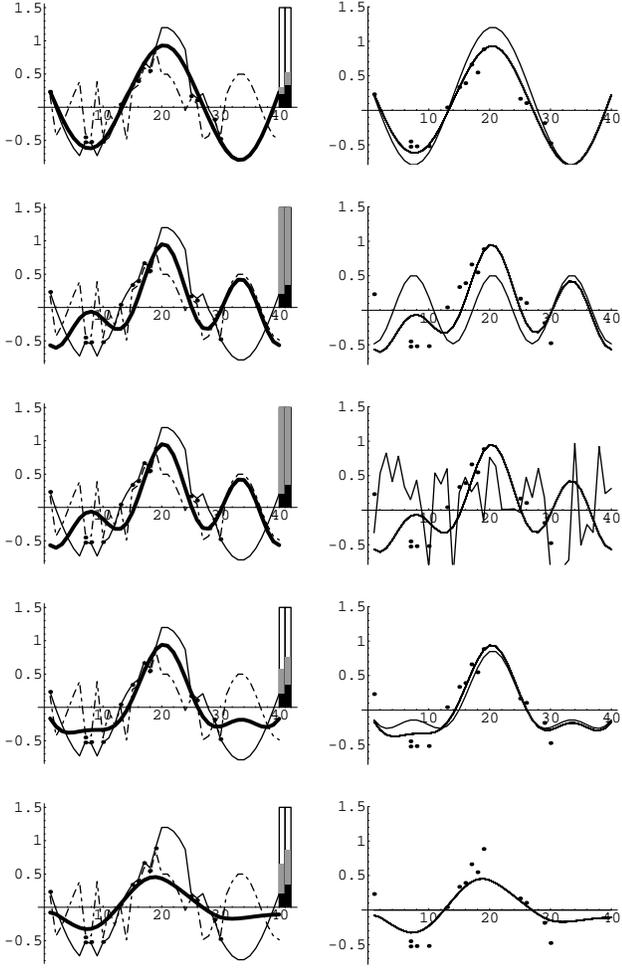
Figure 19: Smoother model, a low temperature case $\beta = 1 > \beta^*$. (Mixture model $L^M$, relaxation, $\eta = 1.0$, 50 iterations, $\lambda_D = 10$, $\lambda_0 = 1$, $\lambda_2 = \alpha^2$, $\lambda_4 = \alpha^4$, with $\alpha = (m-1)/(3\pi)$, m=40 bringing in this case all derivatives in the same order of magnitude.) Rows 1-3: mixture model, starting configuration $T_1$, $T_2$, $T_{random}$ (top to bottom). Rows 4 and 5: one template models with $\overline{T} = (T_1 + T_2)/2$ (row 4) and $T_0 = 0$ (row 5). Because one-template models result for the relaxation method with $\eta = 1$ in a linear equations, only one iteration is needed to obtain the final solution. For $\eta = 1$ the number of iteration steps needed to converge to the final solution can be seen as a measure of the 'nonlinearity' of the equations. For example, the left hand side figures (rows 1–3) show that after only one iteration the solutions hardly change anymore and thus the equations of the mixture model in this parameter range (in contrast to other situations) are nearly linear.
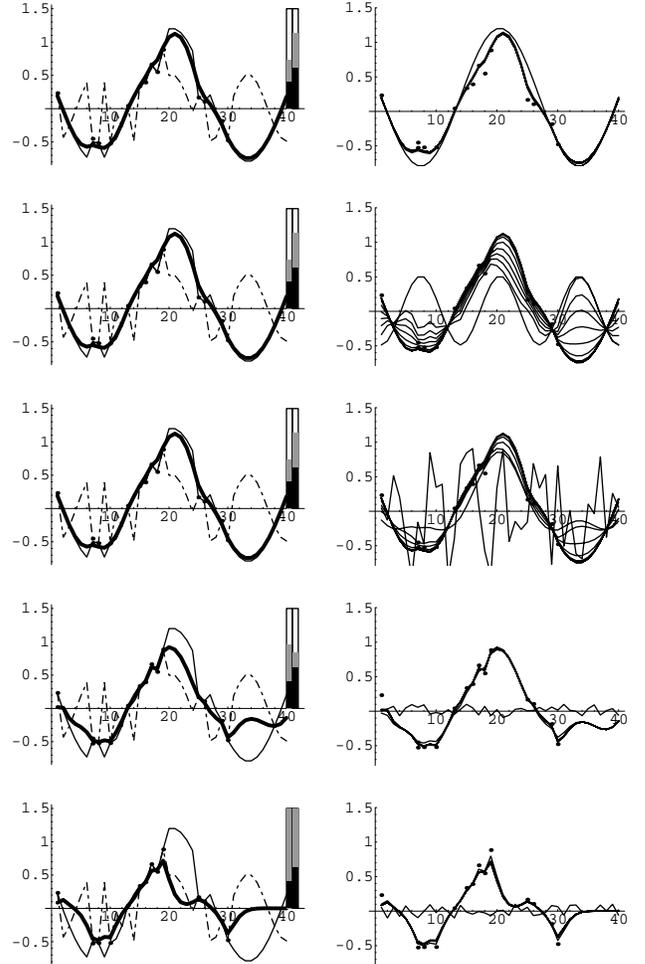
Figure 20: Landau–Ginzburg Regularization. ($L^{I_1}$, relaxation, $\eta = 1.0$, at $\beta = 1$, $\lambda_D = 100$, and $\lambda_0 = 1$, $\lambda_2 = 1$ $\lambda_4 = 1$.) Like for the mixture model rows 1–3 show the solutions evolving from the different starting configurations $T_1$, $T_2$, $T_{random}$. Shown is a high temperature case where only one solution is stable. The figure shows clearly how the nonlinearities of the mean field equation forces the solution $\bar{y}_2$ evolving from $T_2$ (and $T_{random}$) towards the solution $\bar{y}_1$ evolving from $T_1$ (rows 2, 3). For this solution (row 1), being already near the extremum, the nonlinearities are not effective. One sees that also the one template models $T_0$, $\overline{T}$ are nonlinear.
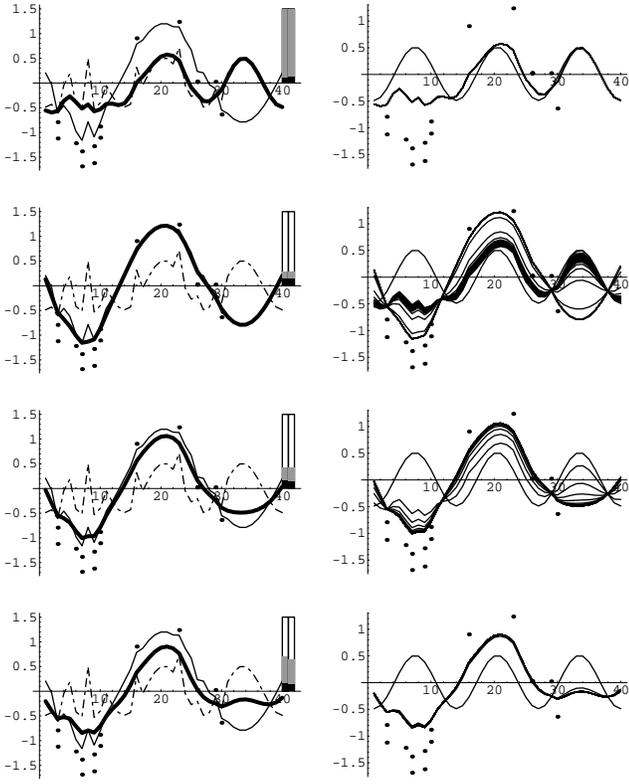
Figure 21: Bifurcation: un– and metastable states. (Mixture model $L^M$, relaxation, $\eta = 1$, $\lambda_D = 1$, $\lambda_0 = 1$, $\lambda_2 = 1$, $\lambda_4 = 1$, $\beta = 1$, 0.482, 0.1, 0.01, top to bottom.) The weaker solution evolving from $T^2$ changes suddenly with $\beta$, with a vanishing gradient at $\beta = \beta^*$. This solution appears as 'shadow' in the iteration picture, and looks stable under a smaller number of iterations. The solution is near the phase transition strongly adapted to the data and quite different from its starting point $T^2$.
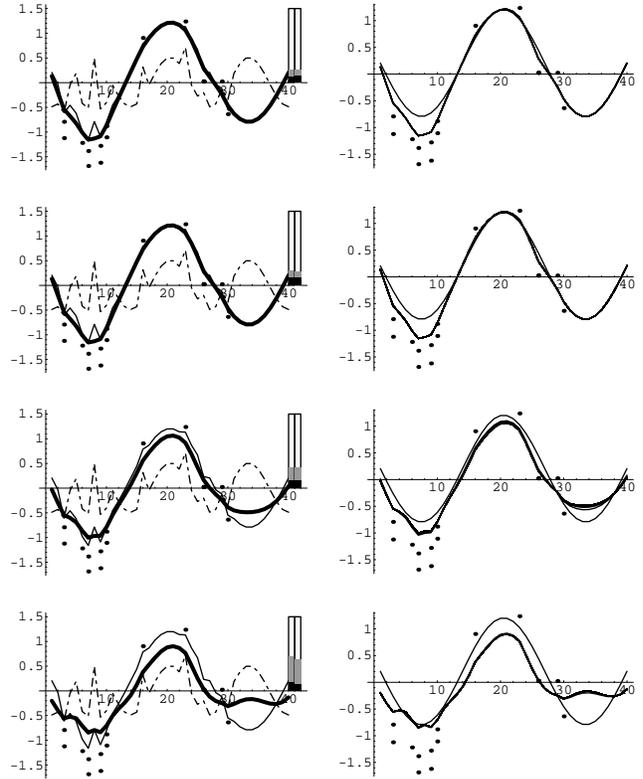
Figure 22: Bifurcation: The stable state. (Mixture model $L^M$, relaxation, $\eta = 1$, $\lambda_D = 1$, $\lambda_0 = 1$, $\lambda_2 = 1$, $\lambda_4 = 1$, $\beta = 1$, 0.482, 0.1, 0.01, top to bottom.) The better solution near $T^1$ remains nearly unchanged. It is also near the phase transition still quite similar to its starting point $T^1$ as the data do not require a strong adaption like for the weaker solution. The deformation becomes larger in the high temperature limit.
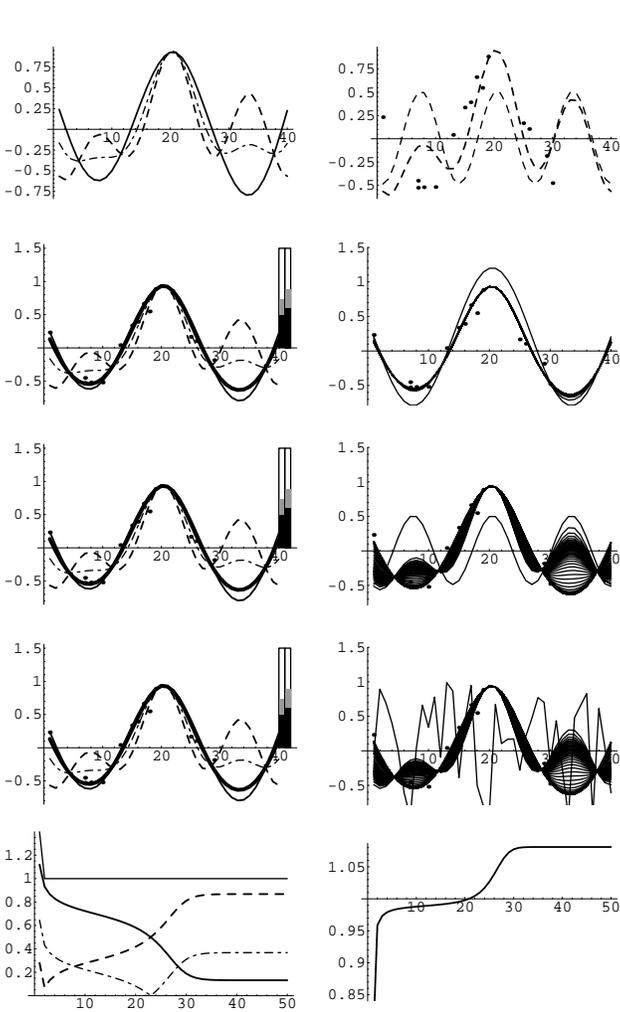
Figure 23: Phase transition for a smooth model (Mixture model $L^M$, relaxation, $\eta = 1.0$, 50 iterations, at $\beta = 0.0105 \approx \beta^*$, $\lambda_D = 10$, $\lambda_0 = 1$, $\lambda_2 = \alpha^2$, $\lambda_4 = \alpha^4$.) Row 1, left: Shown are the two low temperature solutions $\overline{T}^{\mathcal{O}^1}$, $\overline{T}^{\mathcal{O}^1}$ (dashed), and the high temperature limit $\overline{T}^{\mathcal{O}}$ (dot–dashed) in the middle. (Which in this case is similar, but not identical to $\overline{T} = T_1 + T_2$.) Row 1, right: $\overline{T}^{\mathcal{O}^2}$ (thickly, dashed) is shown resulting from the data points and $T_2$ (thinly, dashed). Rows 2-4: (Starting configurations $T_1$, $T_2$, $T_{random}$.) The figure shows that the solution evolving from $T_2$ (and in this case also from $T_{random}$) is nearly stable. The 'shadow' in the right hand figure still shows the corresponding low temperature solution (compare Fig.19.) The amount of iterations needed reflects the high nonlinearity of the mean–field equations at this point. The (linear) one–template models are temperature independent and here not shown (but in Fig.19). The two plots in row 5 are explained in the caption of Fig.24
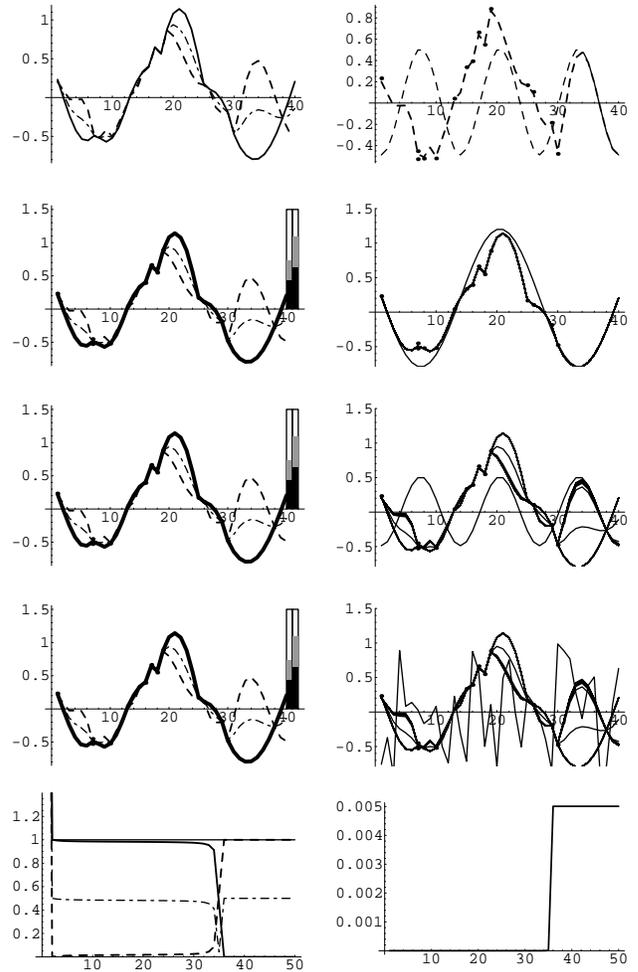
Figure 24: Phase transition for a more data oriented model (Mixture model $L^M$, relaxation, $\eta = 1.0$, 50 iterations, at $\beta = 4.85 \approx \beta^*$, $\lambda_D = 100$, $\lambda_0 = 1$, $\lambda_2 = 1$, $\lambda_4 = 1$.) Plots in rows 1-4 correspond the Fig.23. One sees clearly, that the transition is much sharper, as the higher data orientation of the coefficients favors the solution evolving from $T_1$. Indeed, the solution evolving from $T_2$ seems perfectly stable before their sudden transition. Row 5 left: Shown are the normalized canonical distances (See Eqs.(38, 37)) $d_1 = d_{\mathcal{O},T}(\bar{y} - \overline{T}^{\mathcal{O}^1})$ (thick), $d_2 = d_{\mathcal{O},T}(\bar{y} - \overline{T}^{\mathcal{O}^2})$ (dashed), $d_{HT} = d_{\mathcal{O},T}(\bar{y} - \overline{T}^{\mathcal{O}})$ (dot–dashed), for $\bar{y} = \bar{y}_2$ with starting configuration $\bar{y}_2^0 = T_2$ during iteration. Points with $d_1 + d_2 = 1$ are according to the triangle equality exactly on the line spanned by convex combinations of the two low temperature states, which are solutions of the corresponding limiting linear regularization problems. One sees that the final solutions are on this line, however not the starting configurations which may be anywhere in the high-dimensional space $F^0$. Notice also, that $\bar{y}$ iterates along this line and passes the possible solutions in between. This can be used as a sanity check for numerical calculations. Additionally restricted $\bar{y}$, e.g. with periodic boundary conditions, have in general not $d_1 + d_2 = 1$. Right: The second plot in row 5 shows the increasing $L^M(\bar{y})$ (unnormalized).
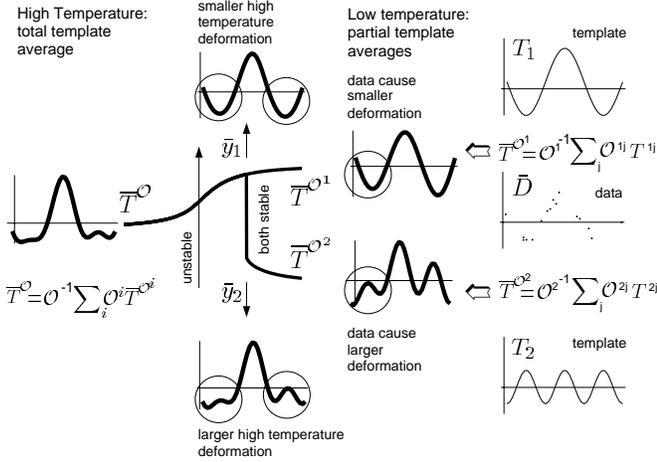
Figure 25: The diagram summarizes schematically the two template example in this Section and the temperature dependence of its solutions $\bar{y}$. Variations of high and low temperature solutions under changing parameters are shown in Figs.26, 27.
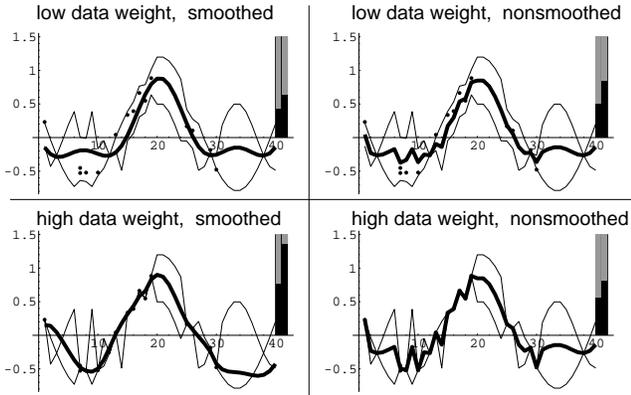


Figure 26: High temperature solutions $\overline{T}^{\mathcal{O}}$. (Mixture model $L^M$, relaxation, $\eta = 1.0$.) The high temperature solutions calculated at $\beta = 0.001$ and (top left: small data and large smoothness influence) $\lambda_D = 1$, $\lambda_0 = 1, \lambda_2 = \alpha^2, \lambda_4 = \alpha^4$ (top right: small data and no smoothness influence) $\lambda_D = 1$, $\lambda_0 = 10$, $\lambda_2 = 0$, $\lambda_4 = 0$ (bottom left: large data and smoothness influence) $\lambda_D = 1000$, $\lambda_0 = 1, \lambda_2 = \alpha^2, \lambda_4 = \alpha^4$ (bottom right: large data and no smoothness influence) $\lambda_D = 1000$, $\lambda_0 = 1$, $\lambda_2 = 0$, $\lambda_4 = 0$
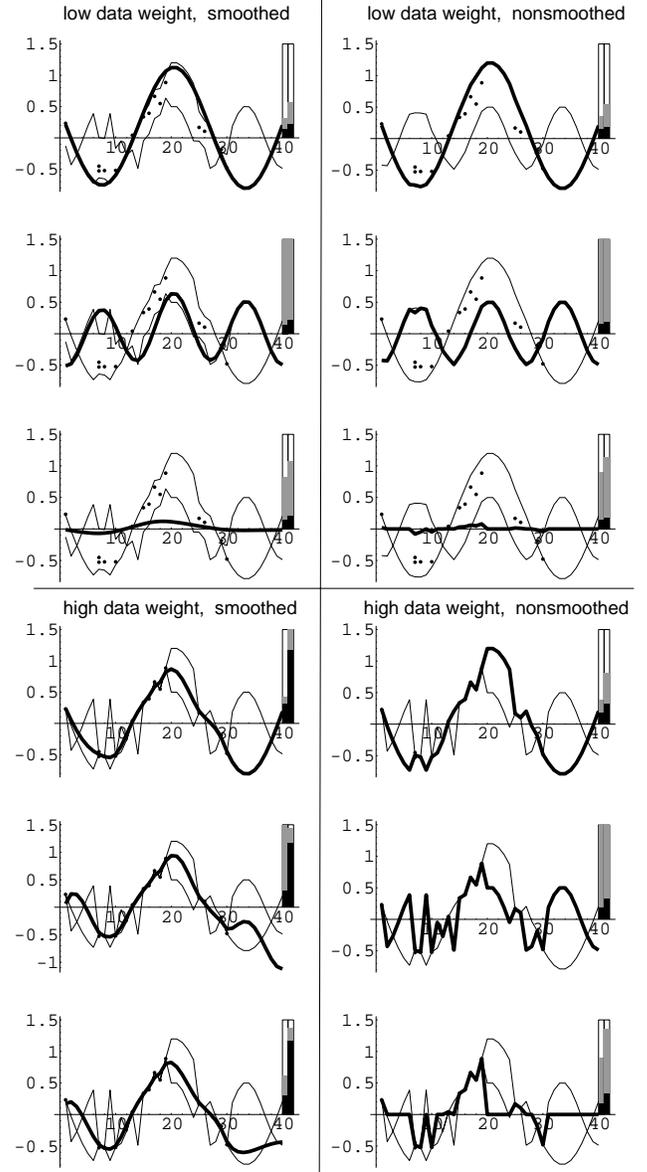.



Figure 27: Low temperature solutions $\overline{T}^{\mathcal{O}^i}$. (Mixture model $L^M$, relaxation, $\eta = 1.0$.) The low temperature solutions calculated $\beta$–independent for the three one template models with $T_1$, $T_2$, and $T_0$ (from top to bottom, with error bars relative within the same parameter values) at the four cases corresponding to Fig.26, (top left) $\lambda_D = 1$, $\lambda_0 = 1, \lambda_2 = \alpha^2, \lambda_4 = \alpha^4$, (top right) $\lambda_D = 1$, $\lambda_0 = 10$, $\lambda_2 = 0$, $\lambda_4 = 0$, (bottom left) $\lambda_D = 1000$, $\lambda_0 = 1, \lambda_2 = \alpha^2, \lambda_4 = \alpha^4$, (bottom right) $\lambda_D = 1000$, $\lambda_0 = 1$, $\lambda_2 = 0$, $\lambda_4 = 0$. One can observe, that solutions evolving from better fitting templates are able to produce smoother functions $\bar{y}$.
.

85

changing $f^0$ randomly (Zero mean Gaussian mutations with $\sigma = 0.02$) at all locations and accepting the change if the smoothness increases. Smoothness is hereby measured by $< \bar{y}^0 | \mathcal{O}^s | \bar{y}^0 >$ with $\lambda_0^{gen} = 0.0$ (i.e. only derivatives of $T^1$ contribute), $\lambda_2^{gen} = 0.5\alpha^2 \approx 8.56$ and $\lambda_4^{gen} = 0.5\alpha^4 \approx 146.6$, with $\alpha = (m-1)/(3\pi)^2 \approx 4.138$ so the derivatives have the same order of magnitude. (In the following the learning models do not have the same coefficients as the generation model, i.e. $\lambda_i \neq \lambda_i^{gen}$.) This smoothing process has been iterated 2000 times. Then data are drawn from $f^0$ with a Gaussian distribution with $\sigma = 0.2$ and mean $\bar{y}_x(f^0)$ from the interval $[1, 30]$. Thus, the task can be seen as a simple two–template prediction or reconstructing problem. with the intervall $[31, 40]$ representing either future values (for time series) or a hidden area (in image reconstruction). See the thickly dashed curve in the upper right picture in Fig.17 for the $f^0$ used for the results discussed in the following.

Figs. 18 – 20 present numerical results for the two–template example and the two prototypical nonlinear regularization methods:

1. the *finite temperature regularization* with mixture model

$$L^M \propto e^{-\frac{1}{2}(\Delta_D + \Delta_1)} + e^{-\frac{1}{2}(\Delta_D + \Delta_2)},$$

and stationarity equations (43, 44)

2. the *Landau–Ginzburg regularization* with an interaction model

$$L^{I_1} = g(\gamma_D \Delta_D + \gamma \Delta_1 \Delta_2).$$

in a naive fuzzy version and stationarity equations (47, 48)

The two iteration schemes from the spectrum of learning algorithms we used are

A. relaxation with $\mathcal{A} = \mathcal{O}$

$$\bar{y}^{i+1} = (1 - \eta)\bar{y}^i + \eta \mathcal{O}^{-1} t,$$

B. and the gradient,

$$\bar{y}^{i+1} = (1 - \eta)\bar{y}^i + \eta(\mathcal{O}\bar{y}^i - t).$$

Intermediate algorithms can for example invert lower dimensional sub-blocks of $\mathcal{O}$. Such a submatrix of $\mathcal{O}$ can be constructed by including from clusters of correlated variables one or a few (prototypical) representatives. Multigrid methods, for example, can be seen as such an approach for functions with approximately homogeneous local correlations. Note that in our case the EM algorithm coincides with the relaxation algorithm. Fig. 29 shows typical problems of the gradient method for (discretized) differential operators.

Fig. 25 summarizes the typical bifurcation or phase transition behaviour for the mixture model.

# 11 Conclusions

The paper is motivated by the fact that *predetermined dependencies* between answers to different questions are necessary for generalization and, thus, responsible for
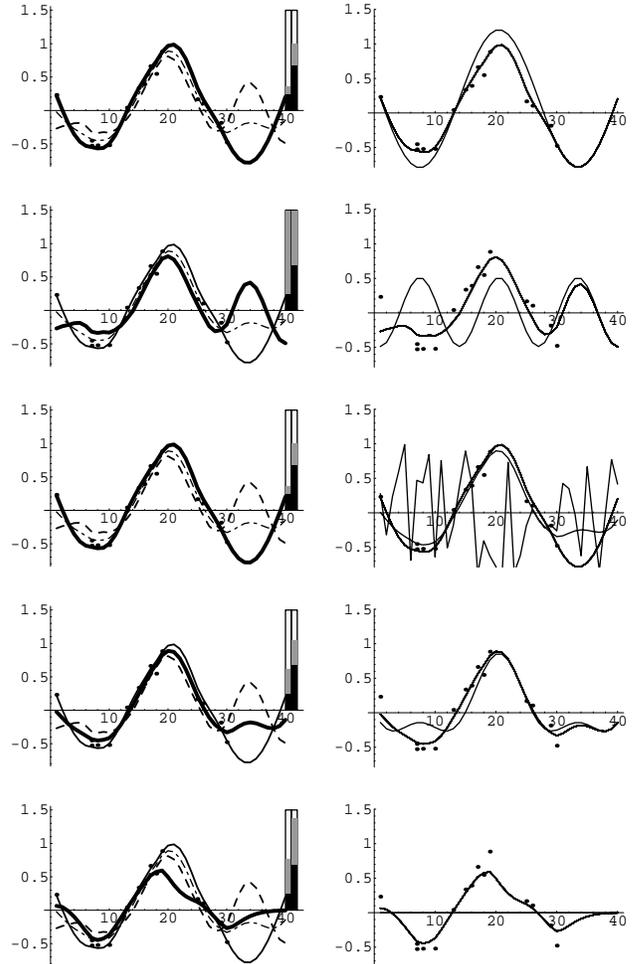


Figure 28: RBF (Radial Basis Function) regularizer. (Mixture model $L^M$, relaxation, $\eta = 1.0$, at $\beta = 1$, $\lambda_D = 10$, and $\lambda_0 = 1$, $\lambda_2 = \sigma_{RBF}^2/2$, $\lambda_4 = \sigma_{RBF}^4/(2!2^2)$.) The $\lambda_i$ are chosen as the first three coefficients of the RBF regularizer (29) with $\sigma_{RBF} = 3$. For the zero template $T_0$, corresponding to the usual linear RBF method, one sees clearly the superposition of Gaussian-like functions centered at the data points.
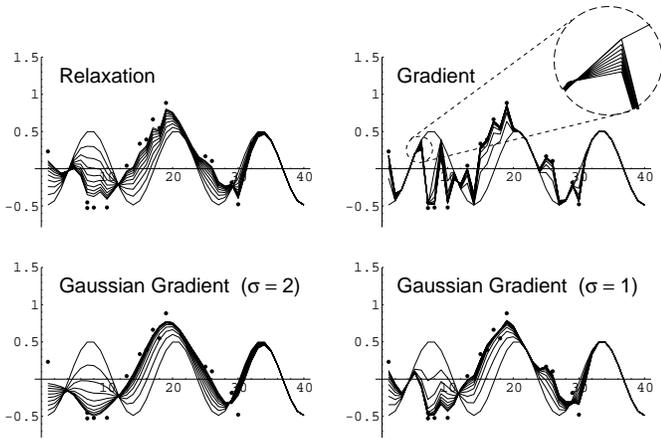
Figure 29: The "Gradient–$\delta$–catastrophe": Row 1, left: Relaxation with full inversion $\mathcal{A}^{-1} = \mathbf{O}^{-1}$, $\eta = 0.2$ (the relaxation method converges with $\eta = 1$ in nearly one step), Row 1, right: Gradient, i.e. $\mathcal{A}^{-1} = \mathcal{I}$ (=Identity) and $\eta = 0.005$. Row 2: A 'Gaussian Gradient' with $\mathcal{A}_{x,x'}^{-1} = (1/\sigma\sqrt{2\pi})e^{-(x-x')^2/(2\sigma^2)}$ for $\sigma = 2$ (left) and $\sigma = 1$ (right) and $\eta = 0.005$. Shown are the first 10 iterations for a situation with high data influence: $\beta = 1$, $\lambda_D = 100$, $\lambda_0 = 1$, $\lambda_2 = 1$, $\lambda_4 = 1$. To obtain an iteration procedure, $\mathbf{O}$ is splitted into two parts $\mathbf{O} = \mathcal{A} + \mathcal{B}$ and a new guess $\bar{y}^{i+1}$ is generated according $\bar{y}^{i+1} = \bar{y}^i - \eta\mathcal{A}^{-1}(\mathbf{O}\bar{y}^i - \mathbf{t})$. The gradient algorithm takes $\mathcal{A} = \mathcal{I}$ equal to the identity matrix. For an operator $\mathbf{O}$ with only diagonal (e.g. data terms) and near diagonal (e.g. differences for a discretized differential operator) matrix elements the update information propagates, besides through the $x$–independent factors $Z_i$ contained in $\mathcal{O}$ and $\sqcup$, only locally between neighbouring $x$. In the limit of continuous $x$ this becomes for differential operators, i.e. a vanishing effective neighborhood, arbitrarily slow. Analogous problems arise if $\mathcal{A}^{-1}$ has an (approximate) block structure. Hence, in practice when $\mathbf{O}$ is no integral operator the update algorithm has to use an $\mathcal{A}^{-1}$ which connects dependent parts of $\bar{y}$. This can in our case be achieved by enlarging the neighborhood for example by choosing (at least for the beginning) a Gaussian $\mathcal{A}^{-1}$ or other forms of local means. The nonlocality may be further increased by using nonlocal or hierarchically organized blocks ('neighborhoods'), like for example multigrid methods. One may also write $\bar{y} = \Delta\bar{y} - \tilde{y}$ with $\tilde{y}_x = \tilde{y}(x, w)$ an approximation of $\bar{y}$ with nonlocal dependencies, and update first $\tilde{y}(x, w)$ and then $\Delta\bar{y}$. In case $\tilde{y}$ approximates the nonlocal dependencies quite well a gradient algorithm (e.g. backpropagation for a neural net) with respect to the parameters $w$ may already nicely converge so adaption of the remaining $\Delta\bar{y}$ is fast enough. Also, the $\overline{T}_{\lambda_2 = \lambda_4 = 0}^{\mathcal{O}^i}$ or if available as result of a first linear regularization step $\overline{T}^{\mathcal{O}^i}$ might be good starting configurations. In general, the relevant nonlocal dependencies can be different in the various low and high temperature regimes.

learning. Standard training examples alone can never lead to any prediction for new data. It seems therefore necessary to concentrate more on informations about the dependencies than it is usually done.

The aim of the paper is to treat those dependencies as explicit as possible, and to discuss possibilities to base information about dependencies upon measurement and control. This is especially important if the objects/situations of interest are complex and/or the amount of available standard training data is small.

For local questions, the predetermined dependencies have the form of stationarity conditions for the answer generating probability distributions. To enable generalization answers to *nonlocal questions* must be available, with the definitions of nonlocal questions representing predetermined dependencies. Commonly implemented nonlocal dependencies correspond to bounds on smoothness or other symmetries. Many forms of prior information can be available in practice. Often they appear as implicit or linguistic concepts defining for example object classes like faces, chairs, pedestrians or cars, and are, because difficult to formalize, not included in the learning algorithm. Such dependencies can be implemented by an interface using *fuzzy priors*. Nonlocal questions are usually not sampled and not directly included in the loss function.

Often *non–approximation aspects* are of interest, like the amount of resources (time, memory, money, understandability, complexity) needed by available alternatives $\hat{f}$. Then empirical risk minimization cannot be interpreted as being equivalent to a Bayesian maximum posterior approximation but can be extended to a two–step procedure (MaP–MiR). Priors often depend formally from an infinite number of function values. Such priors can be implemented by using for the preparation process or for definition and control of the situation under study measuring devices different from those for the training process.

Priors stating that a function is probably similar to a template $T^1$ OR to a template $T^2$ can be implemented by some 'soft OR' or mixture model. This leads for the maximum posterior approximation to stationarity equations with *nonlinear* dependence from the local function values, reflecting the nontrivial interactions between different locations. As nonlinear equations have in general multiple solutions. such priors can, for example, be used to model phenomena like ambiguous illusions in perception.

The results of learning are statements about decision relevant data assuming their dependency on available data. Hence, learning is a reformulation of knowledge, and consists of

1. the algorithmic problem of extracting decision relevant information from given knowledge, usually a list of data and required dependencies, and

2. the (empirical) validity problem of controlling or identifying the relating dependencies.

Consequently, control (e.g. identification of situations appropriate for generalization) is needed to relate the results of past measurements to situations for which learn-

ing is intended, and the ability to generalize is intimately related to the ability to control and compute the required dependencies in (a finite number of) application situations. For example, stationarity of data generating processes and attributes of measurement devices, or in more biological terms of the sensory input, like limiting bounds and averaging processes, have to be established and controlled to guarantee smoothness and other approximate symmetries. Summarizing the active interpretation of learning we can say: *Generalization is control or identification of decision relevant dependencies.*

# References

[1] Aarts, E. & Korts, J. (1989) Simulated Annealing and Boltzmann Machines. New York: Wiley.

[2] Abu–Mostafa, Y. (1990) Learning from Hints in Neural Networks. *Journal of Complexity* **6**, 192–198.

[3] Abu–Mostafa, Y. (1993a) Hints and the VC Dimension. *Neural Computation* **5**, 278–288.

[4] Abu–Mostafa, Y. (1993b) A method for learning from hints. *Advances in Neural Information Processing Systems* **5**, S. Hanson et al (eds), 73–80, San Mateo, CA: Morgan Kauffmann.

[5] Ackley, D.H., Hinton, G.E., & Sejnowski, T.J. (1985) A Learning Algorithm for Boltzmann machines. *Cognitive science* **9**, 147–169.

[6] Alder, B.J. & Wainwright, T.E. (1959) Studies in molecular dynamics. I. General method. *Journal of Chemical Physics* **31**,459–466.

[7] Allgower, E.L. & Georg, K. (1990) Numerical Continuation Methods, Berlin: Springer–Verlag.

[8] Allen, D.M. (1974) The relationship between variable selection and data augmentation and a method of prediction. *Technometrics* **16**, 125.

[9] Amari, S. (1985) Differential geometrical methods in statistics. Springer Lecture Notes in Statistics, 28, Springer–Verlag.

[10] Amari, S. (1995) Information Geometry of the EM and em Algorithms for Neural Networks. *Neural Networks* **8** (9), 1379–1408.

[11] Amit, D., Gutfreund, H., & Sompolinski, H. (1987) Statistical Mechanics of Neural Networks Near Saturation. *Annals of Physics* **173**, 30–67.

[12] Amit, D. (1989) Modeling Brain Function. Cambridge: Cambridge University Press.

[13] Austern, N. (1970) Direct Nuclear Reactions. New York: Wiley.

[14] Balasubramanian, V. (1996) Statistical Inference, Occam's Razor and Statistical Mechanics on the Space of Probability Distributions *cond-mat/9601030*.

[15] Balian, R. (1991) From Microphysics to Macrophysics. Vol. I. Berlin: Springer.

[16] Barber, D. & Williams, C.K.I. (1997) Gaussian processes for Bayesian classification via hybrid Monte Carlo. To appear in *Advances in Neural Information Processing Systems 9.*

[17] Barndorff–Nielsen, O.E. (1978) Information and exponential families in statistical theory. New York: Wiley.

[18] Bazaraa, M.S., Sherali, H.D., & Shetty, C.M. (1993) Nonlinear Programming. (2nd ed.) New York: Wiley.

[19] Beck, C. & Schlögl, F. (1993) Thermodynamics of chaotic systems. Cambridge University Press.

[20] Berger, J.O. (1985) Statistical Decision Theory and Bayesian Analysis. (Second Ed.),New York: Springer.

[21] Bertsekas, D.P. (1995) Nonlinear Programming. Belmont, MA: Athena Scientific.

[22] Berry, M.V. (1966) *Proc. Phys. Soc.* (London) **89**, 479.

[23] Besag, J. & Green, P.J. (1993) Spatial statistics and Bayesian computation. *Journal of the Royal Statistical Society B* **55**, 25–102.

[24] Bialek, W., Callan, C.G., & Strong, S.P. (1996) Field Theories for Learning Probability Distributions. *cond-mat/9607180*.

[25] Binder, K. (ed.) (1986) Monte Carlo methods in statistical physics. (2nd ed.) Berlin: Springer-Verlag.

[26] Binder, K. (ed.) (1987) Applications of the Monte Carlo methods in statistical physics. (2nd ed.) Berlin: Springer-Verlag.

[27] Binder, K. (ed.) (1995) The Monte Carlo method in condensed matter physics. (2nd ed.) Berlin: Springer-Verlag.

[28] Binder, K., Heermann, D.W. (1988) Monte Carlo simulation in statistical physics: an introduction. Berlin: Springer-Verlag.

[29] Binney, J.–J., Dowrick, N.J., Fisher, A.J., & Newman, M.E.J. (1992) The Theory of Critical Phenomena. Oxford Science Publications. Oxford: Clarendon Press.

[30] Bishop, C.M. (1995) Training with noise is equivalent to Tikhonov regularization. *Neural Computation* **7** (1), 108–116.

[31] Bishop, C.M. (1995) Neural Networks for Pattern Recognition. Oxford: Oxford University Press.

[32] Blatt, M., Wiseman, S., & Domany, E. (1997) Data Clustering Using a Model Granular Magnet, Neural Computation (in print)

[33] Bleistein, N. & Handelsman, R.A. (1986) Asymptotic Expansions of Integrals. New York: Dover.

[34] Buntine, W.L. & Weigend, A.S. (1991) Bayesian back-propagation. *Complex Systems* **5**, 603–643.

[35] Connor, J.N.L. & Marcus, R.A. (1971) *J. Chem. Phys.* **55**, 5636.

[36] Davis, L. (ed.) (1987) Genetic Algorithms and Simulated Annealing. Morgan Kaufmann.

[37] Davis, L. (ed.) (1991) Handbook of Genetic Algorithms. Van Nostrand Reinhold.

[38] Dayan, P., Hinton, G.E., Neal, R.M., & Zemel, R.S. (1995) The Helmholtz Machine. *Neural Computation* **7**, 889–904.

[39] De Bruijn, N.G. (1981) Asymptotic Methods in Analysis. New York: Dover.

[40] Dempster, A.P., Laird, N.M., & Rubin, D.B. (1976) Maximum likelihood from incomplete data via the EM algorithm. *Proceedings of the Royal Statistical Society*, 1-38.

[41] Derka, R., Bužek, V., Adam, G., & Knight, P.L. (1996) From quantum Bayesian inference to quantum tomography. quant-ph/9701029.

[42] Devaney, R.L. (1986) An Introduction to Chaotic Dynamical Systems. Menlo Park: Benjamin/-Cummings.

[43] Dudley, R.M. (1984) A course on empirical processes. *Lecture Notes in Mathematics*, 1097:2-142.

[44] Durbin, R., & Willshaw, D. (1987) An analog approach to the traveling salesman problem using an elastic net method. *Nature* (London) **326**, 689–691.

[45] Efron, B. & Tibshirani R.J. (1993) An Introduction to the Bootstrap. New York: Chapman & Hall.

[46] Everitt, B.S. & Hand, D.J. (1981) Finite Mixture Distributions. London: Chapman and Hall.

[47] Feller, W. (1957) An Introduction to Probability Theory and Its Applications. Vol. I (2nd edition) New York: Wiley.

[48] Fernández, R., Fröhlich, J., & Sokal, A.D. (1992) Random Walks, Critical Phenomena, and Triviality in Quantum Field Theory. Berlin: Springer–Verlag.

[49] Ferraro, M. (1992) Invariant Pattern Representations and Lie Groups Theory. *Advances in Electronics and Electron Physics*, **84**, 131.

[50] Gardiner, C.W. (1990) Handbook of Stochastic Methods. (2nd edition) Berlin: Springer Verlag.

[51] Gelfand, A.E. & Smith, A.F.M. (1990) Sampling–based approaches to calculating marginal densities. *Journal of the American Statistical Association* **85**, 398–409.

[52] Gelfand, A.E., Hills, S.E., Racine–Poon, A., & Smith, A.F.M. (1990) Illustration of Bayesian inference in normal data models using Gibbs sampling. *Journal of the American Statistical Association* **85**, 972–985.

[53] Gelman, A., Carlin, J.B., Stern, H.S., & Rubin, D.B. (1995) Bayesian Data Analysis. New York: Chapman & Hall.

[54] Geman, S & Geman, D. (1984) Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images. *IEEE Trans. on Pattern Analysis and Machine Intelligence* **6**, 721–741. Reprinted in Shafer & Pearl (eds.) (1990) Readings in Uncertainty Reasoning. San Mateo, CA: Morgan Kaufmann.

[55] Geyer, C. (1992) Practical Markov chain Monte Carlo (with discussion). *Statistical Science* **7**, 473–511.

[56] Giraud, B.G. & Nagarajan, M.A. (1991) *Ann. Phys. (N.Y.)* **212**, 260.

[57] Girosi, F. & Chan, N.T. (1995) Prior Knowledge and the Creation of Virtual Examples for RBF Networks. *IEEE Workshop on Neural Networks for Signal Processing*, September 1995, Cambridge, MA.

[58] Girosi, F., Jones, M., & Poggio, T. (1995) Regularization Theory and Neural Networks Architectures. *Neural Computation* **7** (2), 219–269.

[59] Glimm, J. & Jaffe, A., (1987) Quntum Physics. A Functional Integral Point of View. New York: Springer–Verlag.

[60] Goldberg, D.E. (1989) Genetic Algorithms in Search, Optimization, and Machine Learning. Addison–Wesley.

[61] Golden, R.M. (1996) Mathematical Methods for Neural Network Analysis and Design. Cambridge, MA: MIT Press.

[62] Goldenfeld, N. (1992) Lectures on Phase Transitions and the Renormalization Group. Frontiers in Physics Series Vol.85, Addison–Wesley.

[63] Green, P.J. & Silverman, B.W. (1994) Nonparametric Regression and Generalized Linear Models. London: Chapman & Hall.

[64] Hackbusch, W. (1989) Integralgleichungen. Teubner Studienbücher. Stuttgart: Teubner.

[65] Härdle, W. (1990) Applied nonparametric regression. Cambridge: Cambridge University Press.

[66] Hammersley, J.M. & Handscomb, D.C. (1964) Monte Carlo Methods. London: Chapman & Hall.

[67] Hastie, T.J. & Tibshirani, R.J. (1990) Generalized Additive Models. London: Chapman & Hall.

[68] Hastings, W.K. (1970) Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* **57**, 97–109.

[69] Haussler, D. (1995) Decision Theoretic Generalizations of the PAC Model for Neural Net and Other Learning Applications. In: Wolpert, D. (Ed.) The Mathematics of Generalization. SFI Proceedings Volume XX, Addison–Wesley.

[70] Hertz, J., Krogh, A. & Palmer, R.G. (1991) Introduction to the Theory of Neural Computation. Santa Fe Institute, Lecture Notes Volume I, Addison–Wesley.

[71] Hinton, G.E., Sejnowski, T.J. (1983) Optimal Perceptual Inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (Washington 1983), 448–453. New York: IEEE.

[72] Hinton, G.E., Sejnowski, T.J. (1986) Learning and Relearning in Boltzmann machines. In Rumelhart, D.E., McClelland, J.L., and the PDP Research Group (1986) *Parallel Distributed Processing*, vol. 1, chap. 7. Cambridge, MA: MIT Press.

[73] Hinton, G.E., Dayan, P., Frey, B.J., & Neal R.M., (1995) The "Wake–Sleep" Algorithm for Unsupervised Neural Networks. *Science* **268**, 1158–1161.

[74] Hofmann, T., & Buhmann, J.M., (1997) Pairwise Data Clustering by Deterministic Annealing. IEEE Transactions on Pattern Analysis and Machine Intelligence, **19** (1), 1–14.

[75] Holland, J.H. (1975) Adaption in Natural and Artificial Systems. University of Michigan Press. (2nd ed. MIT Press, 1992.)

[76] Ivanchenko, Y.M. & Lisyansky, A.A. (1995) Physics of Critical Fluctuations. Springer–Verlag.

[77] Itzkyson, C. & Drouffe, J.–M., (1989) Statistical Field Theory. (Vols. 1 and 2) Cambridge: Cambridge University Press.

[78] Jaakola, T., & Jordan, M.I. (1996) Computing upper and lower bounds on likelihoods in intractable networks. M.I.T., Artificial Intelligence Laboratory, Memo No.1571.

[79] Jacobs, R.A., Jordan, M.I., Nowlan, S.J. & Hinton, G.E. (1991) Adaptive mixture of local experts. *Neural Computation* **3**, 79–87.

[80] Jaynes, E.T., (1996) Probability Theory: The Logic of Science. Fragmentary Edition March 1996. (available under ftp://bayes.wustl.edu/Jaynes.book/)

[81] Jordan, M.I. (1995) Why the logistic function? A tutorial discussion on probabilities and neural networks. Massachusetts Institute of Technology, *Computational Cognitive Science Technical Report 9503*

[82] Jordan, M.I. & Jacobs, R.A. (1994) Hierarchical mixtures of experts and the EM algorithm. *Neural Computation* **6**, 181–214.

[83] Jeffrey, R. (1968) Probable knowledge. In *The problem of inductive logic*, ed. I. Lakatos. Amsterdam: North Holland.

[84] Jensen, F.V. (1996) An Introduction to Bayesian Networks. New York: Springer–Verlag.

[85] van Kampen, N.G. (1992) Stochastic Processes in Physics and Chemistry. Amsterdam: North–Holland.

[86] Kahneman, D. & Tversky, A. (1979) Prospect theory: An Analysis of Decision under Risk. *Econometrica*, **47** (263), 269–291.

[87] Kahneman, D. & Tversky, A. (1982a) On the study of statistical intuitions. *Cognition*, **11**, 123–141.

[88] Kahneman, D. & Tversky, A. (1982b) Variants of uncertainty. *Cognition*, **11**, 143–157.

[89] Kahneman, D. & Tversky, A. (1982c) A reply to Evans. *Cognition*, **12**, 325–326.

[90] Kahneman, D., Slovic, P., & Tversky, A. (1983) Judgment under uncertainty: Heuristics and biases. Cambridge: Cambridge University Press.

[91] Kandel, A. (1982) Fuzzy Techniques in Pattern Recognition. New York: John Wiley and Sons.

[92] Kearns, M.J. & Vazirani, U.V. (1994) An Introduction to Computational Learning Theory. Cambridge, MA: MIT Press.

[93] Klir, J.G. & Yuan, B. (1995) Fuzzy Sets and Fuzzy Logic. Upper Saddle River, NJ: Prentice Hall.

[94] Kohonen, T. (1995) Self–Organizing Maps. Berlin: Springer.

[95] Kolmogorov, A.N. (1957) On the representation of continuous functions of several variables by superposition of continuous functions of one variable and addition. *Doklady Akademiia Nauk SSSR* **114** (5), 953–956.

[96] Kontkanen, P., Myllymäki, P., & Tirri, H. (1997) Constructing Bayesian finite mixture models by the EM algorithm. NeuroCOLT Technical Report Series, NC–TR–97–003.

[97] Landau, L.D. & Lifshitz, E.M. (1980) Statistical Physics Part 1 (3rd edition) New York: Pergamon.

[98] Laplace, P.S. (1810a) Mémoire sur les formules qui sont fonctions de très grands nombres et sur leurs application aux probabilités. *Oeuvres de Laplace* **12**, 301–345.

[99] Laplace, P.S. (1810b) Mémoire sur les intégrales définies et leur application aux probabilités. *Oeuvres de Laplace* **12**, 357–412.

[100] Lauritzen, S.L. (1996) Graphical Models. New York: Oxford University Press.

[101] Le Bellac, M. (1991) Quantum and Statistical Field Theory. Oxford Science Publications. Oxford: Clarendon Press.

[102] Le Cam, L. (1953) On some asymptotic properties of maximum likelihood estimates and related Bayes estimates. *University of California Publications in Statistics* **1** (11), 277–330.

[103] Le Cam, L. (1986) Asymptotic Methods in Statistical Decision Theory. New York: Springer.

[104] Le Cam, L. & Yang, G.L. (1990) Asymptotics in Statistics. New York: Springer.

[105] Leen, T.K. (1995) From Data Distributions to Regularization in Invariant Learning. *Neural Computation* **7**, 974–981.

[106] Lemm, J.C. (1984) Entscheidungsstrategien und Informationskosten. Diplomarbeit Universität Münster.

[107] Lemm, J.C., Giraud, B.G., & Weiguny, A. (1994) *Phys. Rev. Lett.* **73** 420.

[108] Lemm, J.C. (1995a) Inhomogeneous Random Phase Approximation for Nuclear and Atomic Reactions. *Annals of Physics* **244** (1), 136–200.

[109] Lemm, J.C. (1995b) Inhomogeneous Random Phase Approximation: A Solvable Model. *Annals of Physics* **244** (1), 201–238.

[110] Lemm, J.C., Beiu, V., & Taylor, J.G. (1995) Density Estimation as a Preprocessing Step for Constructive Algorithms. In Neural Networks: Artificial Intelligence and Industrial Applications. Proceedings of the Third Annual SNN Symposium on Neural networks, Nijmegen, The Netherlands, 14–15 September 1995 B. Kappen & S. Gielen (eds.) London: Springer Verlag.

[111] Lin, S.–H., Kung, S.–Y., & Lin, L–J (1997) Face recognition/detection by probabilistic Decision–Based Neural Network. *IEEE Trans. Neural Networks* **8** (1), 114–132.

[112] MacKay, D.J.C. (1991) Bayesian Methods for Adaptive Models. Ph.D thesis, California Institute of Technology.

[113] MacKay, D.J.C. (1992a) Bayesian interpolation. *Neural Computation* **4** (3), 415–447.

[114] MacKay, D.J.C. (1992b) A practical Bayesian framework for backpropagation networks. *Neural Computation* **4**, 448–472.

[115] MacKay, D.J.C. (1992c) The evidence framework applied to classification networks. *Neural Computation* **4**, 720–736.

[116] Metropolis, N., Rosenbluth, A.W., Rosenbluth, M.N., Teller, A.H., & Teller, E., (1953) Equation of state calculations by fast computing machines. *Journal of Chemical Physics* **21**, 1087–1092.

[117] Mezard, M., Parisi, G., & Virasoro, M.A. (1987) Spin Glass Theory and Beyond. World Scientific.

[118] Michalewicz, Z. (1992) Genetic Algorithms + Data Structures = Evolution Programs. Springer–Verlag.

[119] Miller, W.H. (1970) *J. Chem. Phys.* **53**, 3578.

[120] Minski, M.L. & Papert, S.A. (1990) (Expanded Edition, Original edition, 1969) Perceptrons. Cambridge, MA: MIT Press.

[121] Mitchell, M. (1996) An Introduction to Genetic Algorithms. Cambridge, MA: MIT Press.

[122] Montvay, I. & Münster, G., (1994) Quantum Fields on a Lattice, Cambridge University Press.

[123] Neal, R.M., & Hinton, G.E., (1993) A new view of the EM algorithm that justifies incremental and other variants. (submitted to *Biometrika*).

[124] Neal, R.M. (1993) Probabilistic inference using Markov chain Monte Carlo methods. Technical Report CRG–TR–93–1, Dept. of Comp.. Sc., Univ. of Toronto, Canada.

[125] Neal, R.M. (1996) Bayesian Learning for Neural Networks. New York: Springer.

[126] Neal, R.M. (1997) Monte Carlo Implementation of Gaussian Process Models for Bayesian Regression and Classification. Technical Report No. 9702, Dept. of Statistics, Univ. of Toronto, Canada.

[127] Negele, J., W. & Orland, H. (1988) Quantum Many–Particle Systems. Frontiers In Physics Series (Vol. 68), Addison–Wesley.

[128] Newton, R.G. (1982) Scattering Theory of Waves and Particles. New York: Springer–Verlag.

[129] Pearl, J. (1988) Probabilistic Reasoning in Intelligent Systems. San Mateo, CA: Morgan Kauffmann.

[130] Pierre, D.A. (1986) Optimization Theory with Applications. New York: Dover. (Original edition Wiley, 1969).

[131] Poggio, T. & Girosi, F. (1990) Networks for Approximation and Learning. Proceedings of the IEEE, Vol 78, No. 9.

[132] Pollard, D. (1984) Convergence of Stochastic Processes. New York: Springer Verlag.

[133] Pomerleau, D. (1991) *Neural Computation* **3**, 88.

[134] Press, W.H., Teukolsky, S.A., Vetterling, W.T., & Flannery, B.P. (1992) Numerical Recipes in C. (2nd ed.) Cambridge University Press.

[135] Ratsaby, J. & Maiorov V. (1996) Learning from Examples and Side Information. NeuroCOLT Technical Report Series, NC–TR–96–050.

[136] Richter, S.L. & DeCarlo, R.A. (1983) Continuation methods: theory and applications. *IEEE Trans. Circuits Syst.* CAS–**30**, 347–352.

[137] Ripley, B.D. (1977) Modelling spatial patterns (with discussion). *Journal of the Royal Statistical Society series B* **39**, 172–212.

[138] Ripley, B.D. (1987) Stochastic Simulation. New York: Wiley.

[139] Ripley, B.D. (1996) Pattern Recognition and Neural Networks. Cambridge University Press.

[140] Rissanen, J. (1989) Stochastic Complexity in Statistical Inquiry. Singapore: World Scientific.

[141] Roepstorff, G. (1991) Pfadintegrale in der Quantenphysik. Braunschweig: Vieweg.

[142] Rose, K., Gurewitz, E., & Fox, G.C. (1990) Statistical mechanics and phase transitions in clustering. *Phys. Rev. Lett.* **65**, 945–948.

[143] Safran, S.A., (1994) Statistical Thermodynamics of Surfaces, Interfaces, and Membranes. Frontiers in Physics Series Vol.90, Addison–Wesley.

[144] Saul, L.K., Jaakola, T., & Jordan, M.I. (1996) Mean Field Theory for Sigmoid Belief Networks. *Journal of Artificial Intelligence Research* **4**, 61–76.

[145] Schervish, M.J. (1995) Theory of Statistics. New York: Springer–Verlag.

[146] Schwefel, H.–P. (1995) Evolution and Optimum Seeking. Wiley.

[147] Seneta, E. (1981) Non–Negative Matrices and Markov Chains. (2nd ed.) New York : Springer-Verlag,

[148] Seung, H.S. (1995) Annealed Theories of Learning. In Neural Networks: The Statistical Mechanics Perspective, Proceedings of the CTP-PBSRI Joint Workshop on Theoretical Physics, World Scientific.

[149] Shawe–Taylor, J., Bartlett, P.L., Williamson, R.C., & Anthony, M. (1996a) A Framework for Structural Risk Minimization. NeuroColt Technical Report Series, NC-TR-96-032.

[150] Shawe–Taylor, J., Bartlett, P.L., Williamson, R.C., & Anthony, M. (1996b) Structural Risk Minimization over Data–Dependent Hierachies. Neuro-Colt Technical Report Series, NC-TR-96-053.

[151] Sietsma, J. & Dow, R.J.F. (1991) Creating artificial neural networks that generalize. *Neural Networks* **4** (1), 67–79.

[152] Silverman, B.W. (1986) Density Estimation for Statistics and Data Analysis. London: Chapman& Hall.

[153] Smith, A.F.M. & Roberts, G.O. (1993) Bayesian computation via Gibbs sampler and related Markov chain Monte Carlo methods (with discussion). *Journal of the Royal Statistical Society series B* **55**, 3–23.

[154] Smola, A.J. & Schölkopf, B. (1997) On a Kernel-based Method for Pattern Recognition, Regression, Approximation and Operator Inversion. Preprint.

[155] Stone, M. (1974) Cross–validation choice and assessment of statistical predictions. *Journal of the Royal Statistical Society B* **36**, 111-147.

[156] Stone, M. (1974) An asymptotic equivalence of choice of model by cross–validation and Akaike's criterion. *Journal of the Royal Statistical Society B* **39**, 44.

[157] Tanner, M.A. (1993) Tools for Statistical Inference. New York: Springer Verlag.

[158] Taylor, J.R. (1972) Scattering Theory. New York: Wiley.

[159] Tierney, L. (1994) Markov chains for exploring posterior distributions (with discussion). *Annals of Statistics* **22**, 1701–1762.

[160] Tikhonov, A.N. (1963) Solution of incorrectly formulated problems and the regularization method. *Soviet Math. Dokl.* **4**, 1035–1038.

[161] Titterington, D.M., Smith, A.F.M., & Makov, U.E. (1985) Statistical Analysis of Finite Mixture Distributions. New York: Wiley.

[162] Tversky, A. (1972) Elimination by Aspects. *Psychological Review*, **79** (4), 281–301.

[163] Tversky, A. & Kahneman, D. (1981) The Framing of Decisions and the Psychology of Choice. *Science*, **211** (4481), 453–458.

[164] Vapnik, V.N. (1982) Estimation of Dependences Based on Empirical Data. New York: Springer–Verlag.

[165] Vapnik, V.N. (1995) The Nature of Statistical Learning Theory. New York: Springer–Verlag.

[166] Verri, A. & Poggio, T. (1986) Regularization Theory and Shape Constraints. M.I.T., Artificial Intelligence Laboratory, Memo No.916.

[167] Vetter, T., Poggio, T., & Bülthoff, H. (1992) 3D Object Recognition: Symmetry and Virtual Views. M.I.T., Artificial Intelligence Laboratory, Memo No.1409.

[168] Wahba, G. (1983) Bayesian confidence intervals for the cross–validated smoothing spline. *J. Roy. Statist. Soc.* **B, 45**, 133–150.

[169] Spline Models for Observational Data. Philadelphia: SIAM.

[170] Wald, A. (1938) Contributions to the theory of statistical estimation and testing hypotheses. *Ann. Math. Stat.* **10**, 299–326.

[171] Watkin, T.L.H., Rau, A., & Biehl, M. (1993) The statistical mechanics of learning a rule. *Rev. Mod. Phys.*, **65**, 499–556.

[172] Webb, A.R. (1994) Functional approximation by feed–forward networks: a least–squares approach to generalization. *IEEE Transactions on Neural Networks* **5**, 363–371.

[173] Whitesitt, J.E. (1995) Boolean Algebra and Its Applications. New York: Dover, (Originally published: Addison–Wesley, 1961).

[174] Williams, C.K.I. & Rasmussen, C.E. (1996) Gaussian processes for regression. In D.S. Touretzki, M.C. Mozer, & M.E. Hasselmo (eds.) *Advances in Neural Information Processing Systems 8* Cambridge, MA: MIT Press.

[175] Wierling, A., Giraud, B., Mekideche, F., Horiuchi, H., Maruyama, T., Ohnishi, A., Lemm, J.C., & Weiguny A. (1994) *Z. Phys. A* **348**, 153.

[176] Wolpert, D.H. (1994) Bayesian Backpropagation over I–O functions rather than weights. In Cowan, J.D., Tesauro, G., & Alspector, J. (eds.) *Advances in Neural Information Processing Systems 6* San Francisco, CA: Morgan Kaufmann, pp. 200–207.

[177] Wolpert, D.H. (1994b) Discussion of Ripley, B.D., Neural networks and related methods for classification. *Journal of the Royal Statistical Society B* **56**, 450–451.

[178] Wolpert, D.H. (1996a) The Lack of A Priori Distinctions between Learning Algorithms. *Neural Computation* **8** (7), 1341-1390.

[179] Wolpert, D.H. (1996b) The Existence of A Priori Distinctions between Learning Algorithms. *Neural Computation* **8** (7), 1391-1420.

[180] Yuille, A.L., (1990) Generalized deformable models, statistical physics and matching problems. *Neural Computation*, **2**, (1) 1–24.

[181] Yuille, A.L. & Kosowsky, J.J., (1994) Statistical Physics Algorithms That Converge. *Neural Computation*, **6**, (3) 341–356.

[182] Yuille, A.L., Stolorz, P., $ Utans, J. (1994) Statistical Physics, Mixtures of Distributions, and EM Algorithm. *Neural Computation*, **6** (2), 334–340.

[183] Zadeh, L.A. (1987) Fuzzy Sets and Applications. Selected Papers. R.R. Yager et al. (eds), New York: Wiley.

[184] Zadeh, L.A. (1996) Fuzzy Sets, Fuzzy Logic, and Fuzzy Systems. Selected Papers. G.J. Klir & B. Yuan (eds.) Singapore: World Scientific.

[185] Zinn–Justin, J. (1989) Quantum Field Theory and Critical Phenomena. Oxford Science Publications. Oxford: Clarendon Press.