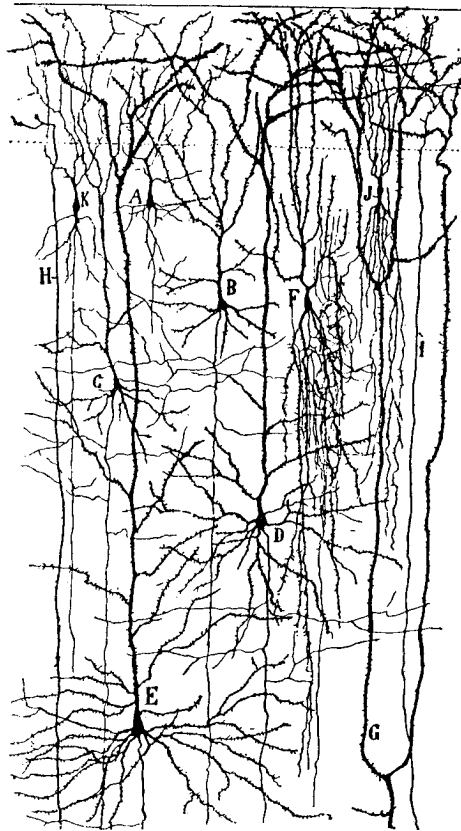


Advanced Neural Networks

Lecture Notes of Course G30/CMNN15

September 2002



ACC Coolen
Department of Mathematics
King's College London

Contents

1	Introduction	3
2	Competitive Non-Supervised Learning Processes	5
2.1	Vector Quantization	5
2.2	Soft Vector Quantization	16
2.3	Time-Dependent Learning Rates	22
2.4	Self-Organizing Maps	24
3	Bayesian Techniques in Supervised Learning	33
3.1	Preliminaries and Introduction	33
3.2	Bayesian Learning of Network Weights	40
3.3	Predictions with Error Bars: Real-Valued Functions	46
3.4	Predictions with Error Bars: Binary Classification	52
3.5	Bayesian Model Selection	54
3.6	Practicalities: Measuring Curvature	59
4	Gaussian Processes	61
4.1	Examples of Networks Reducing to Gaussian Processes	62
4.2	Learning and Prediction with Gaussian Processes	65
4.3	The Choice of Covariance Matrix	68
5	Support Vector Machines for Binary Classification	71
5.1	Optimal Separating Plane for Linearly Separable Tasks	71
5.2	Representation in terms of Support Vectors	74
5.3	Pre-Processing, SVM Kernels	77
A	Probability in a Nutshell	83
B	Gaussian Integrals	87
C	The δ-Distribution	89

Chapter 1

Introduction

This lecture course is the sequel to the core course ‘Neural Networks’, which can be regarded as a prerequisite, in which we expand both in breadth and in depth the material covered so far. Apart from the first of the four chapters in the present notes (and in contrast to the bulk of ‘Neural Networks’), this course ‘Advanced Neural Networks’ deals with relatively novel research results which go back less than ten years.

Expansion in breadth refers to the addition of new neural network classes and new learning rules (Competitive Non-Supervised Learning Processes, this material is not too mathematical because there is little theory on these systems; Support Vector Machines, again a subject with only a modest amount of theory).

Expansion in depth refers to a more solid statistical understanding, interpretation, and increased potential for analysis and prediction of the most popular system types (Bayesian Techniques in Supervised Learning, and its spin-off Gaussian Processes). These latter subjects, although more mathematical in nature, have generated the main progress in industrial and commercial neural network applications over the last ten years, since they have removed in a rigorous way the problems related to data noise and to the quantification of the reliability of neural network decisions.

Chapter 2

Competitive Non-Supervised Learning Processes

In this first chapter we will introduce and study a couple of so-called ‘unsupervised’ learning processes, where the problem is not to learn to associate some ‘correct’ response to each possible incoming signal (specified by a ‘teacher’ or ‘supervisor’ signal, as in e.g. feed-forward networks of the Perceptron or Multi-Layer Perceptron type), but rather to build some alternative (more efficient, more compact or more structured) representation of the data vectors fed into the system. The first two procedures to do this (Vector Quantization and Soft Vector Quantization) aim to achieve data reduction for initially unknown data distributions (with typical applications in communication, where compact representations are obviously cheaper to transmit than non-compact ones, since it will allow more signals to be communicated via the same channel). The third class of systems (Self-Organizing Maps, or Feature Maps, or Kohonen Maps) aim to create a topologically correct but low-dimensional internal representation of a given data distribution. This has applications in e.g. biology (the brain is known to create such maps for sensory signals), data-base mining, medical diagnostics, etc.

2.1 Vector Quantization

Data-Reduction via Code-Book Vectors. Imagine we have a source of real-valued data vectors $\mathbf{x} = (x_1, \dots, x_n) \in \mathfrak{R}^n$, with statistics described by some probability density $p(\mathbf{x})$ (with $\int d\mathbf{x} p(\mathbf{x}) = 1$, and with the definition $\langle f(\mathbf{x}) \rangle = \int d\mathbf{x} p(\mathbf{x}) f(\mathbf{x})$). Alternatively, if the data can only take values from a discrete set, we would have $p(\mathbf{x})$ representing probabilities, with $\sum_{\mathbf{x}} p(\mathbf{x}) = 1$ and $\langle f(\mathbf{x}) \rangle = \sum_{\mathbf{x}} p(\mathbf{x}) f(\mathbf{x})$ (see appendix on elementary probability theory). Especially if n is very large, we would like to represent the real data \mathbf{x} by alternative (and simpler) signals from which the \mathbf{x} can be re-constructed, with some loss of accuracy, but with a reduction in dimensionality. One way to achieve this is the following:

- We ‘cover’ the data space by a suitable (small) set of characteristic data points, the so-called ‘code-book vectors’ $\mathbf{m}_i \in \mathfrak{R}^n$. Here i labels the individual code-book vectors, i.e. if we have N code-book vectors then $i = 1, \dots, N$.
- We then approximate (or ‘round off’) each data point \mathbf{x} by the *nearest* code-book vector, i.e. by that particular \mathbf{m}_i for which $|\mathbf{x} - \mathbf{m}_i|$ is minimal.

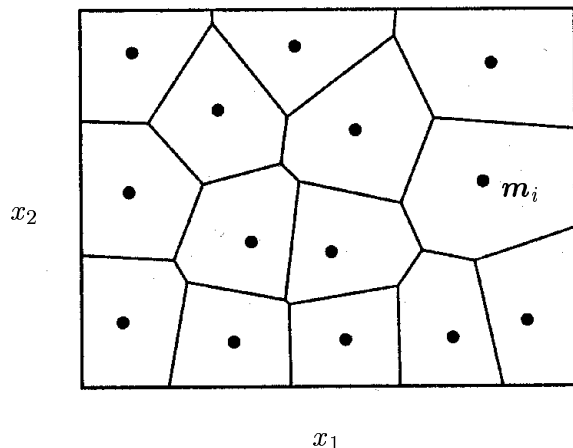


Figure 2.1: Example of a Voronoi tessellation of data space for $n = 2$, i.e. $\mathbf{x} = (x_1, x_2)$. The points \bullet represent the code-book vectors. The ‘Voronoi cell’ associated with code-book vector \mathbf{m}_i , $V_i[\{\mathbf{m}\}]$, is the compartment surrounding point \mathbf{m}_i .

It will be clear why in communication this might be desirable: instead of the n real numbers (x_1, \dots, x_n) we would now need to send only a single integer number: the index i . The price we pay for this is a reduction in accuracy; after all we have approximated \mathbf{x} by \mathbf{m}_i , and have thereby lost information. We can no longer be sure of the exact value of \mathbf{x} , since many more data points \mathbf{x} would have been replaced by \mathbf{m}_i (all those which are close to \mathbf{m}_i). The amount of information lost will obviously be less if we have a larger number or a more efficiently distributed set of code-book vectors.

The above procedure gives us what is known as a Voronoi tessellation of data space. This is a partitioning into convex subsets $V_i[\{\mathbf{m}\}]$, controlled by the choice $\{\mathbf{m}\}$ made for the N code-book vectors, defined as

$$V_i[\{\mathbf{m}\}] = \{\mathbf{x} \in \mathbb{R}^n \mid \forall j \neq i : |\mathbf{x} - \mathbf{m}_i| < |\mathbf{x} - \mathbf{m}_j|\} \quad (2.1)$$

(see figure 2.1).

A good set of code-book vectors is one with the property that the density of code-book vectors in a given region of \mathbb{R}^n is proportional to the density of data-points in that region. In other words, given the statistics of the data and given the number N of code-book vectors we are willing to invest, we aim for a set $\{\mathbf{m}\}$ such that

$$\text{Prob}[\mathbf{x} \in V_i[\{\mathbf{m}\}]] = \text{Prob}[\mathbf{x} \in V_j[\{\mathbf{m}\}]] \quad (2.2)$$

In this way all code-book vectors are used equally often, and no resources are wasted. The problem addressed by the VQ and SVQ algorithms is how to find a proper set of code-book vectors, via adaptive processes; i.e. both are based on gradually ‘learning’ a proper positioning of the code-book vectors by observing the data.

The VQ Algorithm and its Properties. The Vector Quantization (VQ) algorithm is defined as follows. First we initialize the N code-book vectors $\mathbf{m}_i \in \mathbb{R}^n$ (randomly, or according to a recipe to be given below). Then we iterate the following stochastic process:

step 1: pick a data point \mathbf{x} at random, according to the probability density $p(\mathbf{x})$

step 2: find the Voronoi cell containing \mathbf{x} , i.e. find i such that $|\mathbf{x} - \mathbf{m}_i| < |\mathbf{x} - \mathbf{m}_j| \forall j \neq i$

step 3: move \mathbf{m}_i towards \mathbf{x} : $\mathbf{m}_i \rightarrow \mathbf{m}_i + \eta(\mathbf{x} - \mathbf{m}_i)$

step 4: return to 1

The parameter $\eta > 0$ (the learning rate) controls the magnitude of the changes, one takes $\eta \ll 1$ to suppress fluctuations. The uncertainty (stochasticity) of the process is only in the realization of the subsequent data points we draw. Note that there would be a problem with data points \mathbf{x} which are exactly on the border of two Voronoi cells. In practice this is not an issue: firstly, the probability for this to happen is generally very small (unless we have a pathological distribution $p(\mathbf{x})$), and secondly, one could simply decide in those instances not to make a change.

The VQ algorithm can be seen to have the following properties:

- (i) A code-book vector will only become mobile when we pick a data point in its Voronoi cell, i.e. sufficiently close to it. Hence the process of moving the \mathbf{m}_i will proceed slowly in areas where the data density $p(\mathbf{x})$ is low.
- (ii) Unless we reduce η during the process, the code-book vectors will continue to move stochastically, although the *density* of code-book vectors in any given region should become stationary.
- (iii) VQ is very simple and (as we will see below) effective, and has just three parameters to be chosen: N , η , and the duration of the process.

As a consequence of (i) we also conclude that it is not necessarily optimal to initialize the N code-book vectors randomly: those which are initialized in regions where $p(\mathbf{x})$ is zero will *never* be used. Alternatively one could initialize the \mathbf{m}_i by putting them at the locations of the first N observed data points (by construction they can now never be initialized in regions where there are no data).

Examples of VQ in Action. The figures below illustrate the functioning of the VQ algorithm for a number of simple examples with $n = 2$ (i.e. where one has data points $\mathbf{x} = (x_1, x_2)$ in a plane) and $N = 100$ (i.e. a population of 100 code-book vectors), but for different choices of the data distribution $p(\mathbf{x})$ and with different initialization strategies. Firstly, figures 2.2 and 2.3 show examples of simple data distributions and random initialization, where the process is still seen to work fine, simply because in this case no code-book vector happens to have been initialized in data-free regions. Figures 2.4, 2.5 and 2.6 refer to strongly inhomogeneous data distributions, but now with non-random initialization (so that the inhomogeneities cannot disrupt the functioning of VQ). The process is again seen to work fine. In contrast, in figures 2.7, 2.8 and 2.9 the same three inhomogeneous distributions were considered, but now with random (i.e. inappropriate) initialization of the code-book vectors. The resulting locations of the code-book vectors illustrates quite clearly how for random initialization and inhomogeneous data distributions the VQ process fails to use its resources effectively.

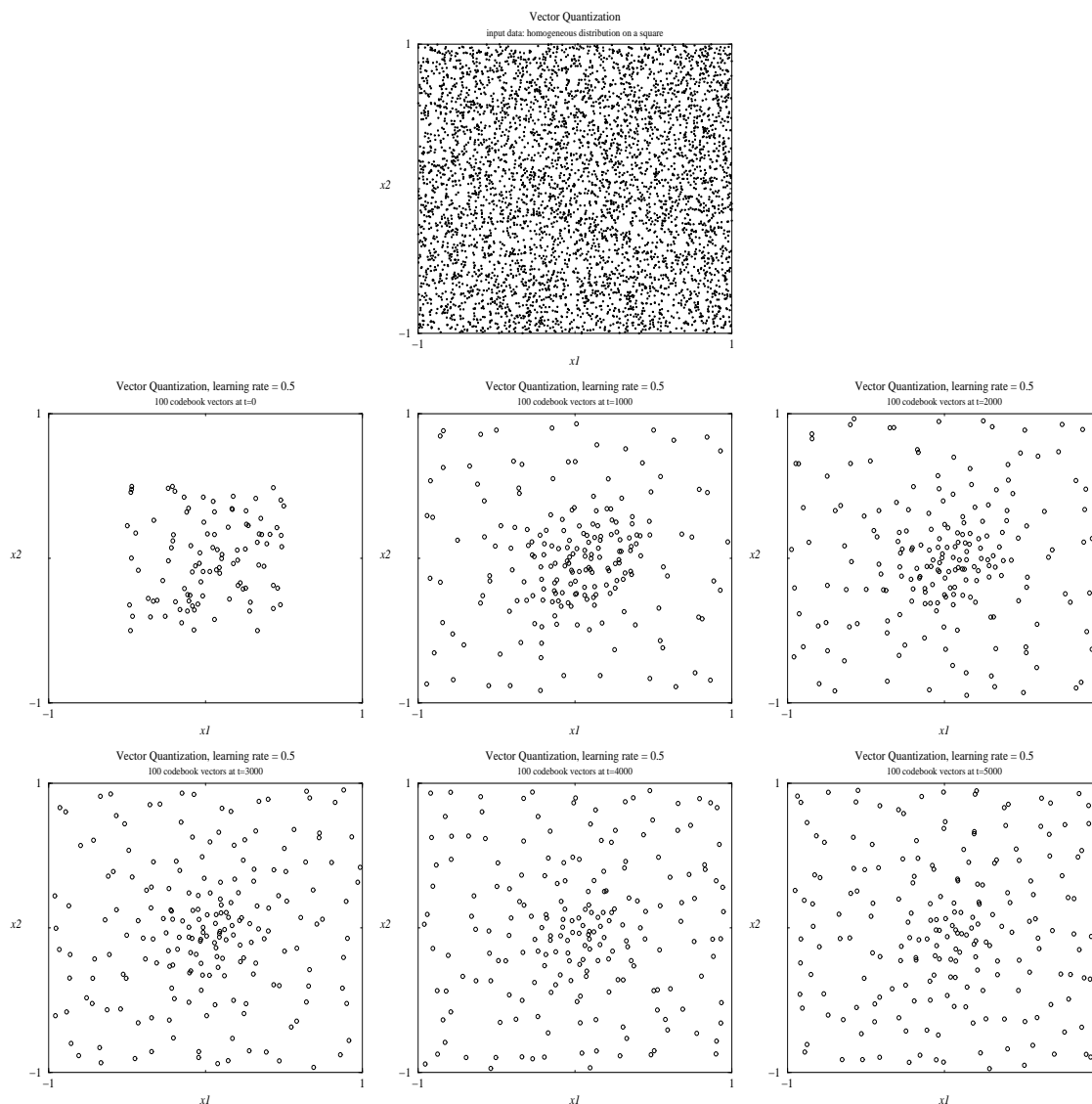


Figure 2.2: Example: numerical simulation of Vector Quantization, with $n = 2$, $\eta = 0.5$ and $N = 100$. Top graph: data distribution $p(\mathbf{x})$, homogeneously distributed over the square $[-1, 1] \times [-1, 1]$. Bottom six graphs: locations of the code-book vectors \mathbf{m}_i during the course of the process ($i = 1, \dots, 100$), at times $t = 0$ (middle row, left), $t = 1000$ (middle row, center), $t = 2000$ (middle row, right), $t = 3000$ (bottom row, left), $t = 4000$ (bottom row, middle) and $t = 5000$ (bottom row, right). Times are measures in the number of iterations. Initialisation of code-book vector allocations: randomly in the square $[-\frac{1}{2}, \frac{1}{2}] \times [-\frac{1}{2}, \frac{1}{2}]$.

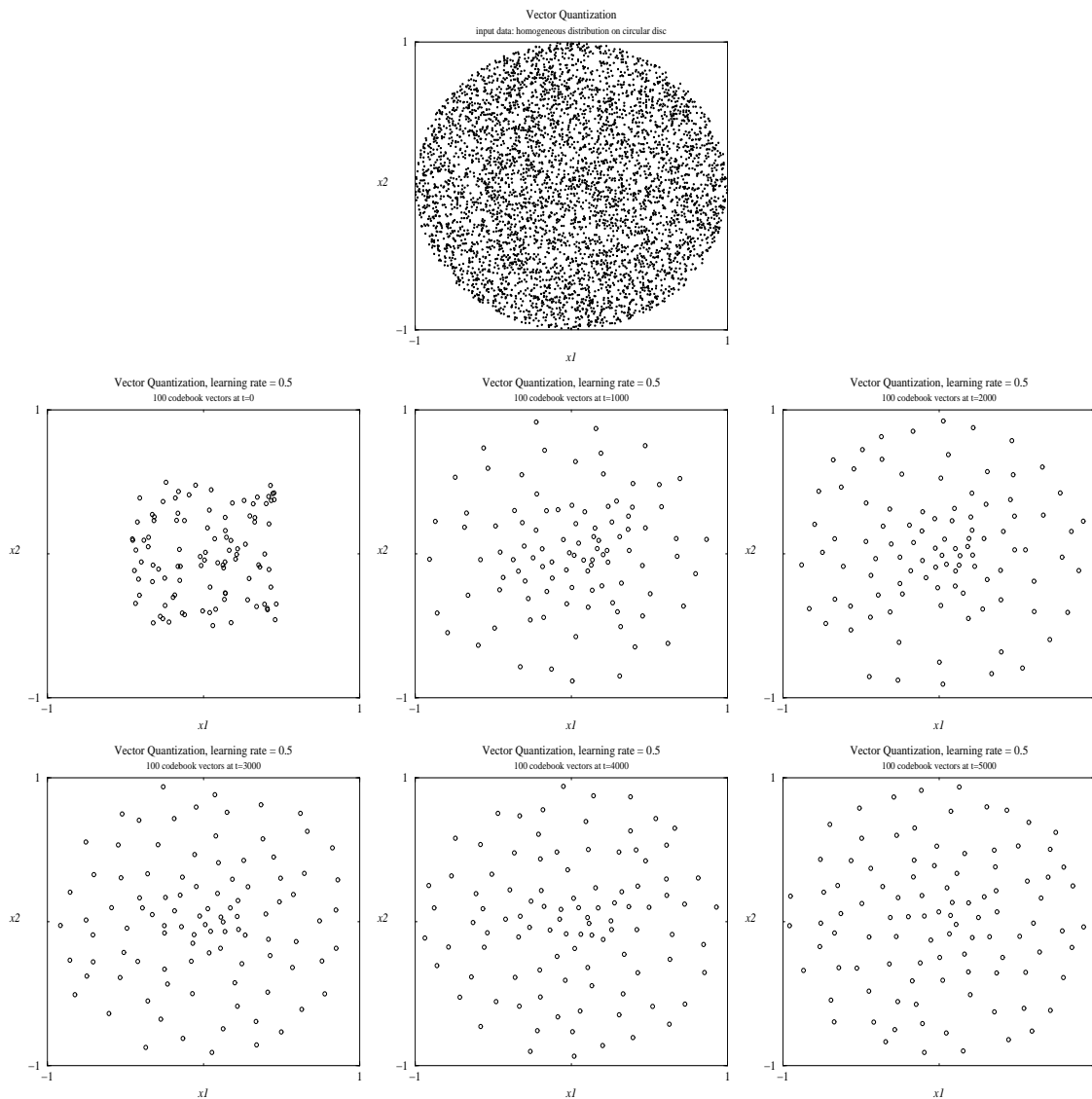


Figure 2.3: Example: numerical simulation of Vector Quantization, with $n = 2$, $\eta = 0.5$ and $N = 100$. Top graph: data distribution $p(\mathbf{x})$, homogeneously distributed over the disk $|\mathbf{x}| < 1$. Bottom six graphs: locations of the code-book vectors \mathbf{m}_i during the course of the process ($i = 1, \dots, 100$), at times $t = 0$ (middle row, left), $t = 1000$ (middle row, center), $t = 2000$ (middle row, right), $t = 3000$ (bottom row, left), $t = 4000$ (bottom row, middle) and $t = 5000$ (bottom row, right). Times are measures in the number of iterations. Initialisation of code-book vector allocations: randomly in the square $[-\frac{1}{2}, \frac{1}{2}] \times [-\frac{1}{2}, \frac{1}{2}]$.

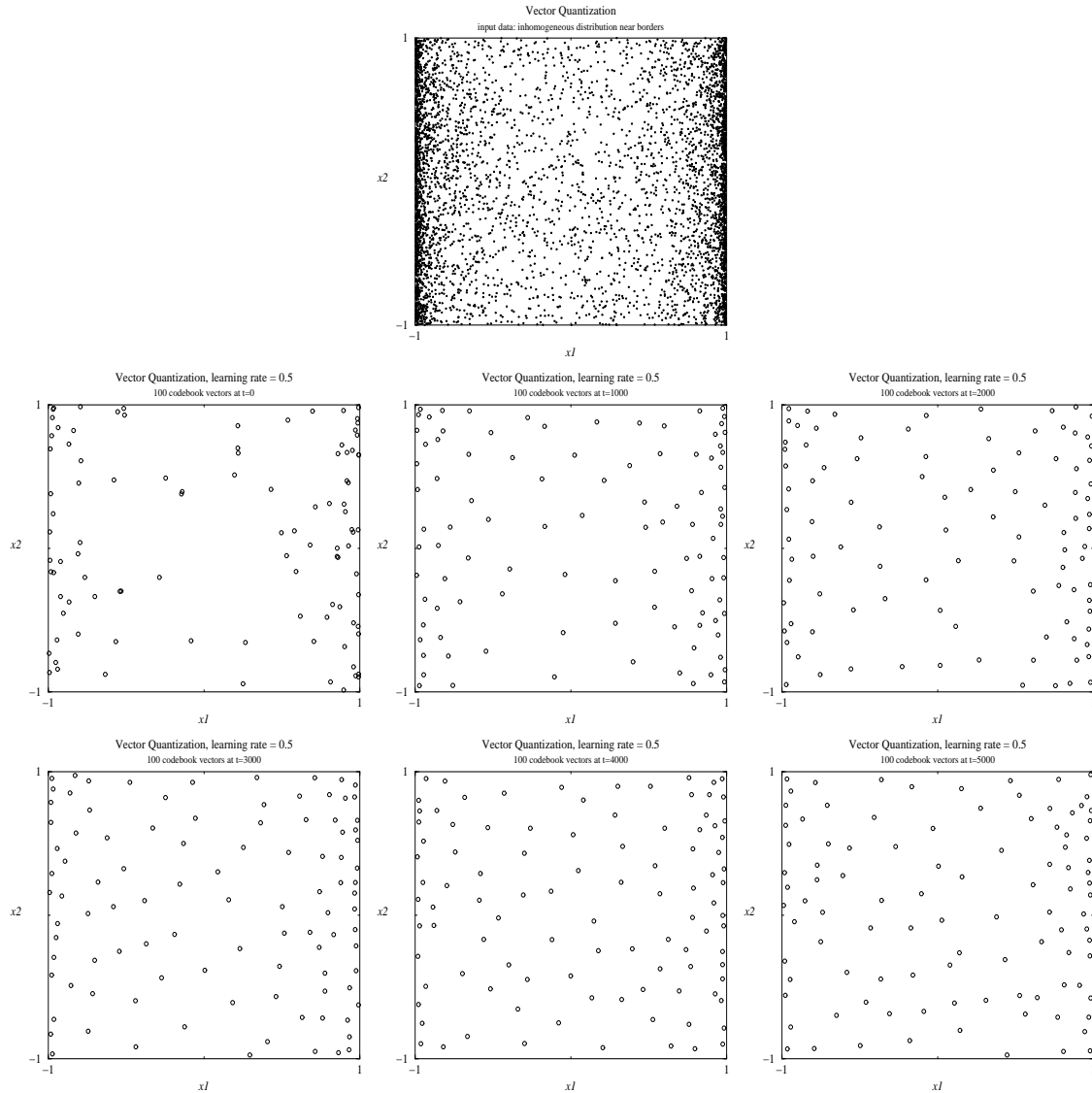


Figure 2.4: Example: numerical simulation of Vector Quantization, with $n = 2$, $\eta = 0.5$ and $N = 100$. Top graph: data distribution $p(\mathbf{x})$, inhomogeneously distributed over the square $[-1, 1] \times [-1, 1]$ (with highest data density near the borders $|x_1| = 1$). Bottom six graphs: locations of the code-book vectors m_i during the course of the process ($i = 1, \dots, 100$), at times $t = 0$ (middle row, left), $t = 1000$ (middle row, center), $t = 2000$ (middle row, right), $t = 3000$ (bottom row, left), $t = 4000$ (bottom row, middle) and $t = 5000$ (bottom row, right). Times are measures in the number of iterations. Initialisation of code-book vector allocations: positioning at the locations of the first 100 data points picked by the process.

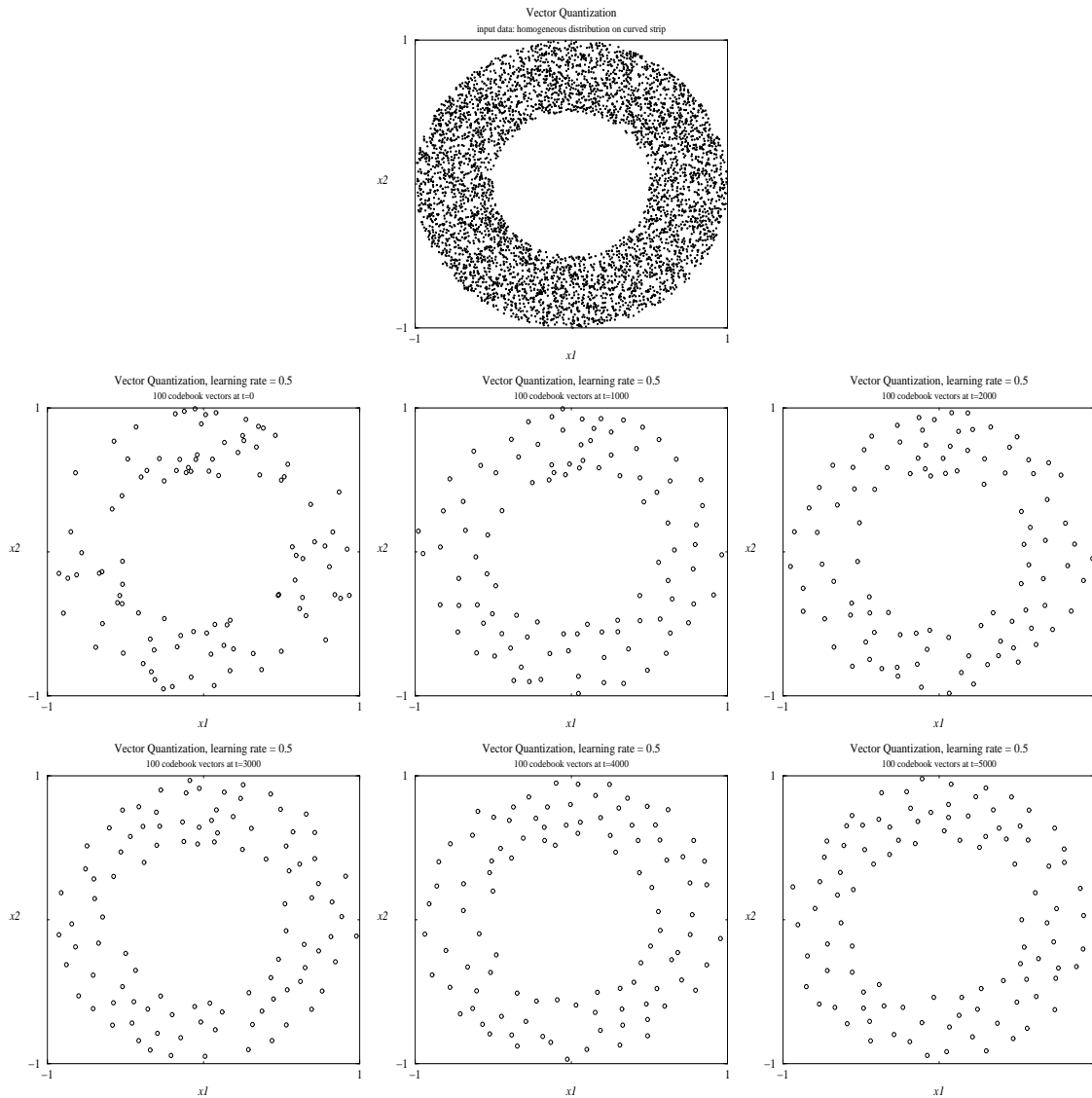


Figure 2.5: Example: numerical simulation of Vector Quantization, with $n = 2$, $\eta = 0.5$ and $N = 100$. Top graph: data distribution $p(\mathbf{x})$, homogeneously distributed over the region $\frac{1}{2} < |\mathbf{x}| < 1$. Bottom six graphs: locations of the code-book vectors \mathbf{m}_i during the course of the process ($i = 1, \dots, 100$), at times $t = 0$ (middle row, left), $t = 1000$ (middle row, center), $t = 2000$ (middle row, right), $t = 3000$ (bottom row, left), $t = 4000$ (bottom row, middle) and $t = 5000$ (bottom row, right). Times are measures in the number of iterations. Initialisation of code-book vector allocations: positioning at the locations of the first 100 data points picked by the process.

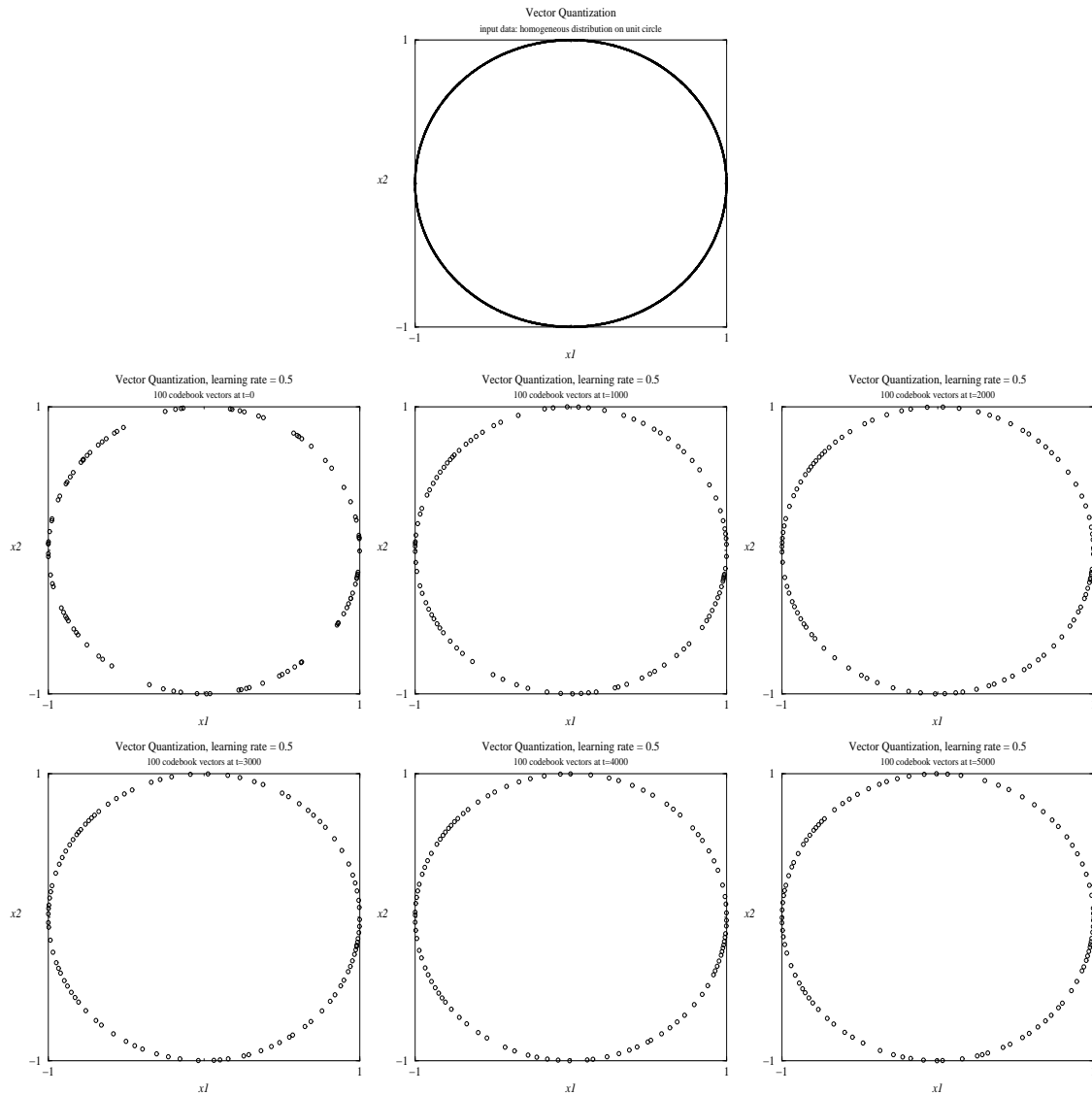


Figure 2.6: Example: numerical simulation of Vector Quantization, with $n = 2$, $\eta = 0.5$ and $N = 100$. Top graph: data distribution $p(\mathbf{x})$, homogeneously distributed over the circle $|\mathbf{x}| = 1$. Bottom six graphs: locations of the code-book vectors \mathbf{m}_i during the course of the process ($i = 1, \dots, 100$), at times $t = 0$ (middle row, left), $t = 1000$ (middle row, center), $t = 2000$ (middle row, right), $t = 3000$ (bottom row, left), $t = 4000$ (bottom row, middle) and $t = 5000$ (bottom row, right). Times are measures in the number of iterations. Initialisation of code-book vector allocations: positioning at the locations of the first 100 data points picked by the process.

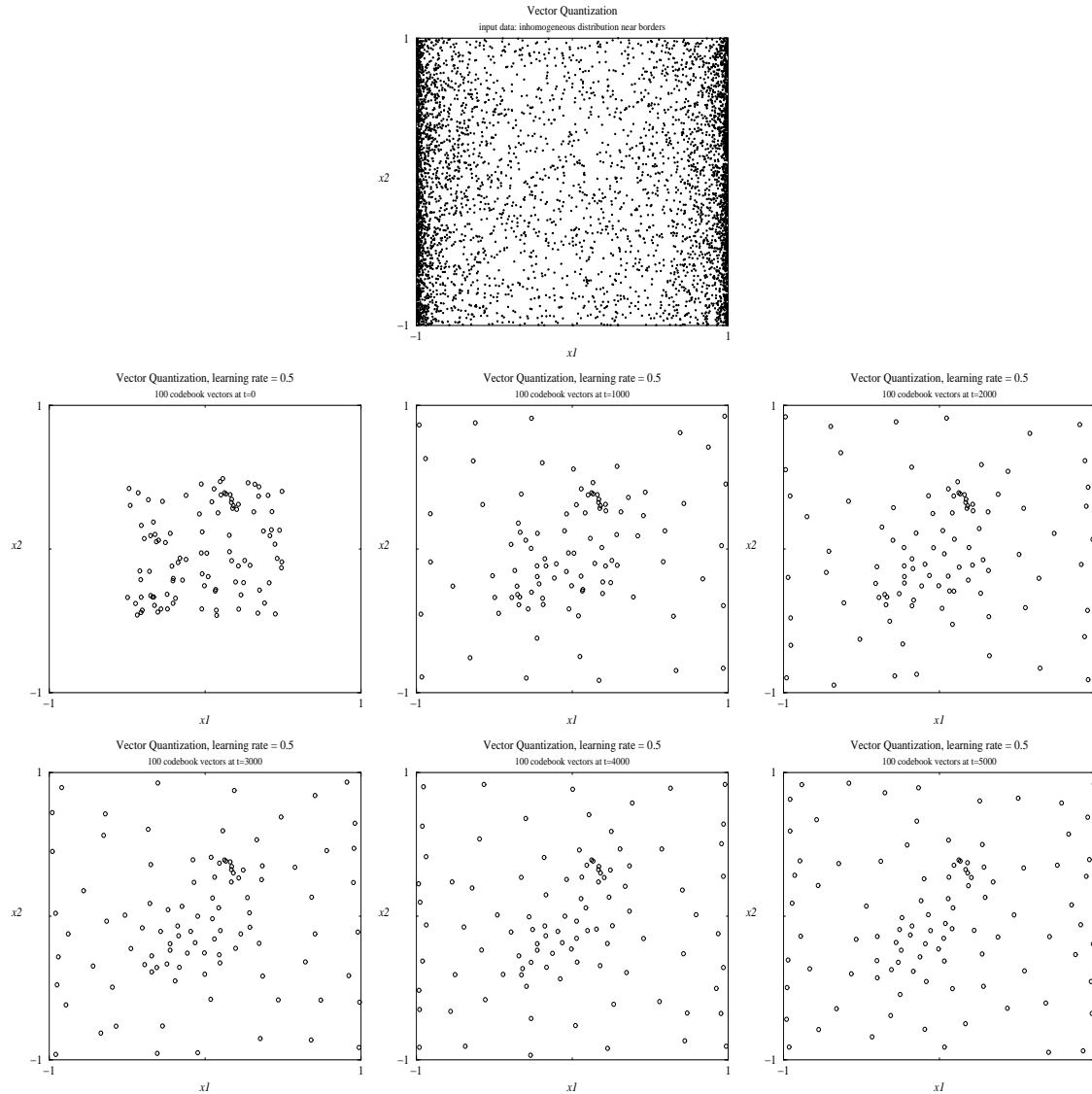


Figure 2.7: Example: numerical simulation of Vector Quantization, with $n = 2$, $\eta = 0.5$ and $N = 100$. Top graph: data distribution $p(\mathbf{x})$, inhomogeneously distributed over the square $[-1, 1] \times [-1, 1]$ (with highest data density near the borders $|x_1| = 1$). Bottom six graphs: locations of the code-book vectors \mathbf{m}_i during the course of the process ($i = 1, \dots, 100$), at times $t = 0$ (middle row, left), $t = 1000$ (middle row, center), $t = 2000$ (middle row, right), $t = 3000$ (bottom row, left), $t = 4000$ (bottom row, middle) and $t = 5000$ (bottom row, right). Times are measures in the number of iterations. Initialisation of code-book vector allocations: randomly in the square $[-\frac{1}{2}, \frac{1}{2}] \times [-\frac{1}{2}, \frac{1}{2}]$. Comparison with figure 2.4 shows that code-book vectors initialized in data regions where $p(\mathbf{x})$ is small tend to get ‘stuck’.

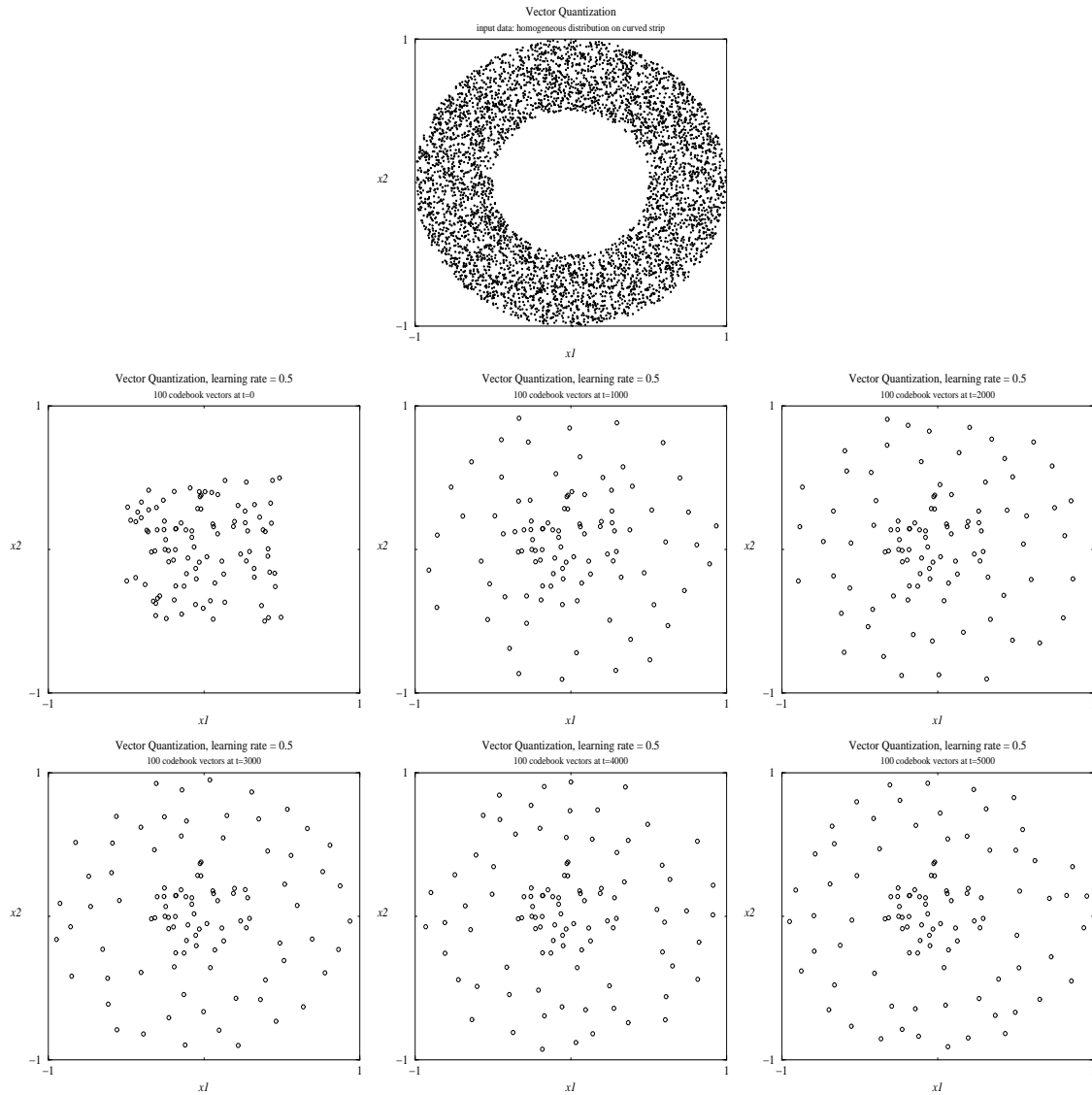


Figure 2.8: Example: numerical simulation of Vector Quantization, with $n = 2$, $\eta = 0.5$ and $N = 100$. Top graph: data distribution $p(\mathbf{x})$, homogeneously distributed over the region $\frac{1}{2} < |\mathbf{x}| < 1$. Bottom six graphs: locations of the code-book vectors \mathbf{m}_i during the course of the process ($i = 1, \dots, 100$), at times $t = 0$ (middle row, left), $t = 1000$ (middle row, center), $t = 2000$ (middle row, right), $t = 3000$ (bottom row, left), $t = 4000$ (bottom row, middle) and $t = 5000$ (bottom row, right). Times are measures in the number of iterations. Initialisation of code-book vector allocations: randomly in the square $[-\frac{1}{2}, \frac{1}{2}] \times [-\frac{1}{2}, \frac{1}{2}]$. Comparison with figure 2.5 shows that code-book vectors initialized in data regions where $p(\mathbf{x})$ is zero tend to get ‘stuck’.

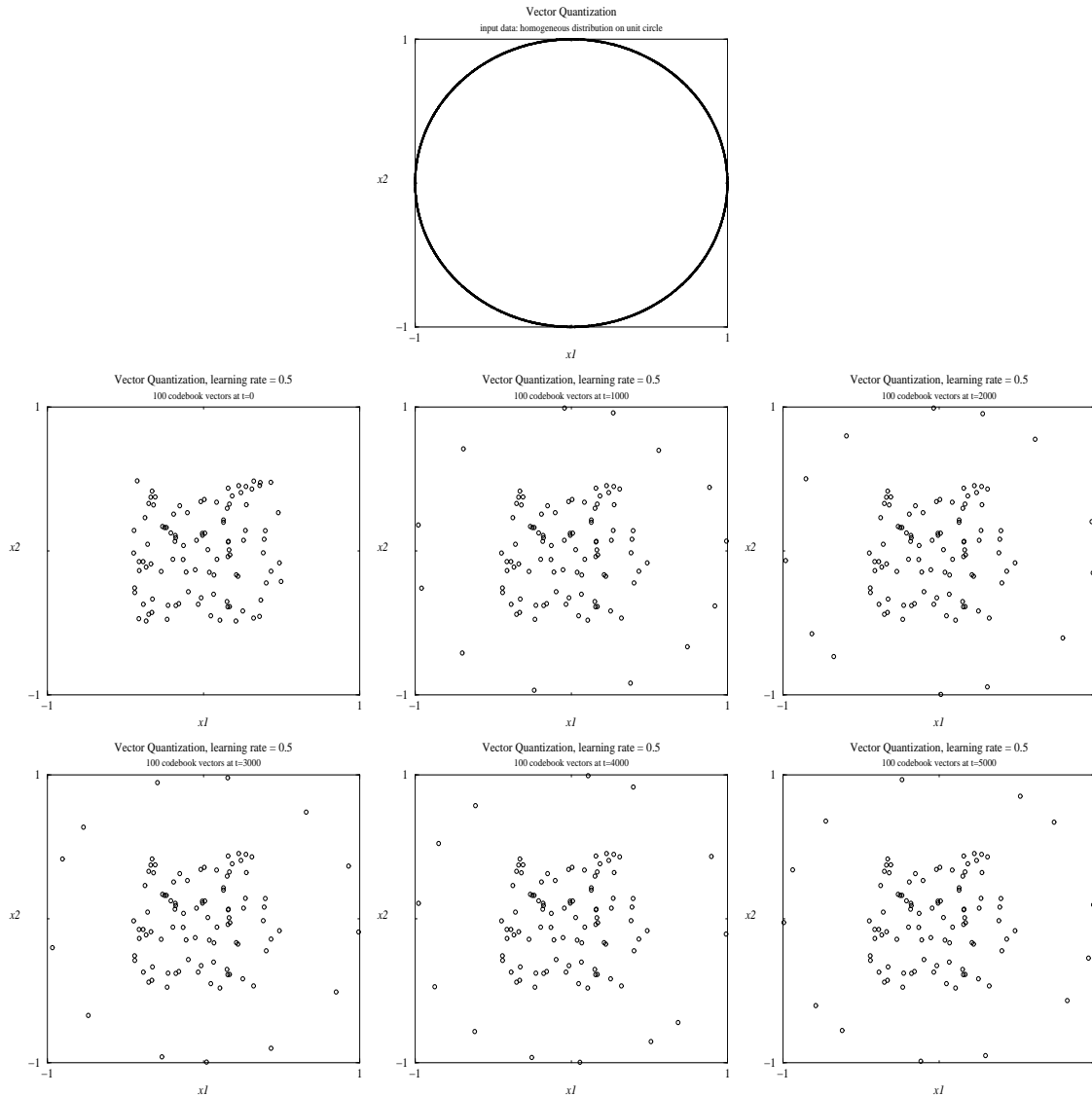


Figure 2.9: Example: numerical simulation of Vector Quantization, with $n = 2$, $\eta = 0.5$ and $N = 100$. Top graph: data distribution $p(\mathbf{x})$, homogeneously distributed over the circle $|\mathbf{x}| = 1$. Bottom six graphs: locations of the code-book vectors \mathbf{m}_i during the course of the process ($i = 1, \dots, 100$), at times $t = 0$ (middle row, left), $t = 1000$ (middle row, center), $t = 2000$ (middle row, right), $t = 3000$ (bottom row, left), $t = 4000$ (bottom row, middle) and $t = 5000$ (bottom row, right). Times are measures in the number of iterations. Initialisation of code-book vector allocations: randomly in the square $[-\frac{1}{2}, \frac{1}{2}] \times [-\frac{1}{2}, \frac{1}{2}]$. Comparison with figure 2.6 shows that code-book vectors initialized in data regions where $p(\mathbf{x})$ is zero tend to get ‘stuck’.

2.2 Soft Vector Quantization

The SVQ Algorithm and its Properties. The Soft Vector Quantization (SVQ) algorithm, which can be regarded as a ‘smooth’ version of VQ, is defined as follows. First we initialize the N code-book vectors $\mathbf{m}_i \in \mathbb{R}^n$ randomly (in contrast to VQ, random initialization poses no problem for SVQ as we will see). Then we iterate the following stochastic process:

step 1: pick a data point \mathbf{x} at random, according to the probability density $p(\mathbf{x})$

step 2: calculate, for all i :

$$F_i[\mathbf{x}, \{\mathbf{m}\}] = \frac{e^{-\beta(\mathbf{x}-\mathbf{m}_i)^2}}{\sum_{j=1}^N e^{-\beta(\mathbf{x}-\mathbf{m}_j)^2}} \quad (2.3)$$

step 3: move all \mathbf{m}_i towards \mathbf{x} : $\mathbf{m}_i \rightarrow \mathbf{m}_i + \eta(\mathbf{x} - \mathbf{m}_i)F_i[\mathbf{x}, \{\mathbf{m}\}]$

step 4: return to 1

The parameter $0 < \eta \ll 1$ (learning rate) again controls the overall magnitude of the changes. Note that, by construction, $\sum_i F_i[\mathbf{x}, \{\mathbf{m}\}] = 1$ and that always $0 \leq F_i[\mathbf{x}, \{\mathbf{m}\}] \leq 1$.

The SVQ algorithm can be seen to have the following general properties:

- (i) All code-book vectors are moved towards \mathbf{x} at every iteration step, but those which are closest to \mathbf{x} are moved most.
- (ii) Unless we reduce η during the process, the code-book vectors will continue to move stochastically, although the *density* of code-book vectors in any given region should become stationary.
- (iii) SVQ is still simple, but it has more parameter more than VQ. This extra parameter, β , defines a characteristic distance in data space: code-book vectors with $|\mathbf{x} - \mathbf{m}_i| > \beta^{-\frac{1}{2}}$ will move only weakly, in contrast to those with $|\mathbf{x} - \mathbf{m}_i| < \beta^{-\frac{1}{2}}$. Hence $1/\sqrt{\beta}$ defines the distance over which code-book vectors tend to exert an influence. Since the characteristic distance between the code-book vectors will also strongly depend on N , one should also expect the optimal value for β to depend on N .
- (iv) SVQ has an advantage over VQ when data distributions $p(\mathbf{x})$ can (slowly) change over time. Although we can ensure in VQ, by appropriate initialization, that the code-book vectors cannot get stuck in data-poor regions, this will no longer be guaranteed if $p(\mathbf{x})$ can change over time: regions which are data-rich now, and which attract code-book vectors, might become data-poor later. In SVQ the system will always be able to adapt to the new data environment, since all code-book vectors are updated all the time.

Let us next investigate the action of the SVQ algorithm for the two extreme choices to be made for the new parameter β : $\beta = \infty$ and $\beta = 0$.

Fact 1: $\lim_{\beta \rightarrow \infty} \text{SVQ} = \text{VQ}$

Proof: Define the Voronoi tessellation of data space induced by the code-book vectors $\{\mathbf{m}\}$. Consider an SVQ iteration step, where we pick data point \mathbf{x} . Assume $\mathbf{x} \in V_k[\{\mathbf{m}\}]$, i.e. \mathbf{x} is found to be in the Voronoi cell of code-book vector k (we again rule out the

pathological cases where \mathbf{x} is at the boundary of Voronoi cells; see the section on VQ). Multiply numerator and denominator of (2.3) by $e^{\beta(\mathbf{x}-\mathbf{m}_k)^2}$, and use the property that $|\mathbf{x}-\mathbf{m}_i| > |\mathbf{x}-\mathbf{m}_k|$ for all $i \neq k$:

$$\begin{aligned} \lim_{\beta \rightarrow \infty} F_i[\mathbf{x}, \{\mathbf{m}\}] &= \lim_{\beta \rightarrow \infty} \frac{e^{-\beta[(\mathbf{x}-\mathbf{m}_i)^2 - (\mathbf{x}-\mathbf{m}_k)^2]}}{\sum_{j=1}^N e^{-\beta[(\mathbf{x}-\mathbf{m}_j)^2 - (\mathbf{x}-\mathbf{m}_k)^2]}} \\ &= \lim_{\beta \rightarrow \infty} \frac{e^{-\beta[(\mathbf{x}-\mathbf{m}_i)^2 - (\mathbf{x}-\mathbf{m}_k)^2]}}{1 + \sum_{j \neq k} e^{-\beta[(\mathbf{x}-\mathbf{m}_j)^2 - (\mathbf{x}-\mathbf{m}_k)^2]}} = \begin{cases} 1 & \text{for } i = k \\ 0 & \text{for } i \neq k \end{cases} \end{aligned}$$

Hence for $\beta \rightarrow \infty$ the SVQ modification of the code-book vectors in a single iteration reduces to

$$\begin{aligned} \mathbf{m}_k &\rightarrow \mathbf{m}_k + \eta(\mathbf{x} - \mathbf{m}_k) \\ \mathbf{m}_i &\rightarrow \mathbf{m}_i \quad \text{for all } i \neq k \end{aligned}$$

which is identical to that of VQ. \square

Fact 2: For $\beta \rightarrow 0$ all code-book vectors \mathbf{m}_i will ultimately collapse to a single point, which will fluctuate around the average data-point $\langle \mathbf{x} \rangle = \int d\mathbf{x} \mathbf{x} p(\mathbf{x})$.

Proof: For $\beta \rightarrow 0$ we find $F_i[\mathbf{x}, \{\mathbf{m}\}] = N^{-1}$ (for any \mathbf{x} , any i and any $\{\mathbf{m}\}$) and the SVQ modifications simply reduce to

$$\mathbf{m}_i \rightarrow \mathbf{m}_i + \frac{\eta}{N}(\mathbf{x} - \mathbf{m}_i) \quad \text{for all } i$$

Now consider the difference between any two code-book vectors. If we label the iterations by $\ell = 0, 1, 2, \dots$, we find

$$\mathbf{m}_i(\ell+1) - \mathbf{m}_j(\ell+1) = \left(1 - \frac{\eta}{N}\right)[\mathbf{m}_i(\ell) - \mathbf{m}_j(\ell)]$$

$$\text{hence} \quad \mathbf{m}_i(\ell) - \mathbf{m}_j(\ell) = \left(1 - \frac{\eta}{N}\right)^\ell [\mathbf{m}_i(0) - \mathbf{m}_j(0)]$$

so $\lim_{\ell \rightarrow \infty} [\mathbf{m}_i(\ell) - \mathbf{m}_j(\ell)] = \mathbf{0}$ \square

To find out where all the code-book vectors will go, we only need to inspect the dynamics of the average $\mathbf{m} = N^{-1} \sum_i \mathbf{m}_i$ (since $\lim_{\ell \rightarrow \infty} [\mathbf{m}(\ell) - \mathbf{m}_i(\ell)] = \mathbf{0}$). Upon writing the data point drawn at iteration ℓ as $\mathbf{x}(\ell)$, we get $\mathbf{m}(\ell+1) = \left(1 - \frac{\eta}{N}\right)\mathbf{m}(\ell) + \frac{\eta}{N}\mathbf{x}(\ell)$. Hence

$$\begin{aligned} \mathbf{m}(1) &= \left(1 - \frac{\eta}{N}\right)\mathbf{m}(0) + \frac{\eta}{N}\mathbf{x}(0) \\ \mathbf{m}(2) &= \left(1 - \frac{\eta}{N}\right) \left[\left(1 - \frac{\eta}{N}\right)\mathbf{m}(0) + \frac{\eta}{N}\mathbf{x}(0) \right] + \frac{\eta}{N}\mathbf{x}(1) \\ &= \left(1 - \frac{\eta}{N}\right)^2 \mathbf{m}(0) + \frac{\eta}{N} \left[\left(1 - \frac{\eta}{N}\right)\mathbf{x}(0) + \mathbf{x}(1) \right] \\ \mathbf{m}(3) &= \left(1 - \frac{\eta}{N}\right) \left[\left(1 - \frac{\eta}{N}\right)^2 \mathbf{m}(0) + \frac{\eta}{N} \left(1 - \frac{\eta}{N}\right)\mathbf{x}(0) + \frac{\eta}{N}\mathbf{x}(1) \right] + \frac{\eta}{N}\mathbf{x}(2) \\ &= \left(1 - \frac{\eta}{N}\right)^3 \mathbf{m}(0) + \frac{\eta}{N} \left[\left(1 - \frac{\eta}{N}\right)^2 \mathbf{x}(0) + \left(1 - \frac{\eta}{N}\right)\mathbf{x}(1) + \mathbf{x}(2) \right] \end{aligned}$$

$$\begin{aligned} & \vdots = \vdots \\ \mathbf{m}(\ell) &= (1 - \frac{\eta}{N})^\ell \mathbf{m}(0) + \frac{\eta}{N} \sum_{k=0}^{\ell-1} (1 - \frac{\eta}{N})^{\ell-1-k} \mathbf{x}(k) \end{aligned}$$

Averaging over the possible choices of data, using $\langle \mathbf{x}(k) \rangle = \langle \mathbf{x} \rangle$, the average of the distribution $p(\mathbf{x})$ (assuming this average to be finite), and using $\sum_{k \geq 0} z^k = 1/(1-z)$ (for $|z| < 1$) then gives us

$$\begin{aligned} \lim_{\ell \rightarrow \infty} \langle \mathbf{m}(\ell) \rangle &= \lim_{\ell \rightarrow \infty} \left\{ (1 - \frac{\eta}{N})^\ell \mathbf{m}(0) + \frac{\eta}{N} \langle \mathbf{x} \rangle \sum_{k \geq 0} (1 - \frac{\eta}{N})^k \right\} \\ &= \frac{\eta}{N} [1 - (1 - \frac{\eta}{N})]^{-1} \langle \mathbf{x} \rangle = \langle \mathbf{x} \rangle \end{aligned}$$

As claimed. \square

One can similarly show that the fluctuations of \mathbf{m} around $\langle \mathbf{x} \rangle$ remain finite, provided the width of the data distribution is finite, i.e. $\langle \mathbf{x}^2 \rangle = \int d\mathbf{x} p(\mathbf{x}) \mathbf{x}^2 < \infty$.

Thus VQ can be seen as a special case of SVQ, obtained upon putting $\beta \rightarrow \infty$. Secondly we infer from inspection of the extreme case $\beta = 0$ that one of the effects of the smoothening of SVQ (relative to VQ) is for code-vectors to ‘drag one another along’.

A Lyapunov Function for Small Learning Rates. Our understanding of SVQ (and thus also of VQ) would greatly improve if we could recognise the process as the minimisation of some error measure. For finite η this is not possible, but for $\eta \rightarrow 0$ it is. Upon labeling the different iterations with $\ell = 0, 1, 2, \dots$, and upon writing the data point drawn at stage ℓ as $\mathbf{x}(\ell)$, we can write both VQ and SVQ in the following compact form:

$$\mathbf{m}_i(\ell + 1) = \mathbf{m}_i(\ell) + \eta F_i[\mathbf{x}(\ell), \{\mathbf{m}(\ell)\}] [\mathbf{x}(\ell) - \mathbf{m}_i(\ell)] \quad (2.4)$$

with

$$F_i^{\text{SVQ}}[\mathbf{x}, \{\mathbf{m}\}] = \frac{e^{-\beta(\mathbf{x} - \mathbf{m}_i)^2}}{\sum_{j=1}^N e^{-\beta(\mathbf{x} - \mathbf{m}_j)^2}} \quad (2.5)$$

$$F_i^{\text{VQ}}[\mathbf{x}, \{\mathbf{m}\}] = \begin{cases} 1 & \text{if } \mathbf{x} \in V_i[\{\mathbf{m}\}] \\ 0 & \text{otherwise} \end{cases} \quad (2.6)$$

For small learning rates we follow the procedure introduced to derive deterministic equations for on-line learning in layered neural networks (see lecture notes ‘Neural Networks’), and define a new time unit $t = \eta\ell$. We then take the $\eta \rightarrow 0$ limit and find the stochastic process (2.4) being replaced by the deterministic equation

$$\frac{d}{dt} \mathbf{m}_i = \langle (\mathbf{x} - \mathbf{m}_i) F_i[\mathbf{x}, \{\mathbf{m}\}] \rangle \quad (2.7)$$

(with $\langle f(\mathbf{x}) \rangle = \int d\mathbf{x} p(\mathbf{x}) f(\mathbf{x})$). Let us now choose $F_i^{\text{SVQ}}[\mathbf{x}, \{\mathbf{m}\}]$, and let us define the following average of Gaussian distributions, each centred at one of the code-book vectors:

$$q(\mathbf{x} | \{\mathbf{m}\}) = \frac{1}{N} \sum_i \frac{e^{-\beta(\mathbf{x} - \mathbf{m}_i)^2}}{(\sqrt{\pi/\beta})^n} \quad (2.8)$$

We can now prove the following:

Fact 3: For $F_i[\mathbf{x}, \{\mathbf{m}\}] = F_i^{\text{SVQ}}[\mathbf{x}, \{\mathbf{m}\}]$, equation (2.7) defines gradient descent:

$$\frac{d}{dt}\mathbf{m}_i = -\frac{1}{2\beta}\nabla_{\mathbf{m}_i}E[\{\mathbf{m}\}] \quad E[\{\mathbf{m}\}] = \int d\mathbf{x} p(\mathbf{x}) \log \left[\frac{p(\mathbf{x})}{q(\mathbf{x}|\{\mathbf{m}\})} \right] \quad (2.9)$$

Proof: Just work out the relevant partial derivatives of $E[\{\mathbf{m}\}]$:

$$\begin{aligned} -\frac{1}{2\beta}\nabla_{\mathbf{m}_i}E[\{\mathbf{m}\}] &= \frac{1}{2\beta} \int d\mathbf{x} p(\mathbf{x}) \nabla_{\mathbf{m}_i} \log q(\mathbf{x}|\{\mathbf{m}\}) \\ &= \frac{1}{2\beta} \int d\mathbf{x} p(\mathbf{x}) \nabla_{\mathbf{m}_i} \log \sum_j e^{-\beta(\mathbf{x}-\mathbf{m}_j)^2} \\ &= \frac{1}{2\beta} \int d\mathbf{x} p(\mathbf{x}) \left\{ \frac{\nabla_{\mathbf{m}_i} e^{-\beta(\mathbf{x}-\mathbf{m}_i)^2}}{\sum_j e^{-\beta(\mathbf{x}-\mathbf{m}_j)^2}} \right\} \\ &= \int d\mathbf{x} p(\mathbf{x}) \left\{ \frac{(\mathbf{x}-\mathbf{m}_i)e^{-\beta(\mathbf{x}-\mathbf{m}_i)^2}}{\sum_j e^{-\beta(\mathbf{x}-\mathbf{m}_j)^2}} \right\} = \langle (\mathbf{x}-\mathbf{m}_i)F_i^{\text{SVQ}}[\mathbf{x}, \{\mathbf{m}\}] \rangle \end{aligned}$$

Hence the desired result. \square

The object in equation (2.9) is a familiar function in information theory, where it would have been written as $D(p||q)$, and where it goes under the name Kullback-Leibler distance. It obeys $D(p||q) \geq 0$ for any two distributions $p(\mathbf{x})$ and $q(\mathbf{x})$, with equality only if $p(\mathbf{x}) = q(\mathbf{x})$ (in a distributional sense). It is an information-theoretic measure of the deviation between the two distributions $p(\mathbf{x})$ and $q(\mathbf{x})$.

From the general inequality $D(p||q) \geq 0$ (and hence $E[\{\mathbf{m}\}] \geq 0$), in combination with the gradient descent equation in (2.9) (which ensures that $\frac{d}{dt}E \leq 0$), we may now conclude that the function in (2.9) is a Lyapunov function for the SVQ process (2.4). Hence we can interpret SVQ, at least for small η , as approximating the data distribution $p(\mathbf{x})$ optimally by a mixture of Gaussians of the form (2.8), via adaptation of the centres \mathbf{m}_i of the Gaussians. Similarly, since $\lim_{\beta \rightarrow \infty} \text{SVQ} = \text{VQ}$ and since $\lim_{\beta \rightarrow \infty} [\beta/\pi]^{\frac{n}{2}} e^{-\beta(\mathbf{x}-\mathbf{m}_i)^2} = \delta[\mathbf{x}-\mathbf{m}_i]$, we may for small η interpret VQ as approximating the data distribution $p(\mathbf{x})$ optimally by a mixture of delta-distributions, via adaptation of the centres of the delta-distributions:

$$\begin{aligned} \eta \ll 1, \text{ SVQ:} \quad & \text{finds } \{\mathbf{m}\} \text{ such that} \quad p(\mathbf{x}) \approx \frac{1}{N} \sum_i \frac{e^{-\beta(\mathbf{x}-\mathbf{m}_i)^2}}{(\sqrt{\pi}/\beta)^n} \\ \eta \ll 1, \text{ VQ:} \quad & \text{finds } \{\mathbf{m}\} \text{ such that} \quad p(\mathbf{x}) \approx \frac{1}{N} \sum_i \delta[\mathbf{x}-\mathbf{m}_i] \end{aligned}$$

The characteristic distance $\sigma = 1/\sqrt{2\beta}$ which we already identified before, is now giving the widths of the individual Gaussians with which SVQ aims to copy the data distribution $p(\mathbf{x})$. For finite η (provided not too large) we may regard SVQ and VQ as noisy versions of the above gradient descent processes.

Finally we briefly note that the stationary state of the $\eta \rightarrow 0$ equation (2.7), given by $\langle (\mathbf{x}-\mathbf{m}_i)F_i[\mathbf{x}, \{\mathbf{m}\}] \rangle = \mathbf{0}$, translates into

$$\mathbf{m}_i = \frac{\int d\mathbf{x} \mathbf{x} F_i[\mathbf{x}, \{\mathbf{m}\}] p(\mathbf{x})}{\int d\mathbf{x} F_i[\mathbf{x}, \{\mathbf{m}\}] p(\mathbf{x})} \quad \text{so for VQ:} \quad \mathbf{m}_i = \frac{\int_{V_i[\{\mathbf{m}\}]} d\mathbf{x} \mathbf{x} p(\mathbf{x})}{\int_{V_i[\{\mathbf{m}\}]} d\mathbf{x} p(\mathbf{x})} \quad (2.10)$$

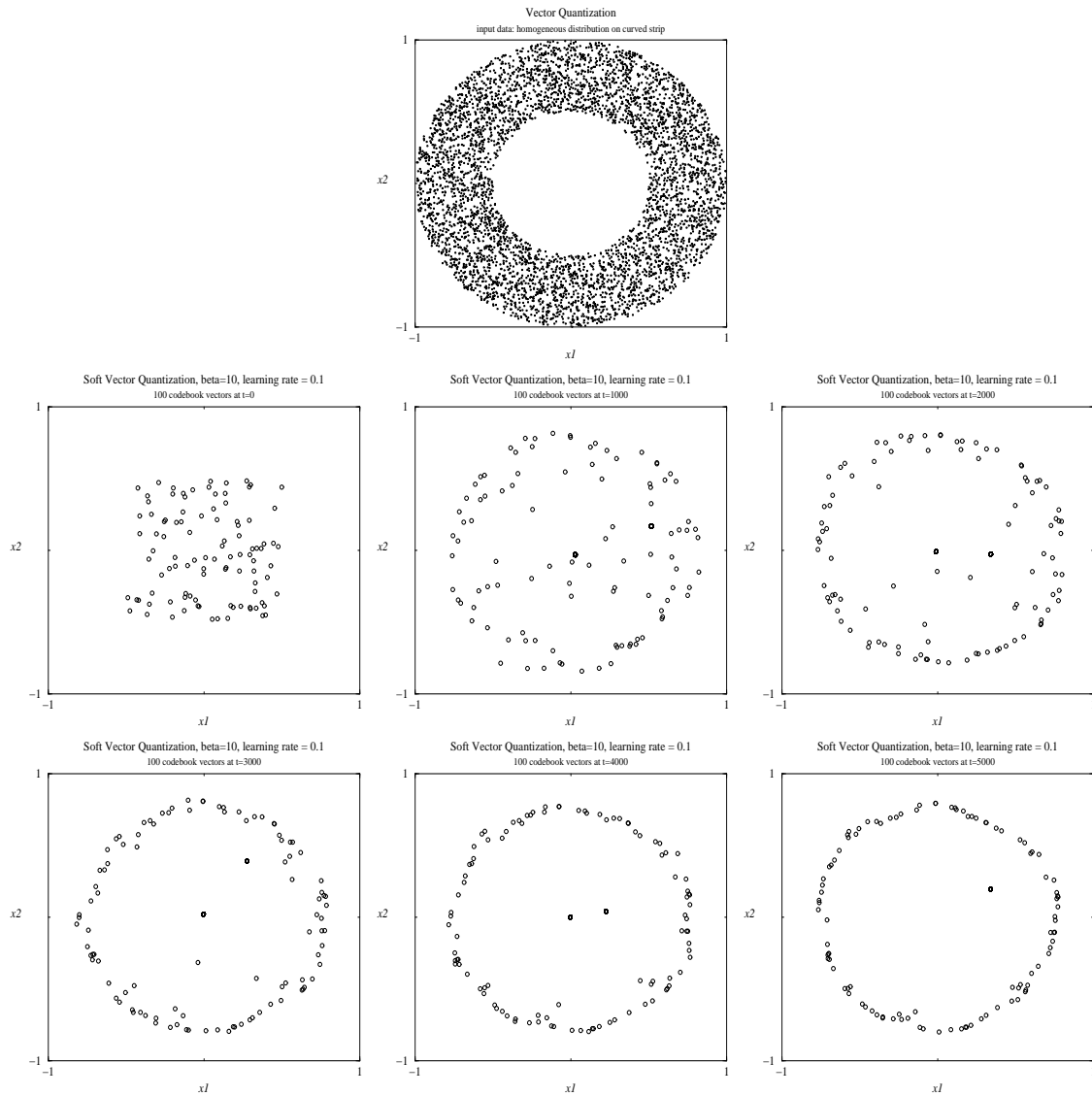


Figure 2.10: Example: numerical simulation of Soft Vector Quantization, with $n = 2$, $\eta = 0.1$, $\beta = 10$ and $N = 100$. Top graph: data distribution $p(\mathbf{x})$, homogeneously distributed over the region $\frac{1}{2} < |\mathbf{x}| < 1$. Bottom six graphs: locations of the code-book vectors \mathbf{m}_i during the course of the process ($i = 1, \dots, 100$), at times $t = 0$ (middle row, left), $t = 1000$ (middle row, center), $t = 2000$ (middle row, right), $t = 3000$ (bottom row, left), $t = 4000$ (bottom row, middle) and $t = 5000$ (bottom row, right). Times are measures in the number of iterations. Initialisation of code-book vector allocations: randomly in the square $[-\frac{1}{2}, \frac{1}{2}] \times [-\frac{1}{2}, \frac{1}{2}]$.

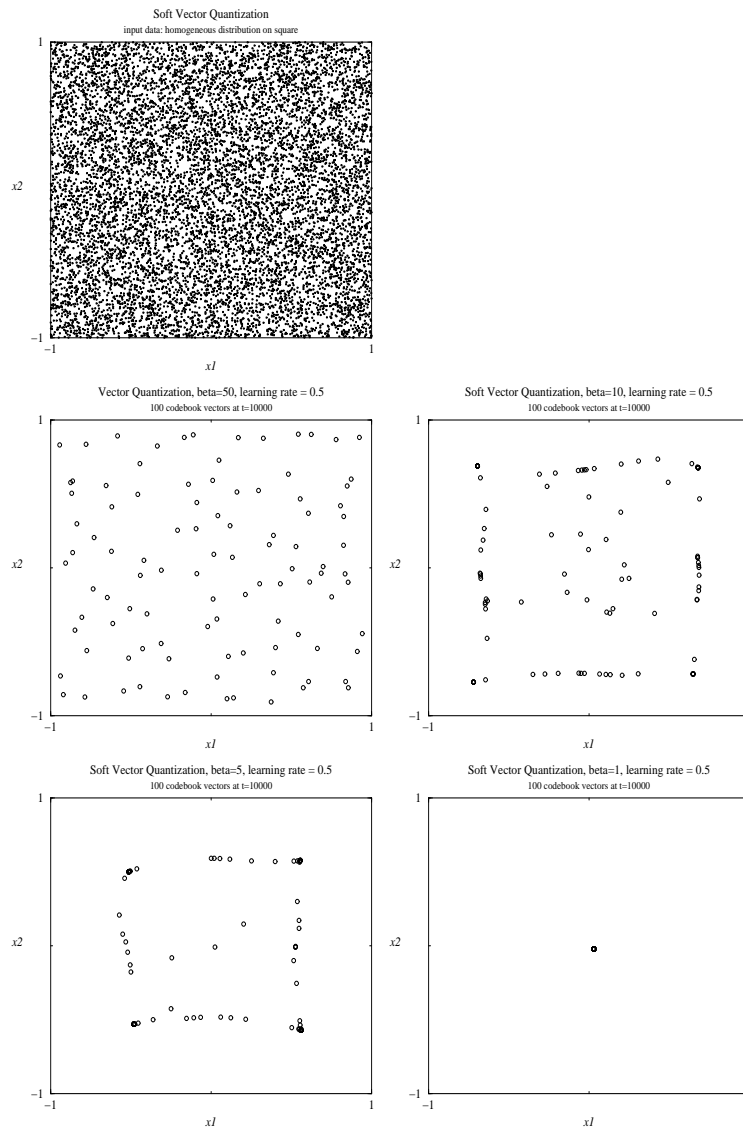


Figure 2.11: Example: numerical simulation of Soft Vector Quantization, with $n = 2$, $\eta = 0.5$ and $N = 100$. Top graph: data distribution $p(\mathbf{x})$, homogeneously distributed over the square $[-1, 1] \times [-1, 1]$. Bottom four graphs: locations of the code-book vectors \mathbf{m}_i after 10,000 iteration steps ($i = 1, \dots, 100$), for different choices of the parameter β : $\beta = 50$ (middle row, left), $\beta = 10$ (middle row, right), $\beta = 5$ (bottom row, left) and $\beta = 1$ (bottom row, right). Initialization of code-book vector allocations: randomly in the square $[-\frac{1}{2}, \frac{1}{2}] \times [-\frac{1}{2}, \frac{1}{2}]$.

For VQ we thus have the transparent result that, in equilibrium and for η sufficiently small, the location of code-book vector \mathbf{m}_i will be exactly the centre of gravity of the data distribution within its associated Voronoi cell.

Examples of SVQ in Action. The figures on the previous two pages illustrate the functioning of the SVQ algorithm for a number of simple examples with $n = 2$ (i.e. where one has data points $\mathbf{x} = (x_1, x_2)$ in a plane) and $N = 100$ (i.e. a population of 100 code-book vectors), but for different choices of the data distribution $p(\mathbf{x})$. Figure 2.10 shows the locations of the code-book vectors at different times, for $\beta = 10$, where the characteristic width of the data covering Gaussians is $\sigma = 1/\sqrt{2\beta} \approx 0.22$. We observe that code-book vectors indeed do not get ‘stuck’ in data-poor regions, and that the code-book vectors indeed keep a distance of the order of σ from the boundaries of the data region (due to the width of the data strip being $\frac{1}{2}$, this more or less sends all code-book vectors to the circle $|\mathbf{x}| = 3/4$). Figure 2.11 shows the asymptotic locations (after 10,000 iterations) of the code-book vectors for four different values of β . The corresponding values of σ are: $\sigma = 0.1$ ($\beta = 50$), $\sigma \approx 0.22$ ($\beta = 10$), $\sigma \approx 0.32$ ($\beta = 5$), and $\sigma \approx 0.71$ ($\beta = 1$). With these values one can understand perfectly the observed clustering properties of the code-book vectors; this underlines the power of results such as those above (i.e. fact 3).

2.3 Time-Dependent Learning Rates

We have seen that for finite η the processes of the type (2.4) can be regarded as ‘noisy’ versions of the deterministic equation (2.7) (for SVQ and VQ, the latter, in turn, minimise a sensible error measure). Hence one would in the initial stages of these processes prefer to have a finite η (since the added randomness may prevent the system from going to a sub-optimal local minimum of the error measure), but one would asymptotically prefer a fluctuation-free and unique (reproducible) final state, i.e. $\eta \rightarrow 0$. It would therefore appear natural to choose a slowly but monotonically decreasing time-dependent learning rate $\eta(\ell)$, where ℓ labels the iterations of the algorithm (2.4). The question is how to determine the optimal rate of decay for $\eta(\ell)$. Too fast a reduction of $\eta(\ell)$ might cause evolution to local minima, or might prevent the code-book vectors from traveling over sufficiently large distances to take up their optimal positions in data space (note that we do not know beforehand the size of data space). Too slow a reduction of $\eta(\ell)$ might cause fluctuations and non-uniqueness of the final state to persist too long for us to achieve the desired stationary state within the time-scales of our experiments. Below we show that sensible criteria to be met by $\eta(\ell)$ are

$$\text{not too fast : } \sum_{\ell=0}^{\infty} \eta(\ell) = \infty \quad \text{not too slow : } \sum_{\ell=0}^{\infty} \eta^2(\ell) < \infty \quad (2.11)$$

Fact 1: The learning rate must obey $\sum_{\ell=0}^{\infty} \eta(\ell) = \infty$, otherwise there will be an a priori bound on the possible distance to be covered by the code-book vectors.

Proof: From $\mathbf{m}_i(\ell + 1) = \mathbf{m}_i(\ell) [1 - \eta F_i[\mathbf{x}(\ell), \{\mathbf{m}(\ell)\}]] + \eta F_i[\mathbf{x}(\ell), \{\mathbf{m}(\ell)\}]\mathbf{x}(\ell)$, which is just an alternative way to write equation (2.4), together with $F_i[\mathbf{x}, \{\mathbf{m}\}] \in [0, 1]$ and $0 \leq \eta \leq 1$, we deduce for time dependent learning rates:

$$|\mathbf{m}_i(\ell + 1)| \leq |\mathbf{m}_i(\ell)| + \eta(\ell)|\mathbf{x}(\ell)|$$

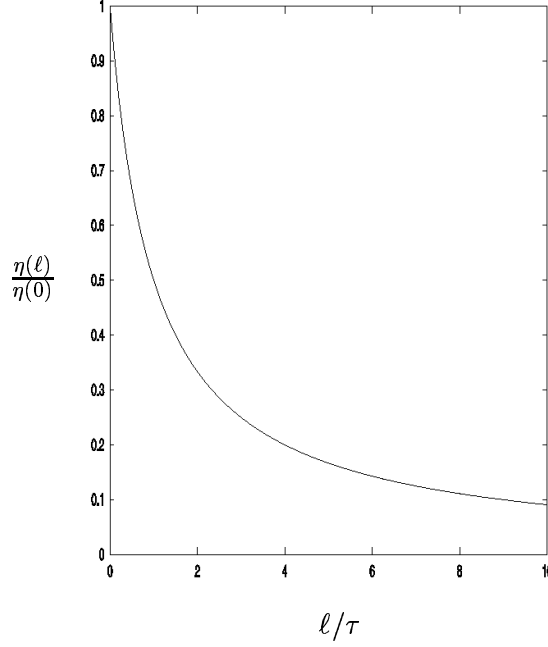


Figure 2.12: The time-dependent learning rate defined in (2.12), which meets the conditions (2.11) on the asymptotic rate of decay to zero.

Hence, upon repeated iteration:

$$|\mathbf{m}_i(\ell)| \leq |\mathbf{m}_i(0)| + \sum_{k=0}^{\ell-1} \eta(k) |\mathbf{x}(k)| \quad \text{so} \quad \lim_{\ell \rightarrow \infty} \langle |\mathbf{m}_i(\ell)| \rangle \leq |\mathbf{m}_i(0)| + \langle |\mathbf{x}| \rangle \sum_{k=0}^{\infty} \eta(k)$$

Thus, unless $\sum_{k=0}^{\infty} \eta(k) = \infty$, we have an a priori bound on the distance which can be traveled by any code-book vector, which could prevent outliers of the data distribution from being reached. \square

Fact 2: The learning rate must obey $\sum_{\ell=0}^{\infty} \eta^2(\ell) < \infty$, otherwise even for the simplest process of the class (2.4) the uncertainty in the location of the average code-book vector will diverge.

Proof: The simplest process of the class (2.4) is the $\beta \rightarrow 0$ limit of SVQ: $F_i[\mathbf{x}, \{\mathbf{m}\}] = 1/N$, so $\mathbf{m}(\ell+1) = \mathbf{m}(\ell) + (\eta(\ell)/N)[\mathbf{x}(\ell) - \mathbf{m}(\ell)]$, where $\mathbf{m}(\ell) = \frac{1}{N} \sum_i \mathbf{m}_i(\ell)$. We now write $\mathbf{m}(\ell) = \langle \mathbf{m}(\ell) \rangle + \mathbf{v}(\ell)$, so that $\langle \mathbf{v}(\ell) \rangle = \mathbf{0}$ and

$$\mathbf{v}(\ell+1) = \left[1 - \frac{\eta(\ell)}{N}\right] \mathbf{v}(\ell) + \frac{\eta(\ell)}{N} [\mathbf{x}(\ell) - \langle \mathbf{x} \rangle]$$

Note that $\mathbf{v}(\ell)$ is statistically independent of all $\mathbf{x}(\ell')$ with $\ell' \geq \ell$. Note also that, since at $\ell = 0$ there are no fluctuations yet: $\mathbf{v}(0) = \mathbf{0}$. Hence

$$\begin{aligned} \mathbf{v}^2(\ell+1) &= \left[1 - \frac{\eta(\ell)}{N}\right]^2 \mathbf{v}^2(\ell) + \frac{\eta^2(\ell)}{N^2} [\mathbf{x}(\ell) - \langle \mathbf{x} \rangle]^2 + \frac{2\eta(\ell)}{N} \left[1 - \frac{\eta(\ell)}{N}\right] \mathbf{v}(\ell) \cdot [\mathbf{x}(\ell) - \langle \mathbf{x} \rangle] \\ \langle \mathbf{v}^2(\ell+1) \rangle &= \left[1 - \frac{\eta(\ell)}{N}\right]^2 \langle \mathbf{v}^2(\ell) \rangle + \frac{\eta^2(\ell)}{N^2} \langle [\mathbf{x} - \langle \mathbf{x} \rangle]^2 \rangle \end{aligned}$$

Further iteration gives

$$\begin{aligned}
\langle \mathbf{v}^2(\ell+2) \rangle &= \left[1 - \frac{\eta(\ell+1)}{N}\right]^2 \left[1 - \frac{\eta(\ell)}{N}\right]^2 \langle \mathbf{v}^2(\ell) \rangle \\
&\quad + \left[1 - \frac{\eta(\ell+1)}{N}\right]^2 \frac{\eta^2(\ell)}{N^2} \langle [\mathbf{x} - \langle \mathbf{x} \rangle]^2 \rangle + \frac{\eta^2(\ell+1)}{N^2} \langle [\mathbf{x} - \langle \mathbf{x} \rangle]^2 \rangle \\
&= \left[1 - \frac{\eta(\ell+1)}{N}\right]^2 \left[1 - \frac{\eta(\ell)}{N}\right]^2 \langle \mathbf{v}^2(\ell) \rangle \\
&\quad + \left[\frac{\eta^2(\ell+1)}{N^2} + \frac{\eta^2(\ell)}{N^2} + \mathcal{O}(\eta^3) \right] \langle [\mathbf{x} - \langle \mathbf{x} \rangle]^2 \rangle \\
&\quad \vdots \\
\langle \mathbf{v}^2(\ell) \rangle &= \langle \mathbf{v}^2(0) \rangle \prod_{k=0}^{\ell-1} \left[1 - \frac{\eta(k)}{N}\right]^2 + \langle [\mathbf{x} - \langle \mathbf{x} \rangle]^2 \rangle \sum_{k=0}^{\ell-1} \left[\frac{\eta^2(k)}{N^2} + \mathcal{O}(\eta^3) \right] \\
&= \langle [\mathbf{x} - \langle \mathbf{x} \rangle]^2 \rangle \sum_{k=0}^{\ell-1} \left[\frac{\eta^2(k)}{N^2} + \mathcal{O}(\eta^3) \right]
\end{aligned}$$

We conclude that the fluctuations \mathbf{v} diverge, i.e. $\lim_{\ell \rightarrow \infty} \langle \mathbf{v}^2(\ell) \rangle = \infty$, as soon as $\sum_{\ell=0}^{\infty} \eta^2(\ell) = \infty$. As claimed. \square

The most commonly used asymptotic form of decay for $\eta(\ell)$, which meets the above two requirements, is $\eta(\ell) \sim \ell^{-1}$ as $\ell \rightarrow \infty$. For instance:

$$\eta(\ell) = \frac{\eta(0)}{1 + \ell/\tau} \quad \text{which obeys} \quad \begin{aligned} \ell \ll \tau &: \eta(\ell) \approx \eta(0)[1 - \ell/\tau] \\ \ell \approx \tau &: \eta(\ell) \approx \frac{1}{2}\eta(0) \\ \ell \gg \tau &: \eta(\ell) \approx \eta(0)\tau/\ell \end{aligned} \quad (2.12)$$

This dependence is drawn in figure 2.12.

2.4 Self-Organizing Maps

Detour: Biological Inspiration. Any flexible and robust autonomous system (whether living or robotic) will have to be able to create, or at least update, an internal ‘map’ or representation of its environment. Information on its environment, however, is usually obtained in an indirect manner, through a redundant set of usually highly non-linear sensors which each provide only partial and indirect information. The system responsible for forming this map needs to be adaptive, as both environment and sensors can change their characteristics during the system’s life-time. Our brain performs recalibration of sensors all the time; e.g. simply because we grow will the neuronal information about limb positions (generated by sensors which measure the stretch of muscles) have to be reinterpreted continually. Anatomic changes, and even learning new skills (like playing an instrument), are found to induce modifications of internal maps. At a more abstract level, the neural system responsible for building the internal representation is confronted with a complicated non-linear mapping from a relatively low-dimensional and more or less flat space (the ‘physical world’ W , with D dimensions) into a high-dimensional one (the space of sensory signals, of which there are $n \gg D$), and the aim is to extract the properties of the original space W from the sensory signals alone, i.e.

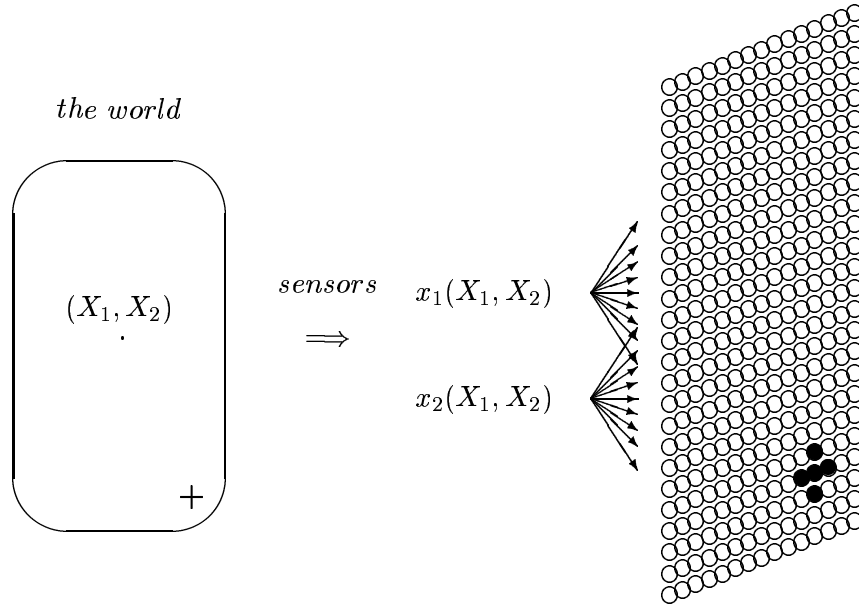


Figure 2.13: Simple 2-dim. example of the internal representation (‘map’) of the environment as built in biological systems. The brain has no direct access to the world coordinates (X_1, X_2) , it is only given sensory signals $x_1(X_1, X_2)$ and $x_2(X_1, X_2)$. Many representations of the coordinates (X_1, X_2) would have been possible; in the specific one aimed for here, each (indirectly) observed position (X_1, X_2) gives rise to a specific localised area of firing activity in a sheet of neurons. Compare this to driving a car with blinded windows, but with electronic sensors observing the environment. The sensory information is to be converted back into world coordinates, and to make an light flash up in a road map inside the car, exactly at the car’s current location. This is the information used by the driver. The question discussed in this section is how to *construct* the road map and the conversion of sensory to world coordinates, on the basis of structure in the sensory signals $x_1(X_1, X_2)$ and $x_2(X_1, X_2)$.

to ‘invert’ the operation which maps world states to sensory signals.

$$\begin{array}{lll}
 \textit{world coordinates} : & \textit{sensory signals} : & \textit{internal representation} : \\
 \mathbf{X} = (X_1, \dots, X_D) \in W & \longrightarrow \mathbf{x} = (x_1, \dots, x_n) & \longrightarrow \text{reconstruction of } W ?
 \end{array}$$

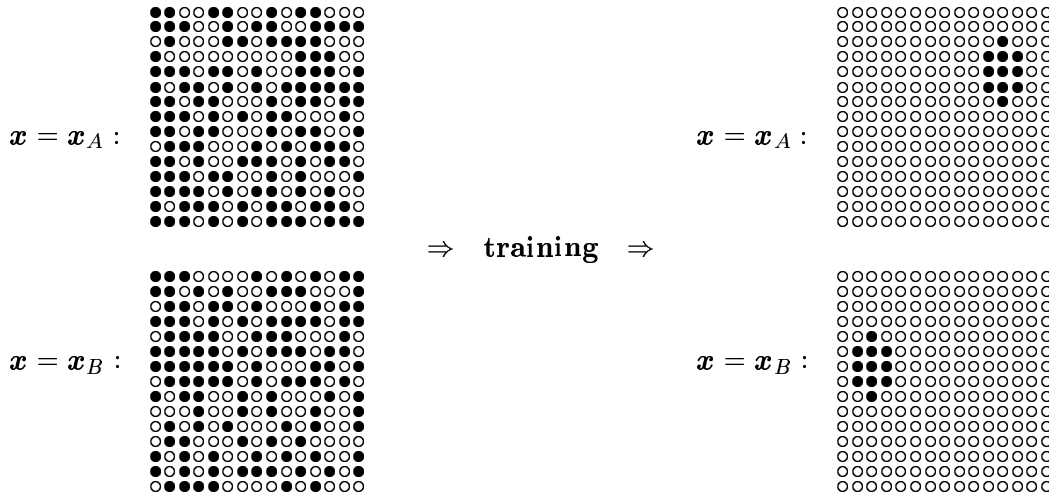
The definition of the ‘world’ depends on the specific degrees of freedom of the system at hand, e.g. D need not be three. The sensory signals are usually non-linear, indirect and noisy, but also highly redundant to compensate for their individual inadequacies. The key to achieving the objective is to exploit continuity and correlations in sensory signals, assuming similar sensory signals to represent similar positions in the environment, which therefore must correspond to similar positions in the internal map: if $\mathbf{X}' = \mathbf{X} + \Delta\mathbf{X}$ then also $\mathbf{x}(\mathbf{X}') \approx \mathbf{x}(\mathbf{X}) + \mathbf{A}(\mathbf{X}) \cdot \Delta\mathbf{X}$, for some $\mathbf{A}(\mathbf{X})$ (the Jacobian matrix of the mapping $\mathbf{X} \rightarrow \mathbf{x}$).

Before turning to the learning process which is to carry out the construction task, we first explain in more detail the type of representation one aims for (which is inspired by biological systems). Let us give a simple example (see also figure 2.13). Image a system operating in

a simple two-dimensional world, where positions are represented by Cartesian coordinates $\mathbf{X} = (X_1, X_2)$, observed by sensors and fed into a neural network as input signals. For simplicity we will consider trivial sensors, and just put $x_i = X_i$. Each neuron i receives information on the input signals (x_1, x_2) in the usual way, via modifiable synapses $\{m_{ij}\}$: $y_i(\mathbf{x}) = g[m_{i1}x_1 + m_{i2}x_2]$, for some monotonic non-linear function $g[z]$ (e.g. $g[z] = \text{erf}[z]$ or $g[z] = \frac{1}{2}(1 + \tanh[z])$). We indicate the (physical) location of neuron i in the array as \mathbf{r}_i . If this network is to become an internal coordinate system, faithfully reflecting the events $\mathbf{x} = (x_1, x_2)$ observed in the outside world (in the present example its topology must accordingly be that of a two-dimensional array), the following objectives are to be met

1. Each neuron y_i is more or less ‘tuned’ to a specific type of signal \mathbf{x}_i , i.e. $y_i(\mathbf{x}) > 0$ only when $|\mathbf{x} - \mathbf{x}_i|$ is sufficiently small.
2. Neighbouring neurons in the array are tuned to similar signals, i.e. if $|\mathbf{r}_i - \mathbf{r}_j|$ is small, than $|\mathbf{x}_i - \mathbf{x}_j|$ is small.
3. External ‘distance’ is monotonically related to internal ‘distance’, i.e. if $|\mathbf{r}_i - \mathbf{r}_j| < |\mathbf{r}_i - \mathbf{r}_k|$, than $|\mathbf{x}_i - \mathbf{x}_j| < |\mathbf{x}_i - \mathbf{x}_k|$.

How to Learn A Topologically Correct Map. It turns out that in order to achieve these objectives one needs learning rules where neurons effectively enter a competition for having input signals \mathbf{x}_i ‘allocated’ to them, whereby neighbouring neurons stimulate one another to develop similar synaptic interactions and distant neurons are prevented from developing similar interactions. Let us try to construct the simplest such learning rule.



Initially, before learning has taken place, input signals just evoke random firing patterns in the neuronal array. After learning we wish to see localised firing patterns, as shown above. Our equations take their simplest form in the case where the input signals are normalised, so we define $(x_1, x_2) \in [-1, 1]^2$ and add a dummy variable $x_3 = \sqrt{1 - x_1^2 - x_2^2}$, together with a dummy synaptic interaction m_{i3} for every neuron. If we write the synapses of neuron i as $\mathbf{m}_i = (m_{i1}, m_{i2}, m_{i3})$, and normalise them according to $|\mathbf{m}_i| = 1$, we simply get

$$y(\mathbf{x}) = g[\mathbf{m}_i \cdot \mathbf{x}], \quad |\mathbf{m}_i| = |\mathbf{x}| = 1 \tag{2.13}$$

Learning implies rotating the vectors \mathbf{m}_i . By construction (since $g[z]$ is a monotonically increasing function), a given neuron is triggered most when $\mathbf{m}_i \cdot \mathbf{x}$ is maximal, which is for the input signal $\mathbf{x} = \mathbf{m}_i$. This has two important implications: one can now interpret \mathbf{m}_i as the input signal to which neuron i is ‘tuned’, and ‘tuning’ neuron i to a given signal \mathbf{x} is found to be equivalent to rotating \mathbf{m}_i towards \mathbf{x} . As a result one finds that a learning rule with the desired effect is, starting from random (normalised) synaptic vectors \mathbf{m}_i , to iterate the following recipe until a (more or less) stable situation is reached:

step 1: pick a data point \mathbf{x} with $|\mathbf{x}| = 1$ at random

step 2: find the most active neuron, i.e. the i such that $|\mathbf{m}_i \cdot \mathbf{x}| > |\mathbf{m}_j \cdot \mathbf{x}|$ for all $j \neq i$

step 3: rotate the synaptic vector of neuron i and its neighbours in the array towards \mathbf{x} , rotate the synaptic vector of all other neurons slightly away from \mathbf{x}

step 4: return to 1

The neuron that was already the one most responsive to the signal \mathbf{x} will be made even more so (together with its neighbours); the other neurons are made less responsive to \mathbf{x} . In biology the above is achieved via Hebbian-type learning rules, and by having short-range excitatory synapses which ensure that neighbouring neurons tend to be simultaneously active, and ‘team up’ in the synaptic adaptation.

We note that in the above formulation the normalisation of synaptic vectors and input vectors was purely introduced in order to retain the connection with the biological picture of firing neurons, via expressions like (2.13). For the purpose of computation and learning in synthetic systems, on the other hand, we can move away from the firing of neurons and work directly in terms of the $\{\mathbf{m}_i\}$ only. Henceforth we will define \mathbf{m}_i simply as the signal in input space to which neuron i is tuned; we forget about how this is actually realised, and we forget about normalisation.

Visualization of the Learning Process: Fishing Nets. The above formulation of the self-organizing maps in terms of the $\{\mathbf{m}_i\}$ has brought us (deliberately) very close the VQ and SVQ algorithms. However, there is an important difference, which is the insistence that the various code-book vectors \mathbf{m}_i are no longer inter-changable (in VQ and SVQ we only cared about the overall distribution of code-book vectors in input space), but now each \mathbf{m}_i is associated with a specific location \mathbf{r}_i in an array, with the condition that those $\{i, j\}$ which belong to nearby locations $\{\mathbf{r}_i, \mathbf{r}_j\}$ in the array must have code-book vectors $\{\mathbf{m}_i, \mathbf{m}_j\}$ which are very similar. Hence, simply drawing the code-book vectors in input space (as we did for VQ and SVQ) is no longer adequate, because it will not inform us about the topological features of the state.

In the case of SOM’s the standard visual representation of a state $\{\mathbf{m}_i\}$ is the following:

- Draw each \mathbf{m}_i as a point (‘knot’) in input space (as in VQ and SVQ)
- Connect with line segments all those points $\{i, j\}$ which correspond to neighbours in the array of the $\{\mathbf{r}_i\}$

We then end up with a graphical representation of the synaptic structure of a network in the form of a ‘fishing net’, with the positions of the knots representing the signals in the world

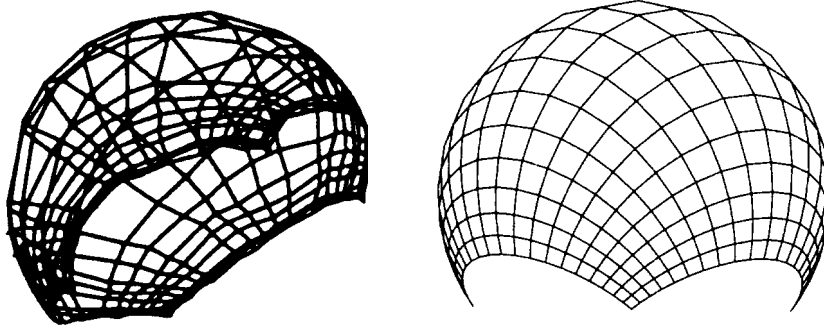


Figure 2.14: Graphical representation of the state of a self-organizing map (SOM) in the form of a ‘fishing net’. Here $\mathbf{r}_i \in [0, 1]^2$ (i.e. the neurons live in a square 2-dimensional array) and $\mathbf{x} \in \mathfrak{R}^3$, with $|\mathbf{x}| = 1$ (i.e. the data live on the surface of a sphere). The positions of the knots represent the signals \mathbf{m}_i to which the neurons are ‘tuned’, and the cords connect the knots of neighbouring neurons in the array. Left: equilibrium result of the numerical simulation of a particular version on the SOM algorithm. Right: corresponding theoretical prediction.

to which the neurons are tuned and with the cords indicating neighbourhood, see figures 2.14. The three objectives of map formation set out at the beginning of this section thereby translate into the following desired properties of the graph:

1. all knots in the graph are separated
2. all cords of the graph are similarly stretched
3. there are no regions with overlapping pieces of graph

Note that this representation is in practice useful only for two-dimensional arrays (i.e. when creating two-dimensional internal representations).

The SOM Algorithm and its Properties. The Self-Organising Map (SOM) algorithm is an abstract realisation of the processes underlying the creation of topologically correct internal representations in the higher brain regions of humans and animals. It is defined in terms of ‘code-book’ vectors \mathbf{m}_i , which represent the specific signals in data space to which neurons are ‘tuned’. It aims to create a topologically correct internal representation of low-dimensional structure (potentially non-linear) hidden in (possibly high-dimensional) data.

Each neuron is located in a physical array, the dimension of which will become the dimension of the internal map created; the vector \mathbf{r}_i indicates the location of neuron i . We define a neighbourhood function h_{ij} as follows:

$$h_{ij} = h\left(\frac{|\mathbf{r}_i - \mathbf{r}_j|}{\sigma}\right), \quad \sigma > 0 \quad (2.14)$$

in which $h(z)$ is a monotonically decreasing function, with $h(0) = 1$, $h(\infty) = 0$, and with a width of order 1. For example:

$$h(z) = e^{-\frac{1}{2}z^2} : \quad h_{ij} = e^{-\frac{1}{2}(\mathbf{r}_i - \mathbf{r}_j)^2 / \sigma^2}$$

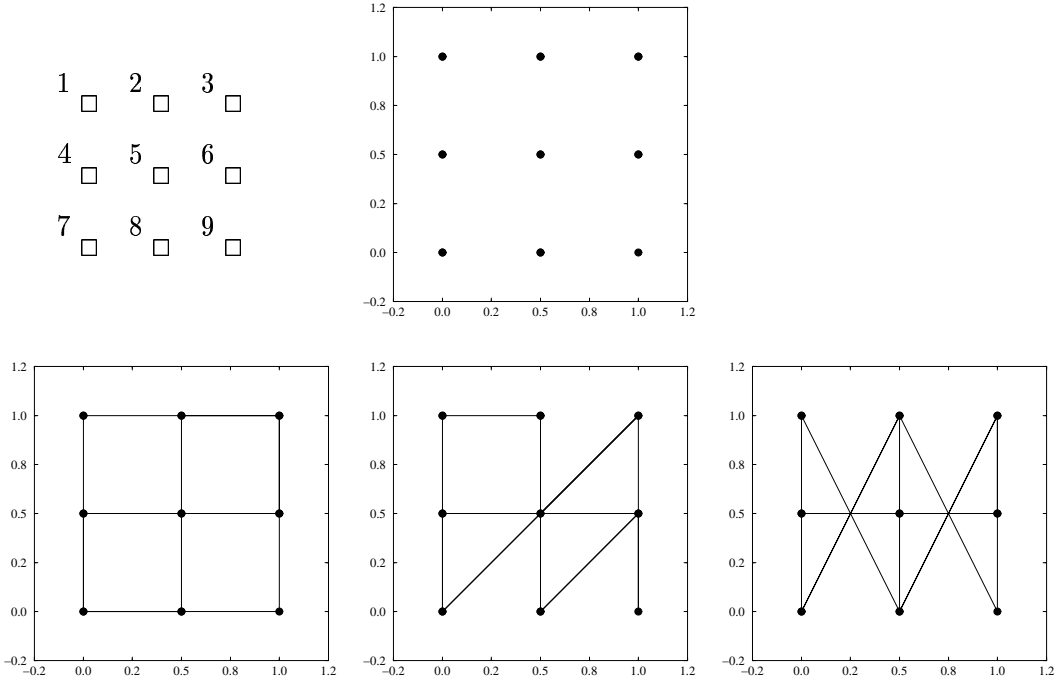


Figure 2.15: Example of the ‘fishing net’ representation of the state of a self-organizing map system. Upper left: the arrangement of 9 neurons in a physical array (defining the $\{\mathbf{r}_i\}$, and hence the notion of neighbours). Upper middle: the positioning of the nine code-book vectors \mathbf{m}_i in data space, corresponding to the state $\{\mathbf{m}_1 = (0, 1)$, $\mathbf{m}_2 = (\frac{1}{2}, 1)$, $\mathbf{m}_3 = (1, 1)$, $\mathbf{m}_4 = (0, \frac{1}{2})$, $\mathbf{m}_5 = (\frac{1}{2}, \frac{1}{2})$, $\mathbf{m}_6 = (1, \frac{1}{2})$, $\mathbf{m}_7 = (0, 0)$, $\mathbf{m}_8 = (\frac{1}{2}, 0)$, $\mathbf{m}_9 = (1, 0)\}$. Bottom graph left: the fishing net graph corresponding to the aforementioned state. Bottom graph middle: the fishing net graph corresponding to the state one obtains by exchanging code-book vectors \mathbf{m}_3 and \mathbf{m}_8 . Bottom graph right: the fishing net graph corresponding to the state one obtains by exchanging code-book vectors \mathbf{m}_2 and \mathbf{m}_8 . Note that all cases would give rise to exactly the same distribution of code-book vectors, so that the incorrect topology of the last two states can indeed only be inferred from the fishing net graph.

$$h(z) = \theta[1 - z] : \quad \begin{cases} h_{ij} = 1 & \text{if } |\mathbf{r}_i - \mathbf{r}_j| < \sigma \\ h_{ij} = 0 & \text{if } |\mathbf{r}_i - \mathbf{r}_j| > \sigma \end{cases}$$

in which $\theta[z]$ is the step function ($\theta[z > 0] = 1$, $\theta[z < 0] = 0$). The SOM algorithm can now be defined as follows. First we initialize the N code-book vectors $\mathbf{m}_i \in \mathfrak{R}^n$ randomly. Then we iterate the following stochastic process:

step 1: pick a data point \mathbf{x} at random, according to the probability density $p(\mathbf{x})$

step 2: find the Voronoi cell containing \mathbf{x} , i.e. find i such that $|\mathbf{x} - \mathbf{m}_i| < |\mathbf{x} - \mathbf{m}_j| \forall j \neq i$

step 3: move *all* \mathbf{m}_j towards \mathbf{x} : $\mathbf{m}_j \rightarrow \mathbf{m}_j + \eta(\mathbf{x} - \mathbf{m}_j)F_j[\mathbf{x}, \{\mathbf{m}(\ell)\}]$, where $F_j[\dots] = h_{ij}$

step 4: return to 1

The learning rate $0 < \eta \ll 1$ controls, as always, the overall magnitude of the changes.

The SOM algorithm can be seen to have the following general properties:

- (i) All code-book vectors are moved towards \mathbf{x} at every iteration step, but those which are closest to the code-book vector in whose Voronoi cell the data point is (i.e. that of the neuron which is triggered most) are moved most.
- (ii) Unless we reduce η during the process, the code-book vectors will continue to move stochastically, although the *density* of code-book vectors in any given region should become stationary.
- (iii) The properties of the $\{h_{ij}\}$ guarantee that neighbouring neurons ‘team up’, and develop similar code-book vectors.
- (iv) SOM has one parameter more than VQ. This extra parameter, σ , defines a characteristic distance in the original physical array of the neurons: code-book vectors with $|\mathbf{r}_i - \mathbf{r}_j| < \sigma$ will be the ones to feel one another’s influence.
- (v) $\lim_{\sigma \rightarrow 0} \text{SOM} = \text{VQ}$ (the proof is trivial). Hence we can see the SOM as VQ plus enforced spatial continuity.

We observe that the SOM is again of the general form (2.4), hence we can apply our results on time-dependent learning rates also to the SOM. The relation between the three processes discussed so far is

$$\lim_{\beta \rightarrow \infty} \text{SVQ} = \text{VQ} = \lim_{\sigma \rightarrow 0} \text{SOM}$$

Applications of the SOM include sensor fusion (in robotics), data visualisation and data-base mining (finding hidden regularities in messy data), data preprocessing, medical diagnostics (finding causal relationships in medical data), etc.

It will be clear that the determination of the right value for σ is delicate. If σ is too small, the system will behave like VQ, and the correct topology will not emerge (or be correct only locally). If σ is too large, all neurons will drag one another along, and all code-book vectors will become identical (see also the simulation examples below). The new parameter σ also introduces boundary effects. Since neurons ‘drag along’ one another’s code-book vectors, it will be clear that those neurons which are at the boundary of the array will feel an effective force dragging their code-book vectors towards those of the inner region of the array (since there are no neurons pulling from the outside). This shows itself in a tendency for the boundaries of the ‘fishing net’ to keep a certain σ -dependent distance from the boundaries of the data region. In practice one therefore chooses a time-dependent σ , similar to the time-dependent learning rate, i.e.

$$\eta(t) = \frac{\eta(0)}{1 + t/\tau_\eta} \quad \sigma(t) = \frac{\sigma(0)}{1 + t/\tau_\sigma}$$

Examples of SOM in Action. The figures below illustrate the functioning of the SOM algorithm for a number of simple examples with $n = 2$ (i.e. where one has data points $\mathbf{x} = (x_1, x_2)$ in a plane) and a 10×10 array of neurons (i.e. a population of 100 code-book vectors), but for different choices of the distance parameter σ and the time-dependence of the learning rate.

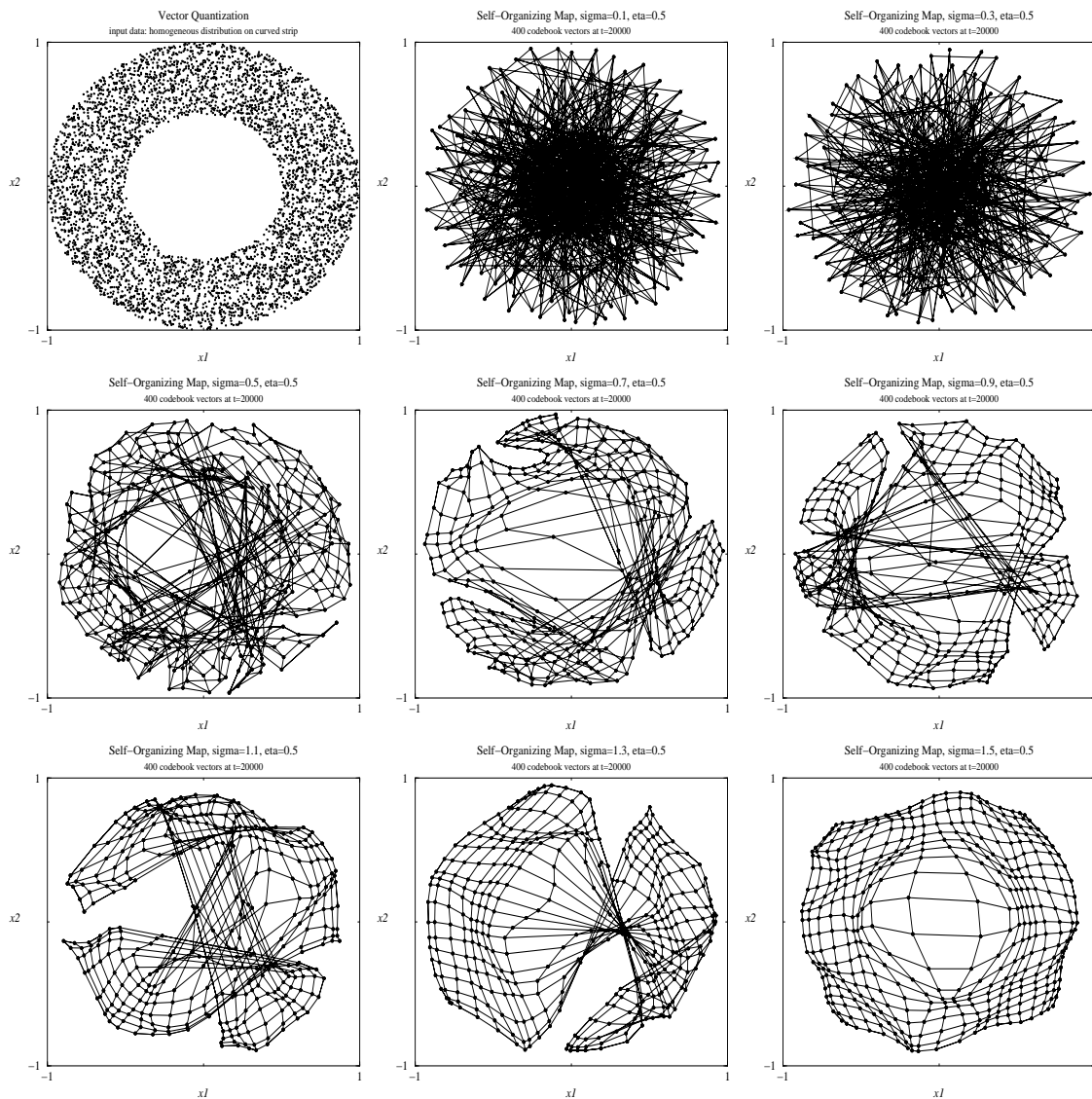


Figure 2.16: Example: numerical simulation of a Self-Organising Map, with $n = 2$, $\eta = 0.5$ and $N = 100$. Top left graph: data distribution $p(\mathbf{x})$, homogeneously distributed over the strip $\frac{1}{2} < |\mathbf{x}| < 1$. Other eight graphs: system state after 20000 iteration steps, for different values of the range parameter σ ($\sigma \in \{0.1, 0.3, 0.5, 0.7, 0.9, 1.1, 1.3, 1.5\}$, from top middle to bottom right). One clearly observes the two effects of losing the correct topology for σ too small, with only local correctness for intermediate σ , and of the code-book vectors of boundary neurons keeping a σ -dependent distance from the data boundary.

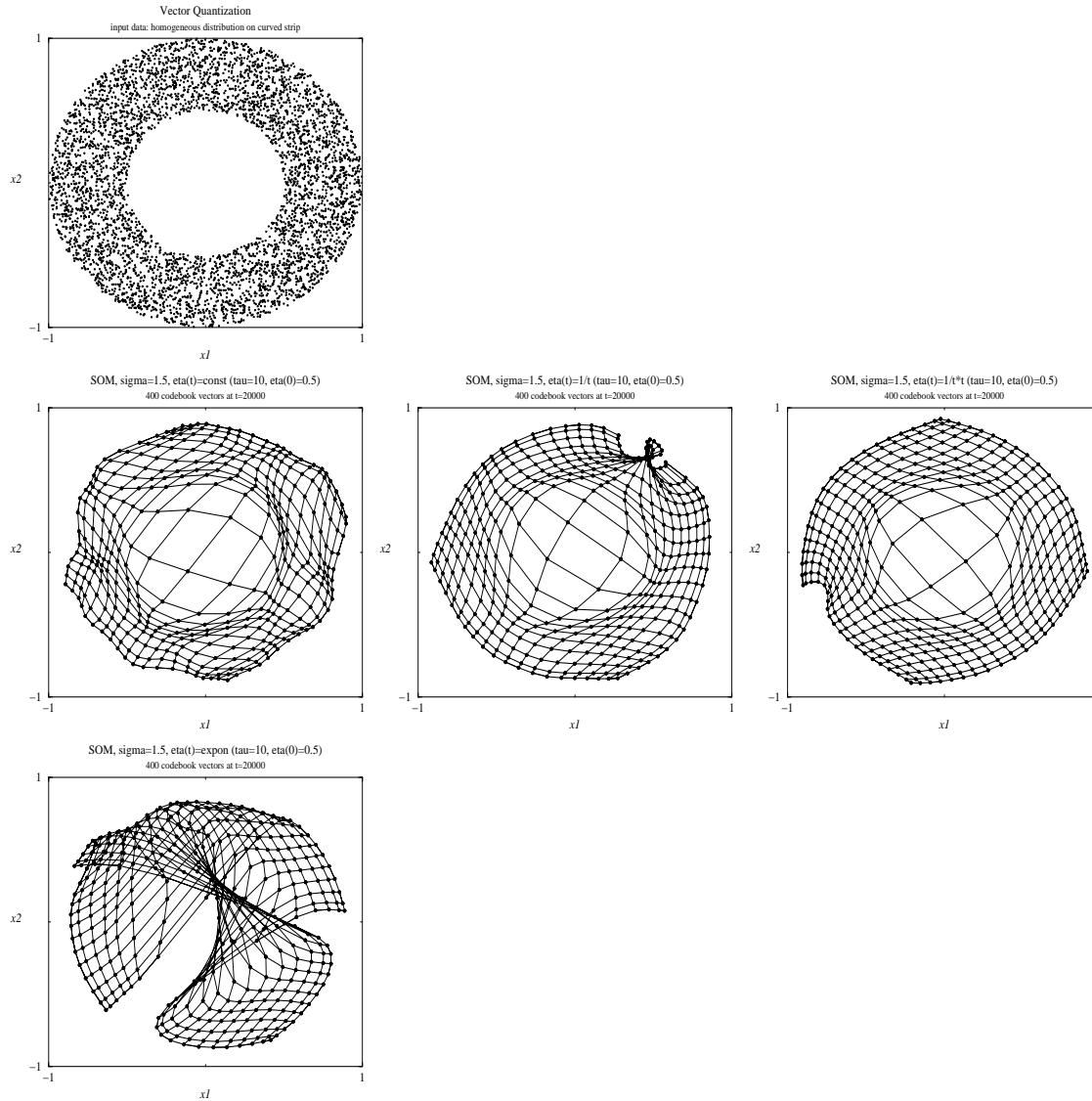


Figure 2.17: Example: numerical simulation of a Self-Organising Map, with $n = 2$, $\sigma = 1.5$, $\eta(0) = 0.5$, and $N = 100$. Top left graph: data distribution $p(\mathbf{x})$, homogeneously distributed over the strip $\frac{1}{2} < |\mathbf{x}| < 1$. Other eight graphs: system state after 20000 iteration steps, for different choices of the time-dependence of the learning rate η . Second row left: $\eta = 0.5$ (constant). Second row, middle: $\eta(t) = \eta(0)/[1 + t/\tau]$, with $\tau = 10$. Second row, right: $\eta(t) = \eta(0)/[1 + t^2/\tau^2]$, with $\tau = 10$. Bottom: $\eta(t) = \eta(0) \exp[-t/\tau]$, with $\tau = 10$. One clearly observes that the system can get ‘stuck’ in a suboptimal state, where the topology is only locally correct.

Chapter 3

Bayesian Techniques in Supervised Learning

In this chapter we discuss the Bayesian approach to supervised learning from examples. Its introduction in the 1990's has been a very welcome development, in that it contributed significantly to converting synthetic neural computation, especially learning in multi-layer networks, from a collection of interesting but ad hoc algorithms ('invent a rule, play with parameters, draw performance graphs, and hope for the best') to a sound, well grounded, and systematic scientific procedure. It also contributed to the area becoming more mathematical in nature.

3.1 Preliminaries and Introduction

Supervised Learning from Examples. Let us first set the scene and define our terminology. In supervised learning from examples we are confronted with a 'task', defined by a collection of questions ξ^μ , drawn randomly from some set $D \subseteq \mathfrak{R}^N$ (with probabilities, or probability density, $p(\xi)$, with corresponding answers t^μ :

$$\begin{array}{l} \text{the task :} \quad \text{'questions' :} \quad \{\xi^1, \dots, \xi^p\}, \quad \xi^\mu \in D \subseteq \mathfrak{R}^N \\ \quad \quad \quad \text{'answers' :} \quad \{t^1, \dots, t^p\}, \quad t^\mu \in \mathfrak{R} \end{array}$$

This task, assumed to have been generated by a 'teacher', is to be learned by a 'student' (the neural network). The student executes a parametrized operation $S : \mathfrak{R}^N \rightarrow \mathfrak{R}$, where $S(\xi) = f(\xi; \mathbf{w})$. The parameters \mathbf{w} determine the details of the operation, and thus represent the 'program'; in multi-layer neural networks they would be the synaptic weights and the neuronal thresholds. If the outputs (or 'targets') are binary, e.g. $t^\mu \in \{0, 1\}$ or $t^\mu \in \{-1, 1\}$ we would call the task *binary classification*. If the outputs can truly take values from a continuous set, we would speak about *regression*.

It is assumed that the data (or 'answers') were not assigned randomly by the teacher (otherwise there would be no point in learning), but that they were generated according to some function $T : \mathfrak{R}^N \rightarrow \mathfrak{R}$. The student has no direct information on the function T (the teacher is a 'black box', or 'oracle'), but has to infer T from the p input-output pairs (ξ^μ, t^μ) ; this is the objective of the learning process. The problem faced by the student is made more difficult by the fact (which is quite realistic in the real world) that the data are not perfect,

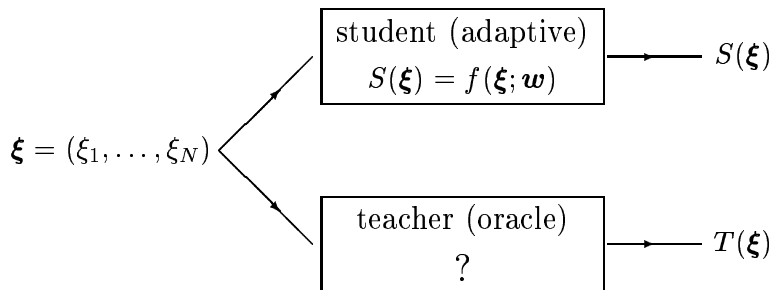


Figure 3.1: The ‘teacher-student’ scenario of on-line supervised learning. The student can be any type of parametrised operation (not necessarily neural or neural-inspired).

i.e. the teacher can make mistakes, or is subject to noise. If, for simplicity, we assume that the inaccuracy or noise is independent for the different outputs (i.e. the teacher is sloppy, rather than consistently using the wrong rule), we have

$$\begin{aligned}
 \text{binary classification :} & \quad \text{Prob}[t^\mu = T(\xi^\mu)] = 1 - \epsilon \quad 0 \leq \epsilon \leq 1 \\
 \text{regression :} & \quad t^\mu = T(\xi^\mu) + z_\mu \quad z_\mu \in \mathfrak{R} \quad \text{drawn randomly} \\
 & \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \text{according to } P(z)
 \end{aligned} \tag{3.1}$$

In the case of binary classification, the parameter ϵ measures the amount of noise in the data, with $\epsilon = 0$ corresponding to perfect data. In the case of regression we may without loss of generality assume $\langle z \rangle = \int dz z P(z) = 0$ (since a non-zero average would be equivalent to a structural modification of the underlying rule T , and could be absorbed in its definition); hence the amount of noise in the data is here measured by the variance of $P(z)$, with perfect data given by $P[z] = \delta[z]$.

Due to the possibility of data noise and the difference between our finite sample $\{\xi^1, \dots, \xi^p\}$ and the full set of possible questions D , there are several performance measures one could define. Here we concentrate on two: the *training error* E_t (measuring how well the student reproduces the given answers t^μ), and the *generalization error* E_g (measuring how well the student has learned the underlying rule T):

$$E_t = \frac{1}{p} \sum_{\mu=1}^p \mathcal{E}[t^\mu, S(\xi^\mu)] \quad E_g = \int_D d\xi p(\xi) \mathcal{E}[T(\xi), S(\xi)]$$

where $\mathcal{E}[t, s]$ is some, as yet unspecified, measure of the difference between t and s , such as $\mathcal{E}[t, s] = |t - s|^\gamma$ (with $\gamma > 0$). Given the dependence of the student’s operation on the parameters \mathbf{w} , via $S(\xi) = f(\xi, \mathbf{w})$, this becomes

$$E_t[\mathbf{w}] = \frac{1}{p} \sum_{\mu=1}^p \mathcal{E}[t^\mu, f(\xi^\mu; \mathbf{w})] \quad E_g[\mathbf{w}] = \int_D d\xi p(\xi) \mathcal{E}[T(\xi), f(\xi; \mathbf{w})] \tag{3.2}$$

There are three important observations to make at this stage:

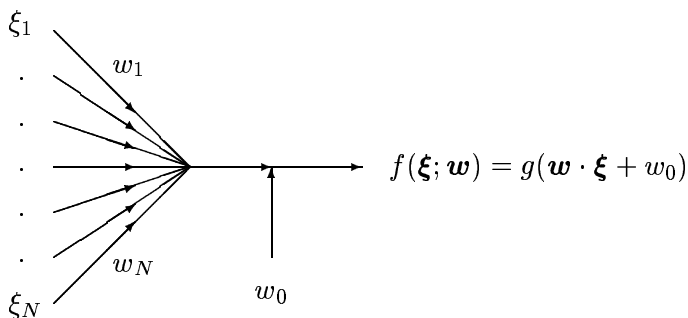
- The above two error measures $E_t[\mathbf{w}]$ and $E_g[\mathbf{w}]$ differ in two aspects: (i) E_t is an average only over the p available questions $\{\xi^1, \dots, \xi^p\}$, whereas E_g involves *all* questions in D . (ii) E_t checks performance relative to the (noisy) answers $\{t^\mu\}$ given, whereas E_g checks performance relative to the correct underlying rule T .

- The real objective of learning is to minimise $E_g[\mathbf{w}]$, via suitable modification of the parameters \mathbf{w} , since we ultimately wish to use our system for predicting the answers to novel questions, rather than just parrot the answers to those of the training stage.
- In the learning process, however, we have no access to $E_g[\mathbf{w}]$, since we know neither the full set D nor the true rule T ; we only have the data $\{(\boldsymbol{\xi}^1, t^1), \dots, (\boldsymbol{\xi}^p, t^p)\} \dots$

Many of the obstacles and problems in neural network learning in the past (i.e. before, say, 1990) had their origin in the fact that one wished to minimize one object (E_g), but one could only measure the other (E_t), and one tried to get away (out of necessity, it seemed) with pretending that the differences between the two were not too important.

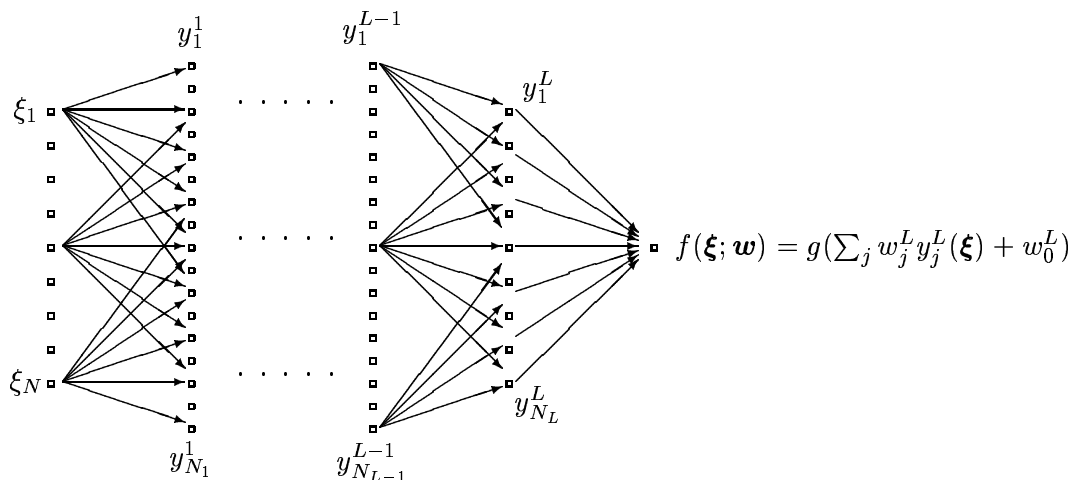
Conventional Network Types. Let us next briefly review the most common types of neural network $S(\boldsymbol{\xi}) = f(\boldsymbol{\xi}, \mathbf{w})$ which are being used in practice for the above purpose. We first consider regression tasks, where the outputs are supposed to be continuous:

- Single layer neural networks (perceptrons):



These are basically single neurons, carrying out a (soft) linear separation in the space D of the input vectors $\boldsymbol{\xi}$, of the form $S(\boldsymbol{\xi}) = g(\mathbf{w} \cdot \boldsymbol{\xi} + w_0)$. Here $g(z)$ is a monotonically increasing sigmoidal function, which saturates for $z \rightarrow \pm\infty$. The modifiable parameters (i.e. the ‘program’) are the synaptic weights (w_1, \dots, w_N) and the threshold w_0 .

- Multi-layer neural networks (MLP’s):

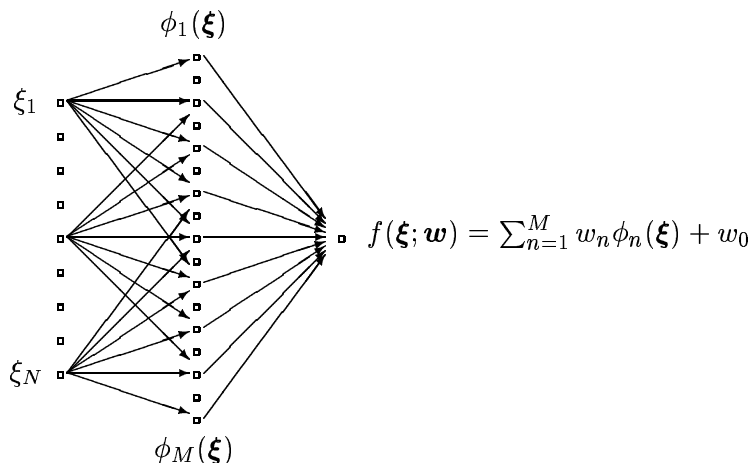


These are feed-forward processing networks of neurons of the type described above (soft linear separators), with monotonically increasing and saturating sigmoidal functions $g(z)$. The overall output is calculated iteratively from the outputs of the preceding layer, according to

$$y_i^1 = g[\sum_{j=1}^N w_{ij}^0 \xi_j + w_0^0] \quad y_i^{\ell+1} = g[\sum_{j=1}^{N_\ell} w_{ij}^\ell y_j^\ell + w_0^\ell] \quad f(\boldsymbol{\xi}; \mathbf{w}) = g[\sum_{j=1}^{N_L} w_j^L y_j^L + w_0^L] \quad (3.3)$$

and the modifiable parameters (i.e. the ‘program’) are the synaptic weights and thresholds of all L layers, i.e. $\mathbf{w} = \{w_{ij}^\ell, w_i^L, w_0^\ell\}$. Given a sufficient number of layers and units within each layer (dependent on the problem to be solved), these networks have been proven to be ‘universal approximators’: any sufficiently well-behaved continuous function $f : \mathfrak{R}^N \rightarrow \mathfrak{R}$ can be approximated to arbitrary accuracy by MLP’s (we will not discuss the precise conditions or the proof here).

- Radial basis function networks (RBF):



The motivation behind these networks is also the principle underlying the previous multi-layer networks: one converts a non-linearly separable problem into a linearly separable one, by suitable pre-processing. In the previous network the final neuron does the linear separation, and the preceding L layers do the pre-processing. In radial basis function networks one chooses specific functions $\phi_n(\boldsymbol{\xi})$ to do a specific type of pre-processing:

$$f(\boldsymbol{\xi}; \mathbf{w}) = \sum_{n=1}^M w_n \phi_n(\boldsymbol{\xi}) + w_0 \quad \phi_n(\boldsymbol{\xi}) = \phi(|\boldsymbol{\xi} - \mathbf{x}^n|) \quad (3.4)$$

The $\phi_n(\boldsymbol{\xi})$ (the ‘radial basis functions’) are allowed to depend only on the distance $|\boldsymbol{\xi} - \mathbf{x}^n|$ between the input vector and a code-book vector \mathbf{x}^n . For the remaining scalar function ϕ to be specified several choices have been proposed; which one to pick is crucial for performance,

but (unfortunately) highly dependent on the problem at hand¹. Popular choices are:

$$\begin{aligned} \text{localised basis functions :} & \quad \phi(u) = e^{-u^2/2\sigma^2} \\ & \quad \phi(u) = (u^2 + \sigma^2)^{-\alpha} \quad (\alpha > 0) \\ \text{non-localised basis functions :} & \quad \phi(u) = u \\ & \quad \phi(u) = (u^2 + \sigma^2)^\beta \quad (0 < \beta < 1) \end{aligned}$$

The modifiable parameters are the weights \mathbf{w} (and sometimes one also allows the code-book vectors \mathbf{x}^n to adapt). Note, finally, that in RBF networks the output does not undergo a non-linear (sigmoidal) operation, hence the output depends linearly on the parameters \mathbf{w} .

Model Complexity and Overfitting. Systems such as those discussed above have been used (and are still being used) extensively as ‘student networks’ $S(\boldsymbol{\xi}) = f(\boldsymbol{\xi}; \mathbf{w})$ in the sense of figure 3.1. The strategy for finding the best parameters \mathbf{w} was, in the early stages of neural computation, based on performing gradient descent on the surface defined by the training error $E_t[\mathbf{w}]$ as given in (3.2), with a quadratic error measure $\mathcal{E}[z] = \frac{1}{2}z^2$:

$$\frac{d}{dt}\mathbf{w} = -\eta \nabla_{\mathbf{w}} E_t[\mathbf{w}] \quad E_t[\mathbf{w}] = \frac{1}{2p} \sum_{\mu=1}^p [t^\mu - f(\boldsymbol{\xi}^\mu; \mathbf{w})]^2 \quad (3.5)$$

with a learning rate $\eta > 0$. The objective is to find \mathbf{w}^* such that $\min_{\mathbf{w}} E_t[\mathbf{w}] = E_t[\mathbf{w}^*]$; the above process would at least be guaranteed to lead us to a local minimum of $E_t[\mathbf{w}]$, since

$$\frac{d}{dt}E_t[\mathbf{w}] = \nabla_{\mathbf{w}} E_t[\mathbf{w}] \cdot \frac{d}{dt}\mathbf{w} = -\eta [\nabla E_t[\mathbf{w}]]^2 \leq 0$$

$$\frac{d}{dt}\mathbf{w} = \mathbf{0} \quad \iff \quad \nabla_{\mathbf{w}} E_t[\mathbf{w}] = \mathbf{0}$$

One problem with this approach has been clear from the beginning: one can end in a local rather than a global minimum. A second problem was appreciated only some ten years later: gradient descent is not the optimal local process to find minima on an error surface (see the chapter on information geometry in the Information Theory course). Here we concentrate on a third problem, called inappropriate model complexity, or ‘overfitting’.

This problem is best explained using a very simple example. Suppose we have a task $T : \mathfrak{X} \rightarrow \mathfrak{X}$ defined by a teacher which carries out some function $g(x)$. The student $S : \mathfrak{X} \rightarrow \mathfrak{X}$ tries to emulate the teacher by some parametrized function $f(x; \mathbf{w})$, for which we choose a finite polynomial (i.e. truncated Taylor expansion):

$$f(x; \mathbf{w}) = w_0 + w_1 x + \dots + w_M x^M \quad \text{or} \quad f(x; \mathbf{w}) = \mathbf{w} \cdot \boldsymbol{\phi}(x), \quad \phi_\ell(x) = x^\ell \quad (3.6)$$

Learning means adapting $\mathbf{w} = \{w_0, \dots, w_M\}$. The data consist of examples of p inputs x_μ and corresponding outputs t^μ (with $p \geq M + 1$), however, the teacher is not perfect (there is a

¹This illustrates how in RBF networks the learning problem is, in a way, simply moved under the carpet, for how do we now choose the function $\phi(u)$, the code-book vectors \mathbf{x}^n , or the number M ? And under which conditions can we be sure that the ad-hoc restriction to RBF units does not render the problem unlearnable? At least in image processing, however, one can set up a rigorous framework, and show that all non-pathological images can be decomposed in terms of a bundle of local expansions (derivatives of Gaussian functions), which can (in turn) be constructed by adding and subtracting Gaussian functions of the type in (3.4).

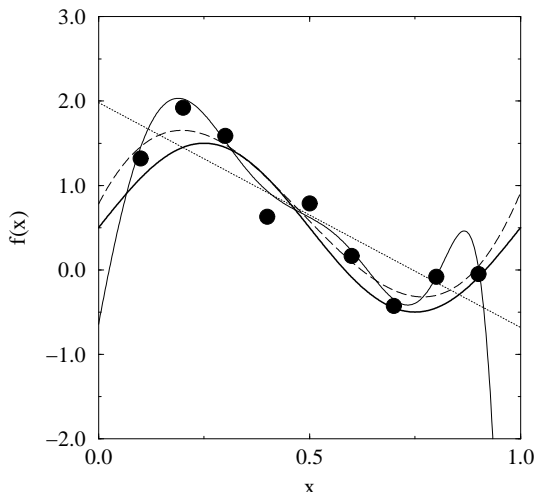


Figure 3.2: Example of a simple learning process where a system ‘learns’ via gradient descent on a quadratic error measure how to approximate a function $g(x) = \frac{1}{2} + \sin(2\pi x)$ (thick solid line) with by an order- M polynomial $f(x; \mathbf{w}) = w_0 + w_1 x + \dots + w_M x^M$, on the basis of nine noisy sample points of this function (circles). The resulting polynomials are shown for $M - 1$ (dotted), $M = 3$ (dashed) and $M = 9$ (thin solid).

noise source in the data generation), so $t_\mu = g(x_\mu) + z_\mu$ where the z_μ are drawn independently according to a zero-average distribution $P(z)$. The training error is given by

$$E_t[\mathbf{w}] = \frac{1}{2} \sum_{\mu=1}^p [t^\mu - \mathbf{w} \cdot \boldsymbol{\phi}(x_\mu)]^2$$

Gradient descent learning then boils down to

$$\frac{1}{\eta} \frac{d}{dt} w_i = \sum_{\mu=1}^p [t^\mu - \mathbf{w} \cdot \boldsymbol{\phi}(x_\mu)] \phi_i(x_\mu) = \sum_{\mu=1}^p t^\mu \phi_i(x_\mu) - \sum_{j=1}^M \left[\sum_{\mu=1}^p \phi_i(x_\mu) \phi_j(x_\mu) \right] w_j$$

This can be written as $\eta^{-1} \frac{d}{dt} \mathbf{w} = \mathbf{u} - \mathbf{A} \mathbf{w}$, with $u_i = \sum_{\mu=1}^p t^\mu (x_\mu)^i$ and $A_{ij} = \sum_{\mu=1}^p (x_\mu)^{i+j}$, with solution (provided \mathbf{A} is invertible)

$$\mathbf{w}(t) = \mathbf{A}^{-1} \mathbf{u} + e^{-\eta t} \mathbf{A} [\mathbf{w}(0) - \mathbf{A}^{-1} \mathbf{u}]$$

The final outcome of the learning process is the weight vector $\mathbf{w}^* = \mathbf{w}(\infty) = \mathbf{A}^{-1} \mathbf{u}$. Insertion into (3.6) gives the associated function extracted by the student from the data. In figure 3.2 we show the result for the example $g(x) = \frac{1}{2} + \sin(2\pi x)$, with $0 \leq x \leq 1$. For small M (e.g. $M = 1$ in the figure), the complexity of the student is insufficient for reducing either training or generalization error. Increasing the complexity of the student (i.e. M) initially improves both (e.g. $M = 3$ in the figure). However, for large M the system increasingly succeeds in reducing the training error by getting the exact locations of the data points right (including the noise), with the side effect of a deterioration of generalization (i.e. $M = 9$ in the figure). Plotting the values of E_t and E_g as functions of time for a student with excessive complexity (M too large) would lead to the behaviour illustrated in figure 3.3.

The most commonly used solutions to the overfitting problem are the following:

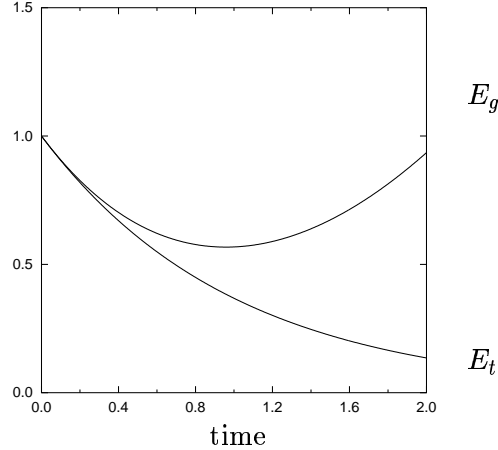


Figure 3.3: Evolution of training and generalization errors during a learning process defined as gradient descent on an error surface, with noisy data, in the regime where the student is too complex (‘overfitting’) and starts reproducing also the noise in the data.

(i) *Cross-validation.*

The idea is to split the p available data points into two sub-sets: a training set $\{(\xi^1, t^1), \dots, (\xi^k, t^k)\}$ and a validation set $\{(\xi^{k+1}, t^{k+1}), \dots, (\xi^p, t^p)\}$. One then uses the training set for learning, and the validation set to estimate the generalization error:

$$E_t[\mathbf{w}] = \frac{1}{k} \sum_{\mu=1}^k \mathcal{E}[t^\mu, f(\xi^\mu; \mathbf{w})] \quad E_g[\mathbf{w}] \approx \frac{1}{p-k} \sum_{\mu=k+1}^p \mathcal{E}[t^\mu, f(\xi^\mu; \mathbf{w})] \quad (3.7)$$

This latter estimate can be used in two possible ways:

(a) *Controlling model complexity (or: finding optimal complexity).*

In terms of the example of figure 3.2: start learning by minimizing the training error E_t with a small value of M , measure the estimate of the generalization error in (3.7), and repeat this procedure for increased values of M until E_g is seen to increase.

(b) *Controlling learning times (or: ‘early stopping’).*

In terms of figure 3.2 (see also figure 3.3): learn by minimizing the training error E_t with a large value of M (a complex student) and monitor the estimate of the generalization error in (3.7) as a function of time. Terminate learning as soon as E_g is seen to increase.

(ii) *Regularization.*

This method is based on the assumption that both the function to be learned and the basis functions of the student are in principle smooth, and that any observed irregularity and discontinuity of data must have been due to noise. Since the only way for the student to reproduce irregular or discontinuous functions is by having large (possibly diverging) components of the parameter vector \mathbf{w} , the idea behind regularization is to prevent \mathbf{w} from becoming too large by adding a suitable ‘penalty term’ to the function

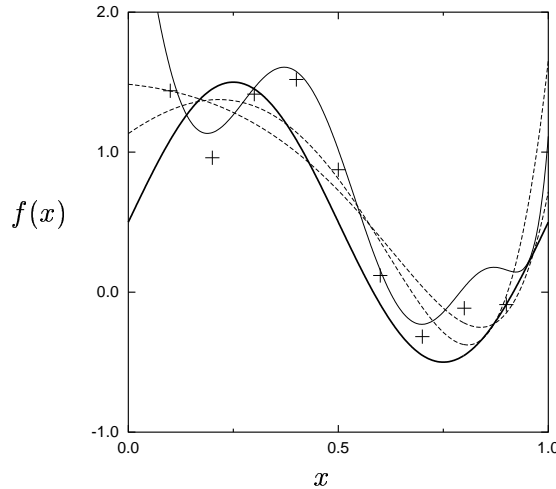


Figure 3.4: Regularization with a quadratic penalty term, as in eqn (3.8), in the example of figure 3.2. Thick solid line: the function to be learned. Crosses: nine noisy sample points of this function. Thin solid line: approximation by order-9 polynomial, without regularization. Dashed lines: similar, but with regularization term ($\lambda = 0.001$ and $\lambda = 0.02$, respectively, from bottom to top at the left of the figure). The $\lambda = 0.001$ curve improves generalization compared to the un-regularized one (i.e. it resembles the true function more). However, for $\lambda = 0.02$ we observe over-regularization: the approximation is becoming too smooth.

to be minimized, such as:

$$\frac{d}{dt} \mathbf{w} = -\eta \nabla \mathbf{w} \left\{ E_t[\mathbf{w}] + \frac{1}{2} \lambda \mathbf{w}^2 \right\} \quad E_t[\mathbf{w}] = \frac{1}{2p} \sum_{\mu=1}^p [t^\mu - f(\boldsymbol{\xi}^\mu; \mathbf{w})]^2 \quad (3.8)$$

Neither cross-validation nor regularization are ideal. In the first case we waste data and CPU time which could have been used for learning. In the second we must tune or guess the form of the penalty term (and its parameter λ); see also figure 3.4. In either case we cannot be sure that we arrive at the lowest possible generalization error. The origin of the problem is clear: we continue to minimize an object which we can easily measure (E_t) which differs from the object which we would really like to minimize (E_g), but which we cannot measure.

3.2 Bayesian Learning of Network Weights

Ideas, Definitions and Benefits. The Bayesian approach deals in a systematic and exact way with the general problem of learning in neural networks and other parametrized information processing systems of the general form $S(\boldsymbol{\xi}) = f(\boldsymbol{\xi}; \mathbf{w}) + \text{noise}$, given the observation of (possibly noisy) data $D = \{(\boldsymbol{\xi}^1, t^1), \dots, (\boldsymbol{\xi}^p, t^p)\}$. In its simplest form, i.e. for a given model, it is based on the following three ideas:

- Consider not just a single parameter vector \mathbf{w} , but an *ensemble* of parameter vectors, characterized by a probability distribution $p(\mathbf{w})$ which evolves during learning.

- Assume that the data D were generated by a system of the form $S(\boldsymbol{\xi}) = f(\boldsymbol{\xi}; \mathbf{w}) + \text{noise}$. Try to calculate the probability $p(\mathbf{w}|D)$ of the parameter vectors, given the data.
- Use the general Bayesian relation $P(A|B)P(B) = P(B|A)P(A)$ to express the desired object $p(\mathbf{w}|D)$ in terms of $p(D|\mathbf{w})$ (the latter can be calculated quite easily).

Learning is regarded as a process during which the arrival of data reduces our uncertainty about the ‘right’ parameter vector \mathbf{w} from a *prior distribution* $p(\mathbf{w})$ (reflecting prior knowledge or prejudices about the problem) to a *posterior distribution* $p(\mathbf{w}|D)$ (reflecting both prior assumptions and the data-generated evidence). Note that we can combine the general statistical relations $p(D|\mathbf{w})p(\mathbf{w}) = p(\mathbf{w}|D)P(D)$ and $p(D) = \int d\mathbf{w} p(\mathbf{w}, D)$, to obtain

$$p(\mathbf{w}|D) = \frac{p(D|\mathbf{w})p(\mathbf{w})}{\int d\mathbf{w}' p(\mathbf{w}')p(D|\mathbf{w}')} \quad (3.9)$$

This is the core identity.

We now put the above ideas in practice. In a nutshell, Bayesian learning works like this:

Stage 1: definitions

Define (i) the parametrized model, assumed responsible for the data, (ii) the prior distribution $p(\mathbf{w})$ of its parameters, and (iii) the data $D = \{(\boldsymbol{\xi}^1, t^1), \dots, (\boldsymbol{\xi}^p, t^p)\}$.

Stage 2: model translation

Convert the model definition into a probabilistic standard form: specify the likelihood of finding output t upon presentation of input $\boldsymbol{\xi}$, given the parameters \mathbf{w} :

$$\text{model definition in standard form} \quad p(t|\boldsymbol{\xi}, \mathbf{w}) \quad (3.10)$$

Stage 3: the posterior distribution

Calculate the *data* likelihood, given the model, $p(D|\mathbf{w}) = \prod_{\mu=1}^p p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w})$. From this, together with identity (3.9), follows the desired posterior parameter distribution

$$p(\mathbf{w}|D) = \frac{p(\mathbf{w}) \prod_{\mu=1}^p p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w})}{\int d\mathbf{w}' p(\mathbf{w}') \prod_{\mu=1}^p p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w}')} \quad (3.11)$$

Stage 4: prediction

The residual uncertainty in the parameters \mathbf{w} generates uncertainty in subsequent data prediction. Prediction of the output t corresponding to a new input $\boldsymbol{\xi}$, given our observation of the data D and our choice of model, thus takes the probabilistic form:

$$p(t|\boldsymbol{\xi}, D) = \int d\mathbf{w} p(t|\boldsymbol{\xi}, \mathbf{w})p(\mathbf{w}|D) \quad (3.12)$$

This concludes the description of Bayesian learning for a single model; the problem has been reduced to doing integrals. The remainder of this section deals with implementation and understanding, and with the generalization of the ideas to model selection.

The Bayesian approach to learning has become very popular within a relatively short period of time. The reason for this can be appreciated by simply listing some of its main appeals and benefits (to be derived on subsequent pages in these notes):

- It provides an interpretation of regularizers and their associated parameters.
- It provides an interpretation of single-point error measures.
- There is no need for cross-validation, so all data can be used for learning.
- It allows for the selection of the optimal complexity within a given model class, and for the selection of the optimal model class from a given set of candidates.
- It provides not only predictions, but also specific confidence estimates (i.e. error bars) corresponding to these predictions.
- Traditional learning via gradient descent training error minimization, with regularization, is recovered as a particular approximation within the Bayesian framework.

Let us gain intuition on how in the Bayesian picture the arrival of data is converted into probabilistic statements about model parameters, by working out the steps (3.10,3.11) for some simple examples. We will write inputs as $\xi \in \mathfrak{R}$ or $\boldsymbol{\xi} \in \mathfrak{R}^K$, adjustable parameters as $w \in \mathfrak{R}$ or $\mathbf{w} \in \mathfrak{R}^N$, outputs as $t \in \mathfrak{R}$, and data as $D = \{(\boldsymbol{\xi}^1, t^1), \dots, (\boldsymbol{\xi}^p, t^p)\}$.

Example 1. Our first example describes a deterministic model with one real input and one real output. Stage one in the process is the definition of model, prior parameter distribution² and data, for which we take:

$$t(\xi) = \tanh(w\xi) \quad p(w) = (2\pi)^{-\frac{1}{2}} e^{-\frac{1}{2}w^2} \quad D = \{(1, -\frac{1}{2})\}$$

(one data point). Next we convert the model into the probabilistic standard form (3.10)³:

$$p(t|\xi, w) = \delta[t - \tanh(w\xi)]$$

We can now calculate the posterior distribution (3.11):

$$\begin{aligned} p(w|D) &= \frac{p(w)p(t^1|\xi^1, w)}{\int dw' p(w')p(t^1|\xi^1, w')} = \frac{\delta[-\frac{1}{2} - \tanh(w)]e^{-\frac{1}{2}w^2}}{\int dw' \delta[-\frac{1}{2} - \tanh(w')]e^{-\frac{1}{2}w'^2}} \\ &= \delta[w - \tanh^{\text{inv}}(-\frac{1}{2})] = \delta[w + \log \sqrt{3}] \end{aligned} \quad (3.13)$$

Here we have used the two identities $\delta[f(w)] = |f'(f^{\text{inv}}(0))|^{-1} \delta[w - f^{\text{inv}}(0)]$ (see Appendix C) and $\tanh^{\text{inv}}(z) = \frac{1}{2} \log[(1+z)/(1-z)]$.

Before we had observed any data our uncertainty about the parameter w was described by the Gaussian prior $p(w)$. The arrival of the data point $(1, -\frac{1}{2})$ induced a collapse of this prior to the δ -distributed posterior $p(w|D) = \delta[w + \log \sqrt{3}]$ (without any uncertainty). This is illustrated in figure 3.5.

Example 2. Let us now make matters slightly more interesting and replace the one-parameter function of the first example by a two-parameter function:

$$t(\xi) = \tanh(w_0 + w_1\xi) \quad p(\mathbf{w}) = (2\pi)^{-1} e^{-\frac{1}{2}\mathbf{w}^2} \quad D = \{(1, -\frac{1}{2})\}$$

²Note that, given a choice made for average and variance, the Gaussian distribution is the maximum entropy (i.e. maximum uncertainty) distribution for real values random variables.

³See appendix C on distributions for the definition and properties of the distribution $\delta[z]$.

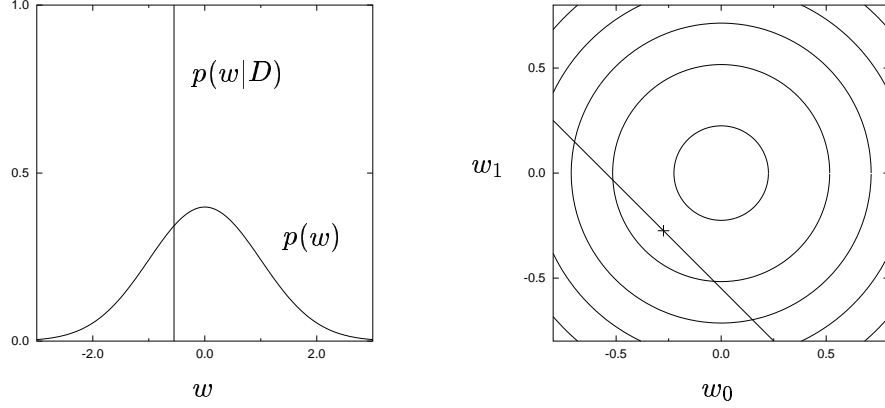


Figure 3.5: Left: reduction of parameter uncertainty due to arrival of data in example 1, from a Gaussian prior distribution $p(w)$ to a δ -peaked posterior distribution $p(w|D)$. Right: reduction of parameter uncertainty in example 2, from a Gaussian prior distribution $p(w_0, w_1)$ (contour lines) to a posterior distribution $p(w_0, w_1|D)$ with support only along the line $w_0 + w_1 = -\log \sqrt{3}$ (with its maximum at the point indicated by + in the figure).

The probabilistic standard form becomes

$$p(t|\xi, \mathbf{w}) = \delta[t - \tanh(w_0 + w_1 \xi)]$$

The posterior distribution becomes:

$$\begin{aligned} p(\mathbf{w}|D) &= \frac{p(\mathbf{w})p(t^1|\xi^1, \mathbf{w})}{\int d\mathbf{w}' p(\mathbf{w}')p(t^1|\xi^1, \mathbf{w}')} = \frac{\delta[-\frac{1}{2} - \tanh(w_0 + w_1)] e^{-\frac{1}{2}\mathbf{w}^2}}{\int d\mathbf{w}' \delta[-\frac{1}{2} - \tanh(w'_0 + w'_1)] e^{-\frac{1}{2}\mathbf{w}'^2}} \\ &= \frac{\delta[w_0 + w_1 + \tanh^{\text{inv}}(\frac{1}{2})] e^{-\frac{1}{2}\mathbf{w}^2}}{\int d\mathbf{w}' \delta[w'_0 + w'_1 + \tanh^{\text{inv}}(\frac{1}{2})] e^{-\frac{1}{2}\mathbf{w}'^2}} = \frac{\delta[w_0 + w_1 + \log \sqrt{3}] e^{-\frac{1}{2}\mathbf{w}^2}}{\int d\mathbf{w}' \delta[w'_0 + w'_1 + \log \sqrt{3}] e^{-\frac{1}{2}\mathbf{w}'^2}} \end{aligned} \quad (3.14)$$

(again using the identities $\delta[f(w)] = |f'(f^{\text{inv}}(0))|^{-1} \delta[w - f^{\text{inv}}(0)]$ and $\tanh^{\text{inv}}(z) = \frac{1}{2} \log[(1+z)/(1-z)]$). Note that the posterior (3.14) is nonzero only along the line $w_0 + w_1 = -\log \sqrt{3}$ in the parameter plane, and that along this line it is maximal for the choice $w_0 = w_1 = -\frac{1}{2} \log \sqrt{3}$. This latter result simply follows upon calculating the minimum of \mathbf{w}^2 along the line $w_0 + w_1 = -\log \sqrt{3}$ (see figure 3.5).

Before we observed any data our uncertainty about the parameters $\mathbf{w} = (w_0, w_1)$ was described by the Gaussian prior $p(\mathbf{w})$. The arrival of the data point $(1, -\frac{1}{2})$ induced a collapse of this prior to an infinitely narrow distribution along the line $w_0 + w_1 = -\log \sqrt{3}$ (a ‘slice’ of the two-dimensional Gaussian prior in the parameter plane). The most likely parameter vector is $w_0 = w_1 = -\frac{1}{2} \log \sqrt{3}$, but we note that in (3.14) we also have exact quantitative information regarding the uncertainty in (w_0, w_1) .

Note, finally, that (in contrast to the above two simple examples) we will generally have to allow our data generating models to incorporate noise, in order to retain a non-zero probability for having generated the (generally noisy) observed data.

The Link Between Bayesian Learning and Traditional Learning. At first sight it appears that Bayesian learning is quite remote from traditional gradient descent learning on an error surface, possibly with regularization, as in

$$\frac{d}{dt}\mathbf{w} = -\eta\nabla_{\mathbf{w}}\{E_t[\mathbf{w}] + \lambda E_w[\mathbf{w}]\} \quad E_t[\mathbf{w}] = \frac{1}{p}\sum_{\mu=1}^p \mathcal{E}[t^\mu - f(\boldsymbol{\xi}^\mu; \mathbf{w})] \quad (3.15)$$

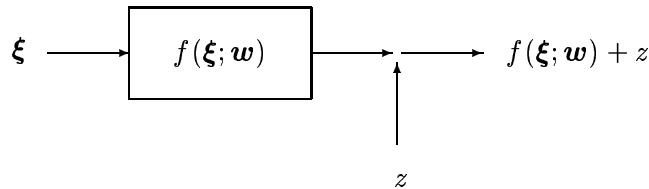
in which $\mathcal{E}[u]$ represents a single-point error measure, the function $E_w[\mathbf{w}]$ represents a regularizer, and with $\mathbf{w} \in \mathfrak{R}^N$. To reveal the link between the two viewpoints we must turn to (3.11), which we write as

$$\log p^{-1}(\mathbf{w}|D) = -\log p(\mathbf{w}) - \sum_{\mu=1}^p \log p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w}) + \log \int d\mathbf{w}' p(\mathbf{w}') \prod_{\mu=1}^p p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w}')$$

Finding the *most probable* parameter vector \mathbf{w}_{MP} , given the data D , is thus equivalent to minimizing $\log p^{-1}(\mathbf{w}|D)$, i.e. to minimizing the quantity $S(\mathbf{w}, D)$:

$$S(\mathbf{w}, D) = -\log p(\mathbf{w}) - \sum_{\mu=1}^p \log p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w}) \quad (3.16)$$

Let us work out this expression for the following simple (and natural) class of data-generating model: $t = f(\boldsymbol{\xi}; \mathbf{w}) + z$,



where f denotes some parametrized function (such as performed by the standard neural networks, e.g. MLP, RBF, etc, which are traditionally used in gradient descent learning of the type (3.15)) and where z is a zero-average random number, representing additive noise in the data. If the noise variable is distributed according to $P[z]$, one has

$$p(t|\boldsymbol{\xi}; \mathbf{w}) = P[t - f(\boldsymbol{\xi}; \mathbf{w})] \quad (3.17)$$

Now the function (3.16) (whose minimum gives the most probable parameters) reduces to

$$\frac{1}{p}S(\mathbf{w}, D) = -\frac{1}{p}\sum_{\mu=1}^p \log P[t^\mu - f(\boldsymbol{\xi}^\mu; \mathbf{w})] + \frac{1}{p}\log[1/p(\mathbf{w})]$$

Comparison with expression (3.15) reveals the following:

- Learning by finding the minimum on a training error surface, with regularizer, as in (3.15), is equivalent to finding the *most probable* parameter vector \mathbf{w}_{MP} ⁴.
- Choosing a specific single-point error measure $\mathcal{E}[u]$ in (3.15) means making the following assumption on the data noise statistics: $P[z] \sim e^{-\mathcal{E}[z]}$.

⁴This is also called the Maximum A-Posteriori Probability (MAP) procedure

- Choosing a specific regularizer $\lambda E_w[\mathbf{w}]$ in (3.15) means deciding on the following specific prior distribution for the parameters \mathbf{w} : $p(\mathbf{w}) = e^{-p\lambda E_w[\mathbf{w}]}$.

Our previously ad-hoc choices of error measure and regularizer are now replaced by direct interpretations, and hence guides for how to make appropriate choices. We also learn *en passant* that the regularizer strength must scale with the number of examples as $\lambda \sim p^{-1}$.

For example, learning by minimisation of the simplest (quadratic) training error measure $\mathcal{E}[u] = u^2$ and quadratic regularizer, viz.

$$\frac{d}{dt}\mathbf{w} = -\eta \nabla \mathbf{w} \left\{ \frac{1}{p} \sum_{\mu=1}^p [t^\mu - f(\boldsymbol{\xi}^\mu; \mathbf{w})]^2 + \frac{1}{2} \lambda \mathbf{w}^2 \right\} \quad (3.18)$$

is found to be fully equivalent to finding the most probable parameter vector \mathbf{w}_{MP} for the model

$$p(t|\boldsymbol{\xi}; \mathbf{w}) = \left[\frac{2\pi}{\beta} \right]^{-\frac{1}{2}} e^{-\frac{1}{2}\beta[t-f(\boldsymbol{\xi}; \mathbf{w})]^2} \quad p(\mathbf{w}) = \left[\frac{2\pi}{\alpha} \right]^{-\frac{N}{2}} e^{-\frac{1}{2}\alpha \mathbf{w}^2} \quad \lambda = \frac{2\alpha}{\beta p} \quad (3.19)$$

For this choice the surface (3.16) simplifies to (modulo an irrelevant constant):

$$S(\mathbf{w}, D) = \frac{1}{2}\beta \sum_{\mu=1}^p [t^\mu - f(\boldsymbol{\xi}^\mu; \mathbf{w})]^2 + \frac{1}{2}\alpha \mathbf{w}^2 \quad (3.20)$$

We need no longer *guess* the value of λ , but can calculate it via $\lambda = 2\alpha/\beta p$ from

$$\alpha = \frac{1}{\langle w_i^2 \rangle_{\text{prior}}} \quad \beta = \frac{1}{\langle t^2 \rangle_{\text{noise}} - \langle t \rangle_{\text{noise}}^2} \quad p = \text{number of data points}$$

Quantities such as α or β in (3.19), which are not themselves adjustable parameters in the sense of \mathbf{w} , but are more global parameters which reflect prior knowledge of the problem and which will influence the leaning process, are called *hyper-parameters*.

Approximation of the Posterior Parameter Distribution. If we were to only calculate the most probable parameter vector \mathbf{w}_{MP} we would gain only interpretations compared to old-fashioned training error minimization. The power of the Bayesian techniques is that they also provide information on the reliability of a learning outcome. This information is embodied in the full posterior distribution $p(\mathbf{w}|D)$, which can be written in terms of (3.16) as

$$p(\mathbf{w}|D) = \frac{e^{-S(\mathbf{w}, D)}}{\int d\mathbf{w}' e^{-S(\mathbf{w}', D)}} \quad (3.21)$$

In those cases where $S(\mathbf{w}, D)$ has only a single relevant local (and thus also global) minimum \mathbf{w}_{MP} , the reliability of the learning outcome is mainly coded in the local curvature (i.e. ‘width’) of $S(\mathbf{w}, D)$ around \mathbf{w}_{MP} . We now expand $S(\mathbf{w}, D)$ in the vicinity of \mathbf{w}_{MP} :

$$S(\mathbf{w}, D) = S(\mathbf{w}_{\text{MP}}, D) + \frac{1}{2}(\mathbf{w} - \mathbf{w}_{\text{MP}}) \cdot \mathbf{A}(\mathbf{w} - \mathbf{w}_{\text{MP}}) + \mathcal{O}(|\mathbf{w} - \mathbf{w}_{\text{MP}}|^3) \quad (3.22)$$

$$A_{ij} = \left. \frac{\partial^2 S}{\partial w_i \partial w_j} \right|_{\mathbf{w}_{\text{MP}}} \quad (3.23)$$

(\mathbf{A} is called the Hessian matrix of S at the point \mathbf{w}_{MP} ; linear terms are absent in (3.22) due to \mathbf{w}_{MP} being a minimum).

Truncation of the expansion (3.22) after the quadratic term leads to a Gaussian approximation of the posterior $p(\mathbf{w}|D)$:

$$S(\mathbf{w}, D) \rightarrow \tilde{S}(\mathbf{w}, D) = S(\mathbf{w}_{\text{MP}}, D) + \frac{1}{2}(\mathbf{w} - \mathbf{w}_{\text{MP}}) \cdot \mathbf{A}(\mathbf{w} - \mathbf{w}_{\text{MP}}) \quad (3.24)$$

$$p(\mathbf{w}|D) \rightarrow \tilde{p}(\mathbf{w}|D) = \left[\frac{\det \mathbf{A}}{(2\pi)^N} \right]^{\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{w} - \mathbf{w}_{\text{MP}}) \cdot \mathbf{A}(\mathbf{w} - \mathbf{w}_{\text{MP}})} \quad (3.25)$$

Average, variance and co-variance matrix of the approximated posterior distribution $\tilde{p}(\mathbf{w}|D)$ are given by the following identities (see appendix B) which involve the inverse of the Hessian (3.23), with the notation $\langle f(\mathbf{w}) \rangle = \int d\mathbf{w} f(\mathbf{w}) \tilde{p}(\mathbf{w}|D)$:

$$\langle \mathbf{w} \rangle = \mathbf{w}_{\text{MP}} \quad \langle w_i^2 \rangle - \langle w_i \rangle^2 = (\mathbf{A}^{-1})_{ii} \quad \langle w_i w_j \rangle - \langle w_i \rangle \langle w_j \rangle = (\mathbf{A}^{-1})_{ij} \quad (3.26)$$

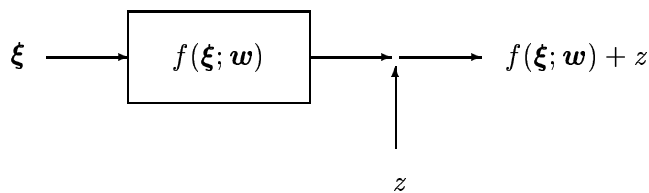
3.3 Predictions with Error Bars: Real-Valued Functions

Learning and predicting the action of continuous real-valued functions from (noisy) data examples is called *regression*. As we have seen in the previous section, a trained network is within the Bayesian framework given by the posterior distribution $p(\mathbf{w}|D)$ for the system parameters, as given by (3.11). Prediction then proceeds via formula (3.12), which can also be written in terms of the function $S(\mathbf{w}, D)$ (3.16) as

$$p(t|\xi, D) = \int d\mathbf{w} p(t|\xi, \mathbf{w}) p(\mathbf{w}|D) \quad p(\mathbf{w}|D) = \frac{e^{-S(\mathbf{w}, D)}}{\int d\mathbf{w}' e^{-S(\mathbf{w}', D)}} \quad (3.27)$$

This formally defines the statistics of the output to be associated with input ξ .

Mean and Variance for Gaussian Output Noise and Gaussian Priors. We now work out output average and variance of (3.27) for the simplest class of systems (3.19): a parametrized function $f(\xi; \mathbf{w})$ with additive zero-average Gaussian data noise, and with a Gaussian prior parameter distribution $p(\mathbf{w})$:



The moments $\langle t \rangle = \int dt t p(t|\xi, D)$ and $\langle t^2 \rangle = \int dt t^2 p(t|\xi, D)$ become (transform integration variables by putting $t = f(\xi; \mathbf{w}) + z/\sqrt{\beta}$):

$$\langle t \rangle = \left[\frac{2\pi}{\beta} \right]^{\frac{1}{2}} \int dt t \int d\mathbf{w} e^{-\frac{1}{2}\beta[t - f(\xi; \mathbf{w})]^2} p(\mathbf{w}|D) \quad (3.28)$$

$$\begin{aligned} &= \int d\mathbf{w} \int \frac{dz}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} [f(\xi; \mathbf{w}) + z/\sqrt{\beta}] p(\mathbf{w}|D) \\ &= \int d\mathbf{w} f(\xi; \mathbf{w}) p(\mathbf{w}|D) \end{aligned} \quad (3.29)$$

$$\begin{aligned}
\langle t^2 \rangle &= \left[\frac{2\pi}{\beta} \right]^{\frac{1}{2}} \int dt t^2 \int d\mathbf{w} e^{-\frac{1}{2}\beta[t-f(\boldsymbol{\xi};\mathbf{w})]^2} p(\mathbf{w}|D) \\
&= \int d\mathbf{w} \int \frac{dz}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} [f(\boldsymbol{\xi};\mathbf{w}) + z/\sqrt{\beta}]^2 p(\mathbf{w}|D) \\
&= \frac{1}{\beta} + \int d\mathbf{w} f^2(\boldsymbol{\xi};\mathbf{w}) p(\mathbf{w}|D)
\end{aligned} \tag{3.30}$$

The prediction variance $\sigma^2 = \langle t^2 \rangle - \langle t \rangle^2$ now follows as

$$\sigma^2 = \frac{1}{\beta} + \underbrace{\int d\mathbf{w} f^2(\boldsymbol{\xi};\mathbf{w}) p(\mathbf{w}|D) - \left[\int d\mathbf{w} f(\boldsymbol{\xi};\mathbf{w}) p(\mathbf{w}|D) \right]^2}_{\text{variance due to uncertainty about system parameters}} \tag{3.31}$$

The first term, β^{-1} , represents the intrinsic uncertainty in the data generation model; the remainder in (3.31) reflects our uncertainty about the right parameters after having learned the data. The Bayesian predicted output $t^*(\boldsymbol{\xi})$ for the input $\boldsymbol{\xi}$, and its associated error margin $\Delta t^*(\boldsymbol{\xi})$, are now defined as $\langle t \rangle$ and σ , respectively:⁵

$$t^*(\boldsymbol{\xi}) = \int d\mathbf{w} f(\boldsymbol{\xi};\mathbf{w}) p(\mathbf{w}|D) \tag{3.32}$$

$$\Delta t^*(\boldsymbol{\xi}) = \sqrt{\frac{1}{\beta} + \int d\mathbf{w} f^2(\boldsymbol{\xi};\mathbf{w}) p(\mathbf{w}|D) - \left[\int d\mathbf{w} f(\boldsymbol{\xi};\mathbf{w}) p(\mathbf{w}|D) \right]^2} \tag{3.33}$$

For Gaussian output noise and priors the function $S(\mathbf{w}, D)$ is given by (3.20), hence

$$p(\mathbf{w}|D) = \frac{e^{-S(\mathbf{w}, D)}}{\int d\mathbf{w}' e^{-S(\mathbf{w}', D)}} \quad S(\mathbf{w}, D) = \frac{1}{2}\beta \sum_{\mu=1}^p [t^\mu - f(\boldsymbol{\xi}^\mu; \mathbf{w})]^2 + \frac{1}{2}\alpha \mathbf{w}^2$$

In (3.32) we now have proper predictions *with error bars*.

Simplification for Approximated Posterior Distribution. The above exact result (3.32,3.33) can be simplified using the approximated posterior distribution $\tilde{p}(\mathbf{w}|D)$ of (3.25). Upon also putting $\mathbf{w} = \mathbf{w}_{\text{MP}} + \mathbf{A}^{-\frac{1}{2}}\mathbf{z}$ in the parameter integrations, we then find

$$t^*(\boldsymbol{\xi}) = \int dz \frac{e^{-\frac{1}{2}\mathbf{z}^2}}{(2\pi)^{N/2}} f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}} + \mathbf{A}^{-\frac{1}{2}}\mathbf{z}) \tag{3.34}$$

$$\Delta t^*(\boldsymbol{\xi}) = \sqrt{\frac{1}{\beta} + \int dz \frac{e^{-\frac{1}{2}\mathbf{z}^2}}{(2\pi)^{N/2}} f^2(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}} + \mathbf{A}^{-\frac{1}{2}}\mathbf{z}) - [t^*(\boldsymbol{\xi})]^2} \tag{3.35}$$

We now make one further and final approximation. If the width of the Gaussian approximation (3.25) is not too large, we may approximate/expand $f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}} + \mathbf{u})$ for small \mathbf{u} :

$$f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}} + \mathbf{A}^{-\frac{1}{2}}\mathbf{z}) = f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}}) + \mathbf{x}(\boldsymbol{\xi}) \cdot \mathbf{A}^{-\frac{1}{2}}\mathbf{z} + \frac{1}{2}(\mathbf{A}^{-\frac{1}{2}}\mathbf{z}) \cdot \mathbf{B}(\boldsymbol{\xi})(\mathbf{A}^{-\frac{1}{2}}\mathbf{z}) + \mathcal{O}([\mathbf{A}^{-\frac{1}{2}}\mathbf{z}]^3) \tag{3.36}$$

⁵Note: $t^*(\boldsymbol{\xi})$ need not be identical to $f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}})$, where \mathbf{w}_{MP} denotes the most probable parameters.

where $x_i(\boldsymbol{\xi}) = \partial f(\boldsymbol{\xi}; \mathbf{w}) / \partial w_i |_{\mathbf{w}_{\text{MP}}}$ and $B_{ij}(\boldsymbol{\xi}) = \partial^2 f(\boldsymbol{\xi}; \mathbf{w}) / \partial w_i \partial w_j |_{\mathbf{w}_{\text{MP}}}$. Insertion into (3.34,3.35) gives:

$$\begin{aligned}
t^*(\boldsymbol{\xi}) &= f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}}) + \frac{1}{2} \sum_{ij} \int d\mathbf{z} \frac{e^{-\frac{1}{2}\mathbf{z}^2}}{(2\pi)^{N/2}} z_i [\mathbf{A}^{-\frac{1}{2}} \mathbf{B}(\boldsymbol{\xi}) \mathbf{A}^{-\frac{1}{2}}]_{ij} z_j + \mathcal{O}(|\mathbf{A}^{-2}|) \\
&= f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}}) + \frac{1}{2} \text{Tr}[\mathbf{A}^{-\frac{1}{2}} \mathbf{B}(\boldsymbol{\xi}) \mathbf{A}^{-\frac{1}{2}}] + \mathcal{O}(|\mathbf{A}^{-2}|) \\
[\Delta t^*(\boldsymbol{\xi})]^2 &= \frac{1}{\beta} + f^2(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}}) - [t^*(\boldsymbol{\xi})]^2 + \mathcal{O}(|\mathbf{A}^{-2}|) \\
&\quad + \int d\mathbf{z} \frac{e^{-\frac{1}{2}\mathbf{z}^2}}{(2\pi)^{N/2}} \left[[\mathbf{x}(\boldsymbol{\xi}) \cdot \mathbf{A}^{-\frac{1}{2}} \mathbf{z}]^2 + f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}}) (\mathbf{A}^{-\frac{1}{2}} \mathbf{z}) \cdot \mathbf{B}(\boldsymbol{\xi}) (\mathbf{A}^{-\frac{1}{2}} \mathbf{z}) \right] \\
&= \frac{1}{\beta} - f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}}) \cdot \text{Tr}[\mathbf{A}^{-\frac{1}{2}} \mathbf{B}(\boldsymbol{\xi}) \mathbf{A}^{-\frac{1}{2}}] + \mathcal{O}(|\mathbf{A}^{-2}|) \\
&\quad + \sum_{ij} \left\{ [\mathbf{A}^{-\frac{1}{2}} \mathbf{x}(\boldsymbol{\xi})]_i [\mathbf{A}^{-\frac{1}{2}} \mathbf{x}(\boldsymbol{\xi})]_j + f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}}) [\mathbf{A}^{-\frac{1}{2}} \mathbf{B}(\boldsymbol{\xi}) \mathbf{A}^{-\frac{1}{2}}]_{ij} \right\} \delta_{ij} \\
&= \frac{1}{\beta} - f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}}) \cdot \text{Tr}[\mathbf{A}^{-\frac{1}{2}} \mathbf{B}(\boldsymbol{\xi}) \mathbf{A}^{-\frac{1}{2}}] \\
&\quad + \mathbf{x}(\boldsymbol{\xi}) \cdot \mathbf{A}^{-1} \mathbf{x}(\boldsymbol{\xi}) + f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}}) \text{Tr}[\mathbf{A}^{-\frac{1}{2}} \mathbf{B}(\boldsymbol{\xi}) \mathbf{A}^{-\frac{1}{2}}] + \mathcal{O}(|\mathbf{A}^{-2}|) \\
&= \frac{1}{\beta} + \mathbf{x}(\boldsymbol{\xi}) \cdot \mathbf{A}^{-1} \mathbf{x}(\boldsymbol{\xi}) + \mathcal{O}(|\mathbf{A}^{-2}|)
\end{aligned}$$

(with the trace of a matrix \mathbf{C} defined as usual: $\text{Tr} \mathbf{C} = \sum_i C_{ii}$). Hence we arrive at:

$$t^*(\boldsymbol{\xi}) = f(\boldsymbol{\xi}; \mathbf{w}_{\text{MP}}) + \frac{1}{2} \text{Tr}[\mathbf{A}^{-\frac{1}{2}} \mathbf{B}(\boldsymbol{\xi}) \mathbf{A}^{-\frac{1}{2}}] + \mathcal{O}(|\mathbf{A}^{-2}|) \quad (3.37)$$

$$\Delta t^*(\boldsymbol{\xi}) = \sqrt{\beta^{-1} + \mathbf{x}(\boldsymbol{\xi}) \cdot \mathbf{A}^{-1} \mathbf{x}(\boldsymbol{\xi})} + \mathcal{O}(|\mathbf{A}^{-2}|) \quad (3.38)$$

For expressions (3.37,3.38) to be accurate⁶ and useful (given our earlier assumptions of Gaussian additive data noise and a Gaussian prior) we need $p(\mathbf{w}|D)$ to be (i) approximately Gaussian, and (ii) sufficiently narrow (in view of (3.26), since this would guarantee that \mathbf{A}^{-1} can indeed be treated as small).

Models for which Truncation is Exact. In order to appreciate the nature of the approximations one would arrive at by simply neglecting the $\mathcal{O}(|\mathbf{A}^{-2}|)$ terms in (3.37,3.38), it is instructive to determine the conditions under which they would have been exact. Given Gaussian data output noise and a Gaussian prior, approximation (3.25) would have been exact only if $S(\mathbf{w}, D)$ (3.16) would depend quadratically on \mathbf{w} , i.e. if $f(\boldsymbol{\xi}; \mathbf{w})$ would depend linearly on \mathbf{w} :

$$p(t|\boldsymbol{\xi}; \mathbf{w}) = \left[\frac{2\pi}{\beta} \right]^{\frac{1}{2}} e^{-\frac{1}{2}\beta[t-f(\boldsymbol{\xi}; \mathbf{w})]^2} \quad p(\mathbf{w}) = \left[\frac{2\pi}{\alpha} \right]^{\frac{N}{2}} e^{-\frac{1}{2}\alpha \mathbf{w}^2} \quad f(\boldsymbol{\xi}; \mathbf{w}) = \sum_i w_i \phi_i(\boldsymbol{\xi}) \quad (3.39)$$

(for some set of functions $\{\phi_i(\boldsymbol{\xi})\}$, as in e.g. RBF networks). Here also (3.36) is exact when truncated after the term linear in \mathbf{z} , so $\mathbf{B}(\boldsymbol{\xi}) = 0$ and the $\mathcal{O}(|\mathbf{A}^{-2}|)$ terms in (3.37,3.38) are

⁶Note: in many textbooks one finds that the second term in $t^*(\boldsymbol{\xi})$ is also neglected.

absent. Furthermore, since $S(\mathbf{w}, D)$ is now truly Gaussian, we can calculate $\mathbf{x}(\boldsymbol{\xi})$ and \mathbf{w}_{MP} explicitly. We define $\boldsymbol{\phi}(\boldsymbol{\xi}) = (\phi_1(\boldsymbol{\xi}), \dots, \phi_N(\boldsymbol{\xi}))$, and get

$$\begin{aligned} S(\mathbf{w}, D) &= \frac{1}{2}\beta \sum_{\mu=1}^p [t^\mu - \mathbf{w} \cdot \boldsymbol{\phi}(\boldsymbol{\xi}^\mu)]^2 + \frac{1}{2}\alpha \mathbf{w}^2 \\ &= \frac{1}{2}\beta \sum_{\mu=1}^p [t^\mu]^2 - \beta \mathbf{w} \cdot \sum_{\mu=1}^p t^\mu \boldsymbol{\phi}(\boldsymbol{\xi}^\mu) + \frac{1}{2}\mathbf{w} \cdot \mathbf{A} \mathbf{w} \end{aligned} \quad (3.40)$$

$$A_{ij} = \alpha \delta_{ij} + \beta \sum_{\mu=1}^p \phi_i(\boldsymbol{\xi}^\mu) \phi_j(\boldsymbol{\xi}^\mu) \quad (3.41)$$

Since the Hessian matrix \mathbf{A} is positive definite (provided $\alpha > 0$), the surface $S(\mathbf{w}, D)$ is convex and has a unique minimum, which is calculated simply from putting $\partial S(\mathbf{w}, D)/\partial w_i = 0$ for all i . This reveals $\mathbf{A} \mathbf{w}_{\text{MP}} = \beta \sum_{\mu=1}^p t^\mu \boldsymbol{\phi}(\boldsymbol{\xi}^\mu)$, or

$$\mathbf{w}_{\text{MP}} = \beta \mathbf{A}^{-1} \sum_{\mu=1}^p t^\mu \boldsymbol{\phi}(\boldsymbol{\xi}^\mu) \quad (3.42)$$

Similarly one finds $x_i(\boldsymbol{\xi}) = \phi_i(\boldsymbol{\xi})$. Insertion into (3.37,3.38) then gives the, for the present models (3.39) fully exact, result

$$t^*(\boldsymbol{\xi}) = \beta \boldsymbol{\phi}(\boldsymbol{\xi}) \cdot \mathbf{A}^{-1} \sum_{\mu=1}^p t^\mu \boldsymbol{\phi}(\boldsymbol{\xi}^\mu) \quad \Delta t^*(\boldsymbol{\xi}) = \sqrt{\beta^{-1} + \boldsymbol{\phi}(\boldsymbol{\xi}) \cdot \mathbf{A}^{-1} \boldsymbol{\phi}(\boldsymbol{\xi})} \quad (3.43)$$

Estimation of Hyper-Parameters by Bootstrapping. The predictions and error margins in (3.37,3.38) and (3.43) obviously depend on the choices made for the hyper-parameters α and β . We will discuss a proper procedure for this later. A simpler (and generally quicker) method for estimating α and β is the following. Provided the number p of data is not too small, it makes sense to require that the trained system will *on average* (i) predict correctly the outputs for the data points used in training, and (ii) make errors on these training data whose magnitude is of the order of the uncertainty it associates with its own predictions:

$$\frac{1}{p} \sum_{\mu=1}^p [t^\mu - t^*(\boldsymbol{\xi}^\mu)] = 0 \quad \frac{1}{p} \sum_{\mu=1}^p \left\{ [t^\mu - t^*(\boldsymbol{\xi}^\mu)]^2 - [\Delta t^*(\boldsymbol{\xi}^\mu)]^2 \right\} = 0 \quad (3.44)$$

These two equations can be used to determine α and β .

Example 1. Let us illustrate the above procedures using the example problem of (3.6), which involves a parametrized function of the form (3.39), with a quadratic regularizer term $\frac{1}{2}\lambda \mathbf{w}^2$. To make the learning problem slightly more difficult we have now also distributed the data inhomogeneously along the x -axis. The non-Bayesian solution to this simple problem (i.e. minimize the training error plus generalizer), is easily calculated and found to be

$$t^*(x) = \sum_{ij=0}^M \phi_i(x) (\mathbf{C}^{-1})_{ij} \sum_{\mu} t^\mu \phi_j(x) \quad C_{ij} = \lambda \delta_{ij} + \sum_{\mu=1}^p \phi_i(x^\mu) \phi_j(x^\mu) \quad \phi_\ell(x) = x^\ell \quad (3.45)$$

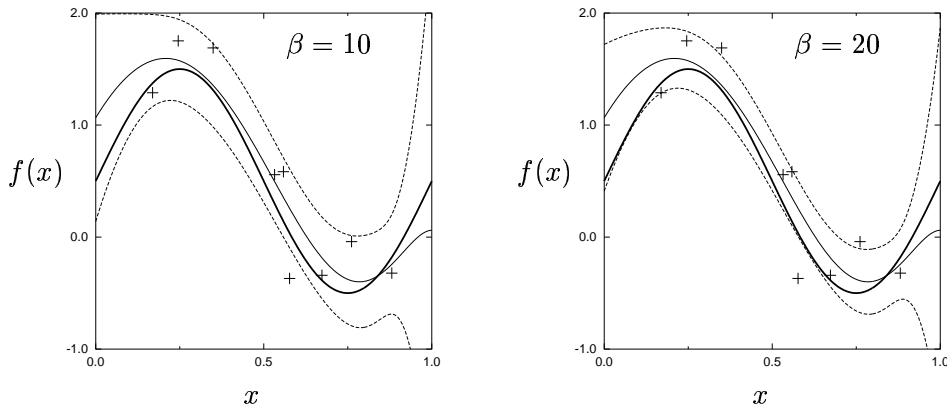


Figure 3.6: Bayesian prediction on the basis of an order-9 polynomial with assumed Gaussian data noise (variance: β^{-1}) and a Gaussian prior (variance: $\alpha^{-1} = (\lambda\beta)^{-1}$, $\lambda = 0.005$). Thick solid line: the actual function $f(x) = \frac{1}{2} + \sin(2\pi x)$ to be learned. Crosses: nine noisy sample points of this function. Solid line: predicted function $t^*(x)$ after learning. Dashed lines: $t^*(x) \pm \Delta t^*(x)$, where $\Delta t^*(x)$ denotes the Bayesian prediction uncertainty. Note: the actual data variance corresponds to $\beta = 12$.

(without error bars). The Bayesian answer (3.43), upon assuming Gaussian output noise and a Gaussian prior, would be exactly the same prediction $t^*(x)$ given in (3.45)⁷, with $\lambda = \alpha/\beta$, but it would also equip this answer with the following error estimate:

$$\Delta t^*(x) = \beta^{-\frac{1}{2}} \sqrt{1 + \sum_{ij=0}^M x^{i+j} (\mathbf{C}^{-1})_{ij}} \quad (3.46)$$

The result is shown in figure 3.6 for $\lambda = \alpha/\beta = 0.005$, $M = 9$, and $\beta \in \{10, 20\}$. The actual value of the variance in the data points (unknown to the network) in this example corresponds to $\beta = 12$. As expected, we see again in the above graphs that it will be very important to have a tool with which to calculate the hyper-parameters α and β .

Example 2. Our second example is also a simple linear system as in (3.39), but now the input vectors are two-dimensional, and we will focus on finding the hyper-parameters via ‘bootstrapping’:

$$p(t|\boldsymbol{\xi}; \mathbf{w}) = \left[\frac{2\pi}{\beta}\right]^{-\frac{1}{2}} e^{-\frac{1}{2}\beta[t-f(\boldsymbol{\xi}; \mathbf{w})]^2} \quad p(\mathbf{w}) = \left[\frac{2\pi}{\alpha}\right]^{-\frac{N}{2}} e^{-\frac{1}{2}\alpha\mathbf{w}^2} \quad f(\boldsymbol{\xi}; \mathbf{w}) = \mathbf{w} \cdot \boldsymbol{\xi} \quad (3.47)$$

with $\boldsymbol{\xi}, \mathbf{w} \in \mathbb{R}^2$. In the notation of (3.39) we here have $\phi_i(\boldsymbol{\xi}) = \xi_i$. The data points to be learned from, four in number, are taken to be the following:

$$(\boldsymbol{\xi}^1, t^1) = \left(\left(1, 0\right), \frac{3}{4}\right), \quad (\boldsymbol{\xi}^2, t^2) = \left(\left(0, 1\right), \frac{1}{4}\right), \quad (\boldsymbol{\xi}^3, t^3) = \left(\left(-1, 0\right), -\frac{1}{4}\right), \quad (\boldsymbol{\xi}^4, t^4) = \left(\left(0, -1\right), -\frac{3}{4}\right)$$

For this choice of data the matrix \mathbf{A} (3.41) becomes

$$\mathbf{A} = \alpha \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} + \beta \left\{ \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} + \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \right\} = (\alpha + 2\beta) \mathbf{I}$$

⁷This is, of course, a direct consequence of our choice of a simple linear example.

Its inversion is trivial, and hence the predictions and error bars (3.43) can be worked out explicitly. Using $\mathbf{A}^{-1} = (\alpha + 2\beta)^{-1} \mathbf{1}$ and $\sum_{\mu=1}^4 t^\mu \boldsymbol{\xi}^\mu = (1, 1)$ one arrives at

$$t^*(\boldsymbol{\xi}) = \frac{\beta(\xi_1 + \xi_2)}{\alpha + 2\beta} \quad \Delta t^*(\boldsymbol{\xi}) = \beta^{-\frac{1}{2}} \sqrt{1 + \frac{\beta(\xi_1^2 + \xi_2^2)}{\alpha + 2\beta}} \quad (3.48)$$

In order to work out the ‘bootstrapping’ relations (3.44) for the present example we next calculate the system’s internal output predictions for the four data points:

$$t^*(\boldsymbol{\xi}^1) = t^*(\boldsymbol{\xi}^2) = \frac{\beta}{\alpha + 2\beta} \quad t^*(\boldsymbol{\xi}^3) = t^*(\boldsymbol{\xi}^4) = \frac{-\beta}{\alpha + 2\beta}$$

$$\Delta t^*(\boldsymbol{\xi}^1) = \Delta t^*(\boldsymbol{\xi}^2) = \Delta t^*(\boldsymbol{\xi}^3) = \Delta t^*(\boldsymbol{\xi}^4) = \beta^{-\frac{1}{2}} \sqrt{\frac{\alpha + 3\beta}{\alpha + 2\beta}}$$

Since also $\sum_{\mu=1}^4 t^\mu = 0$, the first condition of (3.44) now reduces to the trivial identity $0 = 0$. The second condition of (3.44), however, gives us

$$\frac{\alpha + 3\beta}{\beta(\alpha + 2\beta)} = \frac{1}{4} \left\{ \left[\frac{3}{4} - \frac{\beta}{\alpha + 2\beta} \right]^2 + \left[\frac{1}{4} - \frac{\beta}{\alpha + 2\beta} \right]^2 + \left[-\frac{1}{4} + \frac{\beta}{\alpha + 2\beta} \right]^2 + \left[-\frac{3}{4} + \frac{\beta}{\alpha + 2\beta} \right]^2 \right\}$$

$$\frac{\alpha + 3\beta}{\beta(\alpha + 2\beta)} = \frac{1}{32} \left\{ \left[\frac{3\alpha + 2\beta}{\alpha + 2\beta} \right]^2 + \left[\frac{\alpha - 2\beta}{\alpha + 2\beta} \right]^2 \right\}$$

Working out this expression upon putting $\alpha = \lambda\beta$ leads to

$$\beta = 24 \left[\frac{1 + \frac{5}{6}\lambda + \frac{1}{6}\lambda^2}{1 + \lambda + \frac{5}{4}\lambda^2} \right] \quad (3.49)$$

We thus have only one free hyper-parameter left. Upon putting $\alpha \rightarrow 0$ (since $p = 4$ and we have just two adjustable parameters there should in principle be little need for regularization via a prior) gives, via (3.49), the prescription $\beta = 24$, and hence the parameter-free predictions

$$t^*(\boldsymbol{\xi}) = \frac{1}{2}(\xi_1 + \xi_2) \quad \Delta t^*(\boldsymbol{\xi}) = \frac{1}{4} \sqrt{\frac{1}{3}(2 + \xi_1^2 + \xi_2^2)} \quad (3.50)$$

This result is indeed perfectly consistent with the statistics of the training data:

	actual data	$t^*(\boldsymbol{\xi}) \pm \Delta t^*(\boldsymbol{\xi})$
t^1 :	3/4	$1/2 \pm 1/4$
t^2 :	1/4	$1/2 \pm 1/4$
t^3 :	-1/4	$-1/2 \pm 1/4$
t^4 :	-3/4	$-1/2 \pm 1/4$

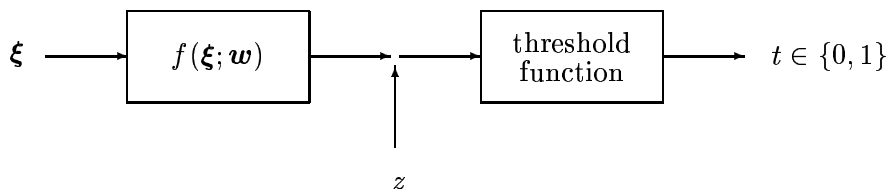
The model is also seen to take the sensible decision to assign increasing error bars when predictions are requested on input vectors which are further away from the origin (around which the data are located).

3.4 Predictions with Error Bars: Binary Classification

Predicting binary classifications on the basis of (noisy) data examples proceeds in a way similar to regression. A trained network is in the Bayesian picture given by the posterior distribution $p(\mathbf{w}|D)$ (3.11) for the system parameters, and prediction is again based on (3.12):

$$p(t|\xi, D) = \int d\mathbf{w} p(t|\xi, \mathbf{w})p(\mathbf{w}|D) \quad p(\mathbf{w}|D) = \frac{p(\mathbf{w}) \prod_{\mu=1}^p p(t^\mu|\xi^\mu, \mathbf{w})}{\int d\mathbf{w}' p(\mathbf{w}') \prod_{\mu=1}^p p(t^\mu|\xi^\mu, \mathbf{w}')} \quad (3.51)$$

The only difference with regression is that, in view of the requirement to extract binary outputs from continuous parametrized functions, the output noise can never be Gaussian.



Without output noise we would define the deterministic operation $t(\xi) = \theta[f(\xi; \mathbf{w})]$, with the step function (i.e. $\theta[z > 0] = 1$, $\theta[z < 0] = 0$). Adding noise implies allowing for $t = 0$ even when $f(\xi; \mathbf{w}) > 0$ and for $t = 1$ even when $f(\xi; \mathbf{w}) < 0$. So $p(1|\xi, \mathbf{w})$ must be a monotonically increasing function of $f(\xi; \mathbf{w})$, such that $p(1|\xi, \mathbf{w}) \rightarrow 0$ when $f(\xi; \mathbf{w}) \rightarrow -\infty$ and $p(1|\xi, \mathbf{w}) \rightarrow 1$ when $f(\xi; \mathbf{w}) \rightarrow \infty$ ⁸. A common choice is:

$$\begin{aligned} p(1|\xi, \mathbf{w}) &= \frac{1}{2} + \frac{1}{2} \tanh[f(\xi; \mathbf{w})] \\ p(0|\xi, \mathbf{w}) &= \frac{1}{2} - \frac{1}{2} \tanh[f(\xi; \mathbf{w})] \end{aligned} \quad (3.52)$$

Decision Boundary and Measure of Uncertainty. Our objective is to classify new inputs ξ into one of two categories. Hence learning from data is here equivalent to deciding on the decision boundary in input space, and to quantify the uncertainty of this decision. If we simply combine (3.51,3.52) we arrive at

$$p(1|\xi, D) = \frac{1}{2} + \frac{1}{2} I(\xi, D) \quad p(0|\xi, D) = \frac{1}{2} - \frac{1}{2} I(\xi, D) \quad (3.53)$$

$$I(\xi, D) = \frac{\int d\mathbf{w} \tanh[f(\xi; \mathbf{w})] p(\mathbf{w}) \prod_{\mu=1}^p p(t^\mu|\xi^\mu, \mathbf{w})}{\int d\mathbf{w} p(\mathbf{w}) \prod_{\mu=1}^p p(t^\mu|\xi^\mu, \mathbf{w})} \quad (3.54)$$

We classify $t^*(\xi) = 1$ if $p(1|\xi, D) > p(0|\xi, D)$, and $t^*(\xi) = 0$ if $p(1|\xi, D) < p(0|\xi, D)$. Hence the decision boundary in ξ -space is defined by $I(\xi, D) = 0$. Our uncertainty is measured by

$$\Delta t^*(\xi) = \text{Prob}[\text{incorrect classification}] = \begin{cases} p(1|\xi, D) & \text{when } t^*(\xi) = 0 \\ p(0|\xi, D) & \text{when } t^*(\xi) = 1 \end{cases}$$

⁸An implicit requirement is that the selected parametrized function $f(\xi; \mathbf{w})$ can actually approach $\pm\infty$ by boosting suitably the parameters \mathbf{w} . This, however, can always be achieved. For example, if we initially have a bounded function $g(\xi; \mathbf{w})$, we can add an extra adjustable parameter w_0 and define $f(\mathbf{x}; \mathbf{w}) = w_0 g(\xi; \mathbf{w})$

These statements can all be compactified, using the definition (3.54), into

$$\begin{aligned} I(\boldsymbol{\xi}, D) > 0 : t^*(\boldsymbol{\xi}) &= 1 & \Delta t^*(\boldsymbol{\xi}) &= \frac{1}{2} - \frac{1}{2}|I(\boldsymbol{\xi}, D)| \\ I(\boldsymbol{\xi}, D) < 0 : t^*(\boldsymbol{\xi}) &= 0 \end{aligned} \quad (3.55)$$

At the decision boundary $I(\boldsymbol{\xi}, D) = 0$ we have $\Delta t^*(\boldsymbol{\xi}) = \frac{1}{2}$, i.e. a mis-classification probability equivalent to random guessing, as it should.

Example. Let us consider a simple task involving the binary classification of two-dimensional vectors, with a linear parametrized function and a Gaussian prior:

$$\boldsymbol{\xi}, \mathbf{w} \in \mathbb{R}^2 \quad f(\boldsymbol{\xi}; \mathbf{w}) = \mathbf{w} \cdot \boldsymbol{\xi} \quad p(\mathbf{w}) = \frac{\alpha}{2\pi} e^{-\frac{1}{2}\alpha \mathbf{w}^2} \quad (3.56)$$

and with the following two data points

$$D = \{((0, -1), 1), ((0, 1), 0)\} \quad (3.57)$$

Working out the key function $I(\boldsymbol{\xi}, D)$ of (3.54) for this example gives

$$\begin{aligned} I(\boldsymbol{\xi}, D) &= \frac{\int d\mathbf{w} \tanh[\mathbf{w} \cdot \boldsymbol{\xi}] e^{-\frac{1}{2}\alpha \mathbf{w}^2} \prod_{\mu=1}^2 p(t^\mu | \boldsymbol{\xi}^\mu, \mathbf{w})}{\int d\mathbf{w} e^{-\frac{1}{2}\alpha \mathbf{w}^2} \prod_{\mu=1}^2 p(t^\mu | \boldsymbol{\xi}^\mu, \mathbf{w})} \\ &= \frac{\int d\mathbf{w} \tanh[w_1 \xi_1 + w_2 \xi_2] e^{-\frac{1}{2}\alpha \mathbf{w}^2} [1 - \tanh(w_2)]^2}{\int d\mathbf{w} e^{-\frac{1}{2}\alpha \mathbf{w}^2} [1 - \tanh(w_2)]^2} \end{aligned} \quad (3.58)$$

The decision boundary in input space is found to be given by the line $\xi_2 = 0$, for if we substitute this into (3.58) we indeed find $I((\xi_1, 0), D) = 0$. Verification of the behaviour of (3.58) for $\xi_2 \rightarrow \pm\infty$ further shows that $I(\boldsymbol{\xi}, D) > 0$ (and hence $t^*(\boldsymbol{\xi}) = 1$) for $\xi_2 < 0$, and $I(\boldsymbol{\xi}, D) < 0$ (and hence $t^*(\boldsymbol{\xi}) = 0$) for $\xi_2 > 0$.

The more interesting (and less trivial) results concern the dependence of the error measure $\Delta t^*(\boldsymbol{\xi})$ in (3.55) on $\boldsymbol{\xi}$. In figure 3.7 (left panel) we show a contour plot of the mis-classification probability $\Delta t^*(\boldsymbol{\xi})$ in the input plane, for $\alpha = 1$. The model (3.56) assumes a linear separation (a line through the origin, and data noise), the slope of the separating line is extracted from the data. Clearly, the further away from the data points, the larger the possible impact of an uncertainty in this slope, which is seen to be reflected in the contours of $\Delta t^*(\boldsymbol{\xi})$.

We continue with the example model (3.56), but now we add one more data point, namely $(\boldsymbol{\xi}^3, t^3) = ((1, \frac{1}{2}), 0)$ and study its effect on the decision boundary and the mis-classification probability:

$$D = \{((0, -1), 1), ((0, 1), 0), ((1, \frac{1}{2}), 0)\} \quad (3.59)$$

Working out $I(\boldsymbol{\xi}, D)$ of (3.54) now gives

$$\begin{aligned} I(\boldsymbol{\xi}, D) &= \frac{\int d\mathbf{w} \tanh[\mathbf{w} \cdot \boldsymbol{\xi}] e^{-\frac{1}{2}\alpha \mathbf{w}^2} \prod_{\mu=1}^3 p(t^\mu | \boldsymbol{\xi}^\mu, \mathbf{w})}{\int d\mathbf{w} e^{-\frac{1}{2}\alpha \mathbf{w}^2} \prod_{\mu=1}^3 p(t^\mu | \boldsymbol{\xi}^\mu, \mathbf{w})} \\ &= \frac{\int d\mathbf{w} \tanh[w_1 \xi_1 + w_2 \xi_2] e^{-\frac{1}{2}\alpha \mathbf{w}^2} [1 - \tanh(w_2)]^2 \left[1 - \tanh(w_1 + \frac{1}{2}w_2)\right]}{\int d\mathbf{w} e^{-\frac{1}{2}\alpha \mathbf{w}^2} [1 - \tanh(w_2)]^2 \left[1 - \tanh(w_1 + \frac{1}{2}w_2)\right]} \end{aligned} \quad (3.60)$$

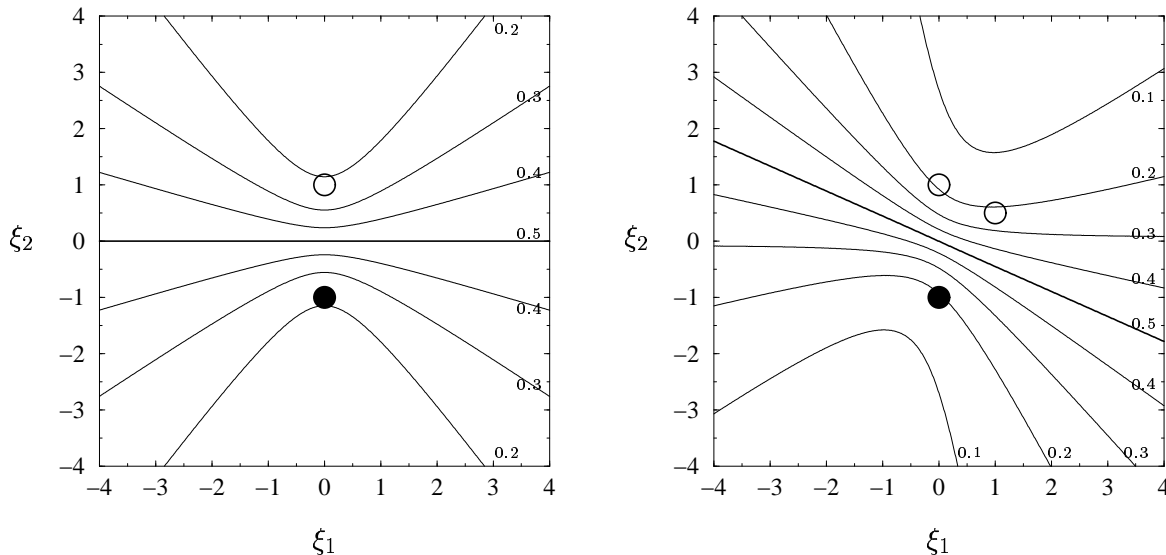


Figure 3.7: Bayesian binary classification of two-dimensional input vectors $\boldsymbol{\xi} = (\xi_1, \xi_2)$ on the basis of a linear parametrized function with (non-Gaussian) data noise and a Gaussian prior (variance: $\alpha^{-1} = 1$). Thick solid line: the decision boundary, where $p(1|\boldsymbol{\xi}, D) = p(0|\boldsymbol{\xi}, D) = \frac{1}{2}$ (equivalently: where $\Delta t^*(\boldsymbol{\xi}) = \frac{1}{2}$). Below this line all points are classified as $t^*(\boldsymbol{\xi}) = 1$, above it as $t^*(\boldsymbol{\xi}) = 0$. Circles: data points $\boldsymbol{\xi}^\mu$ (filled: $t^\mu = 1$, open: $t^\mu = 0$). Thin continuous curves: contour lines of the error probability $\Delta t^*(\boldsymbol{\xi})$. Left panel: $p = 2$. Right panel: $p = 3$ (i.e. one further data point added to those of the left panel; see the main text for details).

The decision boundary is now no longer given by $\xi_2 = 0$; the balance of evidence has changed. In figure 3.7 (right panel) we show a contour plot of the mis-classification probability $\Delta t^*(\boldsymbol{\xi})$ in the input plane for the new situation with the extra data point. Still our model assumes a ‘noisy’ linear separation, but compared to the previous $p = 2$ case we observe both a change in the location of the decision boundary and an overall reduction of the mis-classification probability. Note that, had the new data point been less compatible with the first two, one could also have found an *increase* in the mis-classification probability.

3.5 Bayesian Model Selection

Bayesian Model Comparison. In the formalism described so far one chooses (beforehand) a model assumed responsible for having generated the data, followed by an analysis of the likelihood of parameters for this model. This picture can be generalized to include multiple candidate models. Instead of working with the joint distribution $p(D, \boldsymbol{w})$ of data and model parameters, we must switch to the joint distribution $p(D, \boldsymbol{w}, H)$ to find data D , model H and parameters \boldsymbol{w} for model H ⁹. The generalized picture then becomes

- Consider an *ensemble* of models H with associated parameter vectors \boldsymbol{w} , characterized by a probability distribution $p(H, \boldsymbol{w})$ which evolves during learning.

⁹Note that different models will generally have parameter vectors \boldsymbol{w} with different dimensionality

- Assume that the data D were generated by a system of the form $S(\boldsymbol{\xi}) = f(\boldsymbol{\xi}; \mathbf{w}) + \text{noise}$. Calculate the likelihood $p(\mathbf{w}, H|D)$ of models and their parameters, given the data.
- Express the desired objects $p(\mathbf{w}|D, H)$ in terms of $p(D|\mathbf{w}, H)$ (as before) and $p(H|D)$ in terms of $p(D|H)$, where $p(D|H) = \int d\mathbf{w} p(D|\mathbf{w}, H)p(\mathbf{w}|H)$.

Learning is a process during which the arrival of data reduces our uncertainty about the ‘right’ model H and its ‘right’ parameters \mathbf{w} from the *prior distributions* $p(H)$ and $p(\mathbf{w}|H)$ to the *posterior distributions* $p(H|D)$ and $p(\mathbf{w}|D, H)$. Note that

$$p(\mathbf{w}|D, H) = \frac{p(D|\mathbf{w}, H)p(\mathbf{w}|H)}{\int d\mathbf{w}' p(\mathbf{w}'|H)p(D|\mathbf{w}', H)} \quad (3.61)$$

$$p(H|D) = \frac{p(D|H)p(H)}{\sum_{H'} p(D|H')p(H')} \quad (3.62)$$

Generalized Bayesian learning, with multiple models, now works like this:

Stage 1: definitions

Define (i) the parametrized models H , assumed responsible for the data, (ii) the prior distribution $p(H)$ for these models, (iii) the prior distributions $p(\mathbf{w}|H)$ of their parameters, and (iv) the data $D = \{(\boldsymbol{\xi}^1, t^1), \dots, (\boldsymbol{\xi}^p, t^p)\}$.

Stage 2: model translation

Convert the model definition into a probabilistic standard form: specify the likelihood of finding output t upon presentation of input $\boldsymbol{\xi}$, given model H and parameters \mathbf{w} :

$$\text{model definition in standard form} \quad p(t|\boldsymbol{\xi}, \mathbf{w}, H) \quad (3.63)$$

Stage 3: the posterior distribution

Calculate the *data* likelihood, given model H , $p(D|\mathbf{w}, H) = \prod_{\mu=1}^p p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w}, H)$. From this follow the desired posterior parameter and model distributions

$$p(\mathbf{w}|D, H) = \frac{p(\mathbf{w}|H) \prod_{\mu=1}^p p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w}, H)}{\int d\mathbf{w}' p(\mathbf{w}'|H) \prod_{\mu=1}^p p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w}', H)} \quad (3.64)$$

$$p(H|D) = \frac{p(H) \int d\mathbf{w} \prod_{\mu=1}^p p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w}, H)p(\mathbf{w}|H)}{\sum_{H'} p(H') \int d\mathbf{w} \prod_{\mu=1}^p p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w}, H')p(\mathbf{w}|H')} \quad (3.65)$$

Stage 4: prediction

The residual uncertainty in the choice of model H and parameters \mathbf{w} generates the uncertainty in predictions. Prediction of the output t corresponding to input $\boldsymbol{\xi}$, given our observation of the data D and our choice of model set, takes the probabilistic form:

$$p(t|\boldsymbol{\xi}, D) = \sum_H p(H|D) \int d\mathbf{w} p(t|\boldsymbol{\xi}, \mathbf{w}, H)p(\mathbf{w}|D, H) \quad (3.66)$$

As an alternative to stage 4, which describes the final output statistics as a weighted average over all models under consideration, one could also simply select the most probable model H^* , defined as $p(H^*|D) = \max_H p(H|D)$. This boils down to finding

$$\max_H \{p(H)p(D|H)\} = \max_H \left\{ p(H) \int d\mathbf{w} \prod_{\mu=1}^p p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w}, H) p(\mathbf{w}|H) \right\} \quad (3.67)$$

or, equivalently, can be done by comparing models pair-wise via the ratios

$$\frac{p(H|D)}{p(H'|D)} = \frac{p(D|H)}{p(D|H')} \frac{p(H)}{p(H')} \quad (3.68)$$

This can then be followed by the ordinary Bayesian parameter analysis for a *single* model H^* , as described in the previous sections.

Application to Hyper-parameter Selection. We note that the above reasoning can also be applied to the hyper-parameter selection problem, since systems which differ in the choice of hyper-parameters (rather than in architecture) can be regarded as different candidate models in the sense above. For instance, to describe the family of systems (3.19)

$$p(t|\boldsymbol{\xi}, \mathbf{w}, \beta) = \left[\frac{2\pi}{\beta} \right]^{-\frac{1}{2}} e^{-\frac{1}{2}\beta[t-f(\boldsymbol{\xi};\mathbf{w})]^2} \quad p(\mathbf{w}|\alpha) = \left[\frac{2\pi}{\alpha} \right]^{-\frac{N}{2}} e^{-\frac{1}{2}\alpha\mathbf{w}^2} \quad (3.69)$$

with two hyper-parameters (α, β) , we can simply make the replacement $H \rightarrow (\alpha, \beta)$ in our above formulae. The specific form of (3.69), for which e.g. $p(t|\boldsymbol{\xi}, \mathbf{w}, \alpha, \beta)$ is independent of α , and $p(\mathbf{w}|\alpha, \beta)$ is independent of β , also generates several simplifications:

- Consider an *ensemble* of models of the form (3.69), characterized by a probability distribution $p(\alpha, \beta, \mathbf{w})$ which evolves during learning.
- Assume that the data D were generated by a system of the form (3.69). Calculate the likelihood $p(\mathbf{w}, \alpha, \beta|D)$ of its parameters and hyper-parameters, given the data.
- Express the desired objects $p(\mathbf{w}|D, \alpha, \beta)$ in terms of $p(D|\mathbf{w}, \alpha, \beta)$, and $p(\alpha, \beta|D)$ in terms of $p(D|\alpha, \beta)$, where $p(D|\alpha, \beta) = \int d\mathbf{w} p(D|\mathbf{w}, \beta)p(\mathbf{w}|\alpha)$.

Learning is a process during which the arrival of data reduces our uncertainty about the ‘right’ hyper-parameters (α, β) and the ‘right’ parameters \mathbf{w} from the *prior distributions* $p(\alpha, \beta)$ and $p(\mathbf{w}|\alpha, \beta)$ to the *posterior distributions* $p(\alpha, \beta|D)$ and $p(\mathbf{w}|D, \alpha, \beta)$, according to

$$p(\mathbf{w}|D, \alpha, \beta) = \frac{p(D|\mathbf{w}, \beta)p(\mathbf{w}|\alpha)}{\int d\mathbf{w}' p(D|\mathbf{w}', \beta)p(\mathbf{w}'|\alpha)} \quad (3.70)$$

$$p(\alpha, \beta|D) = \frac{p(D|\alpha, \beta)p(\alpha, \beta)}{\int d\alpha' d\beta' p(D|\alpha', \beta')p(\alpha', \beta')} \quad (3.71)$$

Generalized Bayesian learning, including learning of hyper-parameters, now works like this :

Stage 1: definitions

Define (i) the function $f(\boldsymbol{\xi}; \mathbf{w})$ in the parametrized model (3.69), (ii) the prior distribution $p(\alpha, \beta)$ for its hyper-parameters, and (iii) the data $D = \{(\boldsymbol{\xi}^1, t^1), \dots, (\boldsymbol{\xi}^p, t^p)\}$.

Stage 2: the posterior distribution

Calculate the *data* likelihood, given (α, β) and \mathbf{w} : $p(D|\mathbf{w}, \beta) = \prod_{\mu=1}^p p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w}, \beta)$. From this follow the desired posterior distributions

$$p(\mathbf{w}|D, \alpha, \beta) = \frac{e^{-\frac{1}{2}\alpha\mathbf{w}^2 - \frac{1}{2}\beta \sum_{\mu=1}^p [t^\mu - f(\boldsymbol{\xi}^\mu; \mathbf{w})]^2}}{\int d\mathbf{w}' e^{-\frac{1}{2}\alpha\mathbf{w}'^2 - \frac{1}{2}\beta \sum_{\mu=1}^p [t^\mu - f(\boldsymbol{\xi}^\mu; \mathbf{w}')]^2}} \quad (3.72)$$

$$p(\alpha, \beta|D) = \frac{\frac{p(\alpha, \beta)}{\sqrt{\alpha^N \beta^p}} \int d\mathbf{w} e^{-\frac{1}{2}\alpha\mathbf{w}^2 - \frac{1}{2}\beta \sum_{\mu=1}^p [t^\mu - f(\boldsymbol{\xi}^\mu; \mathbf{w})]^2}}{\int d\alpha' d\beta' \frac{p(\alpha', \beta')}{\sqrt{\alpha'^N \beta'^p}} \int d\mathbf{w} e^{-\frac{1}{2}\alpha'\mathbf{w}^2 - \frac{1}{2}\beta' \sum_{\mu=1}^p [t^\mu - f(\boldsymbol{\xi}^\mu; \mathbf{w})]^2}} \quad (3.73)$$

Stage 3: prediction

Prediction of the output t corresponding to a new input $\boldsymbol{\xi}$, given our observation of the data D and our choice of model family (3.69), takes the probabilistic form:

$$p(t|\boldsymbol{\xi}, D) = \int d\alpha d\beta p(\alpha, \beta|D) \left[\frac{(2\pi)^{N+1}}{\alpha^N \beta} \right]^{-\frac{1}{2}} \int d\mathbf{w} p(\mathbf{w}|D, \alpha, \beta) e^{-\frac{1}{2}\alpha\mathbf{w}^2 - \frac{1}{2}\beta[t - f(\boldsymbol{\xi}; \mathbf{w})]^2} \quad (3.74)$$

Again the learning problem has in principle been reduced to calculating integrals. However, one will now have to think about how to choose the hyper-parameter prior $p(\alpha, \beta)$.

Model Selection – Occam’s Razor. The Bayesian model comparison procedure, followed by model selection according to (3.67) or (3.68), will automatically lead to the selection of the simplest possible model H^* which can account for the data D (modulo prejudices embodied in the prior $p(H)$). This desirable action is referred to as implementing ‘Occam’s Razor’¹⁰ To see how it comes about we return to (3.68), and remove the effect of prejudice by putting $p(H) = p(H')$ for all models $\{H, H'\}$ under consideration. This gives

$$\frac{p(H|D)}{p(H'|D)} = \frac{p(D|H)}{p(D|H')} \quad (3.75)$$

To simplify the argument we will take the collection of possible data sets D to be discrete and countable. Let us compare the following three models:

$$\text{model } H_1 : \quad p(D|H_1) = 0$$

$$\text{model } H_2 : \quad p(D|H_2) > 0, \quad p(D'|H_2) = 0 \text{ for all } D' \neq D$$

$$\text{model } H_3 : \quad p(D|H_3) > 0, \quad p(D'|H_3) > 0 \text{ for some } D' \neq D$$

Model 1 cannot explain the data D . Both model 2 and model 3 can explain these data, but 3 is more complex than 2 because it can explain (by a suitable choice of parameters) a greater variety of possible data (not just D). Hence model 2 is the *simplest* model which can account for data D . Insertion into (3.75) reveals

$$\frac{p(H_1|D)}{p(H_2|D)} = \frac{p(H_1|D)}{p(H_3|D)} = 0 \quad \frac{p(H_2|D)}{p(H_3|D)} = \frac{p(D|H_2)}{p(D|H_3)} \quad (3.76)$$

¹⁰After a monk, William of Occam, who is claimed to have first proposed the general philosophical principle that one should always select the *simplest* possible explanation for an observed phenomenon.

One now always has $\sum_{D'} p(D'|H) = 1$ (marginal distributions are normalized, for any model H). From this it follows for the above examples that $p(D|H_2) = 1 - \sum_{D' \neq D} p(D'|H_2) = 1$ and $p(D|H_3) = 1 - \sum_{D' \neq D} p(D'|H_3) < 1$. Hence we find in (3.76) that $p(H_2|D) > p(H_3|D) > p(H_1|D)$, and that the Bayesian procedure instructs us to select model 2.

Example. We close with a simple example, to illustrate the action of ‘Occam’s Razor’. The task is to learn a binary classification of the inputs $\xi \in \{0, 1\}$ to the targets $t \in \{0, 1\}$. There are four such classifications possible, i.e. four possible data sets with $p = 2$:

$$D_A = \{(0, 1), (1, 0)\} \quad D_B = \{(0, 0), (1, 0)\} \quad D_C = \{(0, 1), (1, 1)\} \quad D_D = \{(0, 0), (1, 1)\}$$

Given one of these data sets, we have to choose between two deterministic parametrized candidate models H_1 and H_2 , without initial prejudice, i.e. $p(H_1) = p(H_2) = \frac{1}{2}$:

$$H_1 : \quad t(\xi) = \theta[w_1] \tag{3.77}$$

$$H_2 : \quad t(\xi) = \theta[w_1 + w_2\xi] \tag{3.78}$$

with $\mathbf{w} = (w_1, w_2) \in \mathfrak{R}^2$ and $\theta[z]$ denoting the step function. As a prior parameter distribution we take $p(\mathbf{w}) = (2\pi)^{-1} e^{-\frac{1}{2}\mathbf{w}^2}$. Clearly model 2 is more complex than 1. We will select our model by working out the ratio (3.68), this requires calculating $p(D|H_1) = \int d\mathbf{w}_1 p(w_1) p(D|w_1, H_1)$ and $p(D|H_2) = \int d\mathbf{w} p(\mathbf{w}) p(D|\mathbf{w}, H_2)$.

First we focus on model H_1 . Here we have $t(\xi) = 1$ for all ξ if $w_1 > 0$, and $t(\xi) = 0$ for all ξ if $w_1 < 0$. Hence we simply find

$$p(D_A|w_1, H_1) = p(D_D|w_1, H_1) = 0, \quad p(D_B|w_1, H_1) = \theta[-w_1], \quad p(D_C|w_1, H_1) = \theta[w_1]$$

This gives

$$\begin{aligned} p(D_A|H_1) &= 0 \\ p(D_B|H_1) &= \int_{-\infty}^0 dw_1 p(w_1) = \frac{1}{2} \\ p(D_C|H_1) &= \int_0^{\infty} dw_1 p(w_1) = \frac{1}{2} \\ p(D_D|H_1) &= 0 \end{aligned} \tag{3.79}$$

Next we turn to the two-parameter model H_2 . Here we have $t(0) = \theta[w_1]$ and $t(1) = \theta[w_1 + w_2]$. Hence we find

$$\begin{aligned} p(D_A|\mathbf{w}, H_2) &= \theta[w_1] \theta[-(w_1 + w_2)] & p(D_B|\mathbf{w}, H_2) &= \theta[-w_1] \theta[-(w_1 + w_2)] \\ p(D_C|\mathbf{w}, H_2) &= \theta[w_1] \theta[w_1 + w_2] & p(D_D|\mathbf{w}, H_2) &= \theta[-w_1] \theta[w_1 + w_2] \end{aligned}$$

This leads to, with the short-hand $\chi = (2\pi)^{-1} \int_0^{\infty} dw_1 e^{-\frac{1}{2}w_1^2} \int_{w_1}^{\infty} dw_2 e^{-\frac{1}{2}w_2^2} > 0$ and upon using the symmetries of the prior $p(\mathbf{w})$:

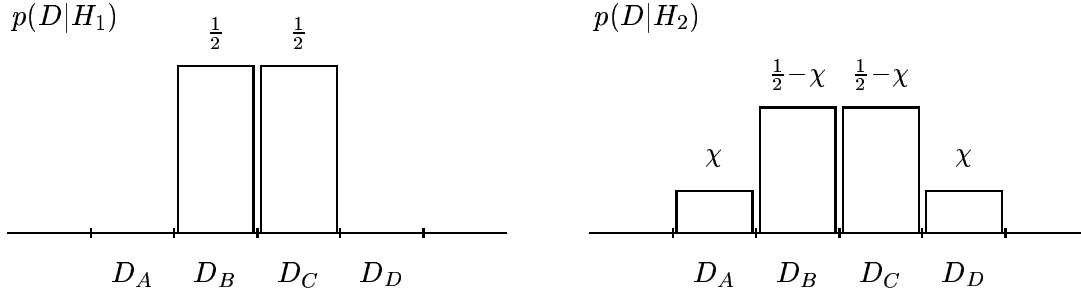
$$p(D_A|H_2) = \int_0^{\infty} dw_1 \int_{-\infty}^{-w_1} dw_2 p(\mathbf{w}) = \chi \tag{3.80}$$

$$p(D_B|H_2) = \int_{-\infty}^0 dw_1 \int_{-\infty}^{-w_1} dw_2 p(\mathbf{w}) = \int_0^{\infty} dw_1 \int_{-\infty}^{w_1} dw_2 p(\mathbf{w}) = \frac{1}{2} - \chi \tag{3.81}$$

$$p(D_C|H_2) = \int_0^\infty dw_1 \int_{-w_1}^\infty dw_2 p(\mathbf{w}) = \frac{1}{2} - \chi \quad (3.82)$$

$$p(D_D|H_2) = \int_{-\infty}^0 dw_1 \int_{-w_1}^\infty dw_2 p(\mathbf{w}) = \int_0^\infty dw_1 \int_{w_1}^\infty dw_2 p(\mathbf{w}) = \chi \quad (3.83)$$

Finally we combine the results (3.79) and (3.80-3.83) into the following picture:



It follows from (3.68) that

$$\frac{p(H_1|D_A)}{p(H_2|D_A)} = \frac{p(H_1|D_D)}{p(H_2|D_D)} = 0, \quad \frac{p(H_1|D_B)}{p(H_2|D_B)} = \frac{p(H_1|D_C)}{p(H_2|D_C)} = \frac{1/2}{1/2 - \chi} > 1$$

We conclude that when observing data D_A or D_D we must select model H_2 (which would be the only candidate to explain these data), but that when observing data D_B or D_C we must select model H_1 (both models can explain these data, but H_1 wins simply because it is the simpler explanation).

3.6 Practicalities: Measuring Curvature

When learning by gradient descent on the surface $S(\mathbf{w}, D)$ (3.16), and using simple expressions such as (3.38) to assign error bars to neural network predictions, one needs to know the curvature matrix \mathbf{A} as given in (3.23) at the minimum \mathbf{w}_{MP} of $S(\mathbf{w}, D)$. In principle this matrix can be calculated directly from the model definitions, but for many-parameter and/or multi-layer networks this is difficult. Here we discuss an alternative, based on the idea that the degree of curvature around the minimum will have a direct effect on the gradient descent dynamics at the minimum when such dynamics is equipped with additive noise. As a result the curvature matrix can be extracted from a measurement of the fluctuations.

Let us assume we have arrived by gradient descent at the minimum \mathbf{w}_{MP} of $S(\mathbf{w}, D)$. We now replace the (discretized) gradient descent dynamics by the following noisy version:

$$w_i(t + \delta) = w_i(t) - \delta \cdot \frac{\partial S(\mathbf{w}, D)}{\partial w_i} + \sqrt{2\delta} \eta_i(t) \quad (3.84)$$

where the $\eta_i(t)$ are independently distributed zero-average Gaussian random variables, with $\langle \eta_i^2(t) \rangle = 1$, and with $0 < \delta \ll 1$. The parameter dynamics is now a stochastic process. Let

us define averages and covariances of this process as

$$\bar{w}_i(t) = \langle w_i(t) \rangle \quad C_{ij}(t) = \langle [w_i(t) - \bar{w}_i(t)][w_j(t) - \bar{w}_j(t)] \rangle \quad (3.85)$$

Their dynamics follow directly from (3.84):

$$\begin{aligned} \bar{w}_i(t + \delta) &= \bar{w}_i(t) - \delta \left\langle \frac{\partial S(\mathbf{w}, D)}{\partial w_i} \right\rangle \\ C_{ij}(t + \delta) &= \langle w_i(t + \delta)w_j(t + \delta) \rangle - \bar{w}_i(t + \delta)\bar{w}_j(t + \delta) \\ &= \left\langle \left[w_i(t) - \delta \frac{\partial S(\mathbf{w}, D)}{\partial w_i} \right] \left[w_j(t) - \delta \frac{\partial S(\mathbf{w}, D)}{\partial w_j} \right] \right\rangle + 2\delta \delta_{ij} - \bar{w}_i(t + \delta)\bar{w}_j(t + \delta) \\ &= C_{ij}(t) + \delta \left\{ 2\delta_{ij} - \langle w_i(t) \frac{\partial S(\mathbf{w}, D)}{\partial w_j} \rangle + \bar{w}_i(t) \left\langle \frac{\partial S(\mathbf{w}, D)}{\partial w_j} \right\rangle \right. \\ &\quad \left. - \langle w_j(t) \frac{\partial S(\mathbf{w}, D)}{\partial w_i} \rangle + \bar{w}_j(t) \left\langle \frac{\partial S(\mathbf{w}, D)}{\partial w_i} \right\rangle \right\} + \mathcal{O}(\delta^2) \end{aligned}$$

Since $\delta \ll 1$, the system can only make infinitesimal excursions $\mathbf{w} = \mathbf{w}_{\text{MP}} + \mathcal{O}(\sqrt{\delta})$ away from the minimum \mathbf{w}_{MP} , hence we may approximate $S(\mathbf{w}, D)$ safely by the quadratic approximation (3.22)

$$S(\mathbf{w}, D) = S(\mathbf{w}_{\text{MP}}, D) + \frac{1}{2}(\mathbf{w} - \mathbf{w}_{\text{MP}}) \cdot \mathbf{A}(\mathbf{w} - \mathbf{w}_{\text{MP}}) + \mathcal{O}(\delta^{3/2})$$

This implies that

$$\frac{\partial S(\mathbf{w}, D)}{\partial w_i} = \sum_j A_{ij}(w_j - w_{\text{MP},j}) + \mathcal{O}(\delta^{3/2})$$

Thus:

$$\begin{aligned} \bar{\mathbf{w}}(t + \delta) &= \bar{\mathbf{w}}(t) - \delta \mathbf{A}(\bar{\mathbf{w}} - \mathbf{w}_{\text{MP}}) + \mathcal{O}(\delta^{3/2}) \\ C_{ij}(t + \delta) &= C_{ij}(t) + \delta \left\{ 2\delta_{ij} - \langle w_i(t) \sum_k A_{jk}(w_k(t) - w_{\text{MP},k}) \rangle + \bar{w}_i(t) \sum_k A_{jk}(\bar{w}_k - w_{\text{MP},k}) \right. \\ &\quad \left. - \langle w_j(t) \sum_k A_{ik}(w_k(t) - w_{\text{MP},k}) \rangle + \bar{w}_j(t) \sum_k A_{ik}(\bar{w}_k - w_{\text{MP},k}) \right\} + \mathcal{O}(\delta^{3/2}) \end{aligned}$$

We now arrive at the simple result (using the symmetry of both \mathbf{A} and \mathbf{C}):

$$\mathbf{C}(t + \delta) = \mathbf{C}(t) + \delta [2\mathbf{I} - \mathbf{C}\mathbf{A} - \mathbf{A}\mathbf{C}] + \mathcal{O}(\delta^{3/2}) \quad (3.86)$$

One can show that, provided $\delta \ll 1$ (more specifically: provided $\delta < 2/a_i$ for all eigenvalues a_i of the matrix \mathbf{A}), the leading order of (3.86) will evolve towards the stationary state $\mathbf{C} = \mathbf{A}^{-1} + \mathcal{O}(\sqrt{\delta})$, and the vector $\bar{\mathbf{w}}$ will evolve towards \mathbf{w}_{MP} . Hence the desired curvature matrix \mathbf{A} can be measured in equilibrium, by choosing δ sufficiently small, according to

$$\delta \rightarrow 0: \quad \mathbf{A} = \mathbf{C}^{-1} \quad C_{ij} = \lim_{t \rightarrow \infty} \langle [w_i - w_{\text{MP},i}][w_j - w_{\text{MP},j}] \rangle \quad (3.87)$$

Chapter 4

Gaussian Processes

We return to the general Bayesian formalism for a single model, characterized by the likelihood $p(t|\boldsymbol{\xi}, \mathbf{w})$ to find output t , given input $\boldsymbol{\xi}$ and parameters \mathbf{w} . The data are, as always, given by p pairs of example inputs $\boldsymbol{\xi}^\mu$ with corresponding (possibly noisy) outputs t^μ :

$$D = \{(\boldsymbol{\xi}^1, t^1), \dots, (\boldsymbol{\xi}^p, t^p)\}$$

For reasons to become clear below we will now write the data set D in expressions such as $p(t|\boldsymbol{\xi}, D) = \int d\mathbf{w} p(t|\boldsymbol{\xi}, \mathbf{w})p(\mathbf{w}|D)$ as the full specification of the constituent pairs $(\boldsymbol{\xi}^\mu, t^\mu)$:

$$p(t|\boldsymbol{\xi}, \boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, t^1, \dots, t^p) = \int d\mathbf{w} p(t|\boldsymbol{\xi}, \mathbf{w})p(\mathbf{w}|\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, t^1, \dots, t^p) \quad (4.1)$$

Multiplication of this expression by $p(t^1, \dots, t^p|\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p)$ gives:

$$\begin{aligned} p(t, t^1, \dots, t^p|\boldsymbol{\xi}, \boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p) &= \int d\mathbf{w} p(t|\boldsymbol{\xi}, \mathbf{w})p(\mathbf{w}, t^1, \dots, t^p|\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p) \\ &= \int d\mathbf{w} p(\mathbf{w})p(t|\boldsymbol{\xi}, \mathbf{w})p(t^1, \dots, t^p|\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \mathbf{w}) \end{aligned}$$

Upon renaming the target input and output as $(\boldsymbol{\xi}^0, t^0) = (\boldsymbol{\xi}, t)$, we arrive at the simple and transparent expression

$$p(t^0, \dots, t^p|\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p) = \int d\mathbf{w} p(\mathbf{w}) \prod_{\mu=0}^p p(t^\mu|\boldsymbol{\xi}^\mu, \mathbf{w}) \quad (4.2)$$

We now see that the desired prediction output distribution (4.1) can be simply obtained as the fraction

$$p(t^0|\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p, t^1, \dots, t^p) = \frac{p(t^0, \dots, t^p|\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p)}{p(t^1, \dots, t^p|\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p)} \quad (4.3)$$

or, written slightly differently:

$$p(t^0|\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p, t^1, \dots, t^p) = \frac{p(t^0, t^1, \dots, t^p|\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p)}{\int dt p(t, t^1, \dots, t^p|\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p)}$$

Gaussian processes are information processing systems of the above form $p(t|\boldsymbol{\xi}, \mathbf{w})$ such that the key object (4.2) is a (multivariate) Gaussian distribution of $\{t^0, \dots, t^p\}$, with moments generally depending on $\{\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p\}$ in a non-trivial and model-specific way¹.

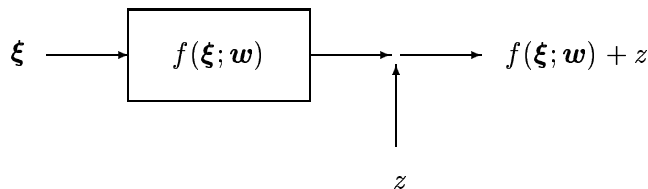
¹Note: in literature one can sometimes find slightly more broad definitions, where upon adding non-Gaussian output noise to Gaussian-distributed targets one would still speak about a Gaussian process.

4.1 Examples of Networks Reducing to Gaussian Processes

In order to show that Gaussian processes are not just found upon choosing far-fetched and/or trivial models, we will first discuss two examples of systems reducing to Gaussian processes. Both are of the form

$$p(t|\boldsymbol{\xi}, \mathbf{w}) = \left[\frac{\beta}{2\pi} \right]^{\frac{1}{2}} e^{-\frac{1}{2}\beta[t-f(\boldsymbol{\xi};\mathbf{w})]^2}, \quad p(\mathbf{w}) = \left[\frac{\alpha}{2\pi} \right]^{\frac{N}{2}} e^{-\frac{1}{2}\alpha\mathbf{w}^2} \quad (4.4)$$

i.e. composed of a deterministic parametrized function $f(\mathbf{w}; \boldsymbol{\xi})$, with a Gaussian prior and Gaussian output noise.



Radial Basis Function Networks. RBF systems with Gaussian output noise and priors, which are members of the class (4.4), are described by a deterministic part of the form $f(\mathbf{w}; \boldsymbol{\xi}) = \mathbf{w} \cdot \boldsymbol{\phi}(\boldsymbol{\xi})$, with $\mathbf{w} \in \Re^N$ and with N functions $\phi_i(\boldsymbol{\xi})$ to be specified². We now work out the joint output distribution (4.2) for data and target. With a modest amount of foresight we define

$$A_{ij} = \alpha\delta_{ij} + \beta \sum_{\mu=0}^p \phi_i(\boldsymbol{\xi}^\mu)\phi_j(\boldsymbol{\xi}^\mu) \quad b_i = \beta \sum_{\mu=0}^p t^\mu \phi_i(\boldsymbol{\xi}^\mu) \quad (4.5)$$

Insertion of (4.4) into (4.2), with $f(\boldsymbol{\xi}; \mathbf{w}) = \mathbf{w} \cdot \boldsymbol{\phi}(\boldsymbol{\xi})$, and subsequent usage of the short-hands (4.5) and the Gaussian integrals in appendix B, gives

$$\begin{aligned} p(t^0, \dots, t^p | \boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p) &= \left[\frac{\alpha}{2\pi} \right]^{\frac{N}{2}} \left[\frac{\beta}{2\pi} \right]^{\frac{p+1}{2}} \int d\mathbf{w} e^{-\frac{1}{2}\alpha\mathbf{w}^2 - \frac{1}{2}\beta \sum_{\mu=0}^p [t^\mu - \mathbf{w} \cdot \boldsymbol{\phi}(\boldsymbol{\xi}^\mu)]^2} \\ &= \left[\frac{\alpha}{2\pi} \right]^{\frac{N}{2}} \left[\frac{\beta}{2\pi} \right]^{\frac{p+1}{2}} e^{-\frac{1}{2}\beta \sum_{\mu=0}^p [t^\mu]^2} \int d\mathbf{w} e^{-\frac{1}{2}\mathbf{w} \cdot \mathbf{A} \mathbf{w} + \mathbf{w} \cdot \mathbf{b}} \\ &= \left[\frac{\alpha}{2\pi} \right]^{\frac{N}{2}} \left[\frac{\beta}{2\pi} \right]^{\frac{p+1}{2}} e^{-\frac{1}{2}\beta \sum_{\mu=0}^p [t^\mu]^2} \frac{(2\pi)^{N/2}}{\sqrt{\det \mathbf{A}}} e^{\frac{1}{2}\mathbf{b} \cdot \mathbf{A}^{-1} \mathbf{b}} \\ &= \left[\frac{\alpha^N \beta^{p+1}}{(2\pi)^{p+1} \det \mathbf{A}} \right]^{\frac{1}{2}} e^{-\frac{1}{2}\beta \sum_{\mu=0}^p [t^\mu]^2 + \frac{1}{2}\beta^2 [\sum_{\mu=0}^p t^\mu \boldsymbol{\phi}(\boldsymbol{\xi}^\mu)] \cdot \mathbf{A}^{-1} [\sum_{\mu=0}^p t^\mu \boldsymbol{\phi}(\boldsymbol{\xi}^\mu)]} \\ &= \left[\frac{\alpha^N \beta^{p+1}}{(2\pi)^{p+1} \det \mathbf{A}} \right]^{\frac{1}{2}} e^{-\frac{1}{2} \sum_{\mu, \nu=0}^p t^\mu t^\nu [\beta \delta_{\mu\nu} - \beta^2 \boldsymbol{\phi}(\boldsymbol{\xi}^\mu) \cdot \mathbf{A}^{-1} \boldsymbol{\phi}(\boldsymbol{\xi}^\nu)]} \end{aligned} \quad (4.6)$$

The output distribution (4.6) is clearly a (zero-average) Gaussian one, characterized fully by the covariance matrix $\mathbf{C}[\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p]$ which is defined as:

$$C_{\mu\nu}[\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p] = \int dt^0 \dots dt^p t^\mu t^\nu p(t^0, \dots, t^p | \boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p)$$

²Note: the model family (4.4) with $f(\boldsymbol{\xi}; \mathbf{w}) = \mathbf{w} \cdot \boldsymbol{\phi}(\boldsymbol{\xi})$ is in fact more general than just the set of all RBF networks, since we need not put RBF-type restrictions on the form of the functions $\phi_i(\boldsymbol{\xi})$.

According to (4.6) we have (see also appendix B):

$$(\mathbf{C}^{-1})_{\mu\nu}[\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p] = \beta\delta_{\mu\nu} - \beta^2\boldsymbol{\phi}(\boldsymbol{\xi}^\mu) \cdot \mathbf{A}^{-1}\boldsymbol{\phi}(\boldsymbol{\xi}^\nu)$$

Inversion of the above matrix gives

$$C_{\mu\nu}[\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p] = \frac{1}{\beta}\delta_{\mu\nu} + \frac{1}{\alpha}\boldsymbol{\phi}(\boldsymbol{\xi}^\mu) \cdot \boldsymbol{\phi}(\boldsymbol{\xi}^\nu) \quad (4.7)$$

This can be verified by insertion:

$$\begin{aligned} \beta \sum_{\rho} C_{\mu\rho}[\dots] [\delta_{\rho\nu} - \beta\boldsymbol{\phi}(\boldsymbol{\xi}^\rho) \cdot \mathbf{A}^{-1}\boldsymbol{\phi}(\boldsymbol{\xi}^\nu)] &= \sum_{\rho} \left[\delta_{\mu\rho} + \frac{\beta}{\alpha}\boldsymbol{\phi}(\boldsymbol{\xi}^\mu) \cdot \boldsymbol{\phi}(\boldsymbol{\xi}^\rho) \right] [\delta_{\rho\nu} - \beta\boldsymbol{\phi}(\boldsymbol{\xi}^\rho) \cdot \mathbf{A}^{-1}\boldsymbol{\phi}(\boldsymbol{\xi}^\nu)] \\ &= \delta_{\mu\nu} - \beta\boldsymbol{\phi}(\boldsymbol{\xi}^\mu) \cdot \mathbf{A}^{-1}\boldsymbol{\phi}(\boldsymbol{\xi}^\nu) + \frac{\beta}{\alpha}\boldsymbol{\phi}(\boldsymbol{\xi}^\mu) \cdot \boldsymbol{\phi}(\boldsymbol{\xi}^\nu) - \frac{\beta^2}{\alpha} \sum_{\rho} [\boldsymbol{\phi}(\boldsymbol{\xi}^\mu) \cdot \boldsymbol{\phi}(\boldsymbol{\xi}^\rho)] [\boldsymbol{\phi}(\boldsymbol{\xi}^\rho) \cdot \mathbf{A}^{-1}\boldsymbol{\phi}(\boldsymbol{\xi}^\nu)] \\ &= \delta_{\mu\nu} - \frac{\beta}{\alpha} \sum_{ij} \phi_i(\boldsymbol{\xi}^\mu) \phi_j(\boldsymbol{\xi}^\nu) \left[\alpha(\mathbf{A}^{-1})_{ij} - \delta_{ij} + \beta \sum_k \sum_{\rho} \phi_i(\boldsymbol{\xi}^\rho) \phi_k(\boldsymbol{\xi}^\rho) (\mathbf{A}^{-1})_{kj} \right] \\ &= \delta_{\mu\nu} - \frac{\beta}{\alpha} \sum_{ij} \phi_i(\boldsymbol{\xi}^\mu) \phi_j(\boldsymbol{\xi}^\nu) \left[\alpha(\mathbf{A}^{-1})_{ij} - \delta_{ij} + \sum_k [A_{ik} - \alpha\delta_{ik}] (\mathbf{A}^{-1})_{kj} \right] = \delta_{\mu\nu} \end{aligned}$$

This confirms (4.7).

Linear-Output Two-Layer Perceptrons. Let us next turn a less trivial model example, still of the form (4.4), which is also found to reduce to a Gaussian process in a specific limit. We take a two-layer network with sigmoidal transfer functions $g(u)$ in the hidden layer (of size L) and a linear output neuron, i.e. $f(\boldsymbol{\xi}; \mathbf{w}, \mathbf{J}) = \sum_{i=1}^L w_i g(\sum_{j=1}^N J_{ij}\xi_j)$:

$$p(t|\boldsymbol{\xi}, \mathbf{w}, \mathbf{J}) = \left[\frac{\beta}{2\pi} \right]^{\frac{1}{2}} e^{-\frac{1}{2}\beta[t - \sum_{i=1}^L w_i g(\sum_{j=1}^N J_{ij}\xi_j)]^2} \quad (4.8)$$

$$p(\mathbf{w}) = \left[\frac{\alpha_w L}{2\pi} \right]^{\frac{L}{2}} e^{-\frac{1}{2}\alpha_w L \mathbf{w}^2}, \quad p(\mathbf{J}) = \left[\frac{\alpha_J}{2\pi} \right]^{\frac{NL}{2}} e^{-\frac{1}{2}\alpha_J \sum_{i=1}^L \sum_{j=1}^N J_{ij}^2} \quad (4.9)$$

The width of the prior for the L hidden-to-output weights $\{w_i\}$ has been re-scaled by a factor $1/\sqrt{L}$, in order to find a well-defined limit $L \rightarrow \infty$ (corresponding to an infinitely large hidden layer) later. This model has three hyper-parameters: α_w , α_J and β . Insertion of the above definitions into (4.2) is now found to give

$$\begin{aligned} p(t^0, \dots, t^p | \boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p) &= \left[\frac{\alpha_w L}{2\pi} \right]^{\frac{L}{2}} \left[\frac{\alpha_J}{2\pi} \right]^{\frac{NL}{2}} \left[\frac{\beta}{2\pi} \right]^{\frac{p+1}{2}} \\ &\times \int d\mathbf{w} d\mathbf{J} e^{-\frac{1}{2}\alpha_w L \mathbf{w}^2 - \frac{1}{2}\alpha_J \sum_{i=1}^L \sum_{j=1}^N J_{ij}^2 - \frac{1}{2}\beta \sum_{\mu=0}^p [t^\mu - \sum_{i=1}^L w_i g(\sum_{j=1}^N J_{ij}\xi_j^\mu)]^2} \end{aligned}$$

It will be advantageous to introduce $p+1$ additional integrals to get rid of the squares in the exponent, via the identity

$$e^{-\frac{1}{2}\sum_{\mu=0}^p K_\mu^2} = (2\pi)^{-\frac{p+1}{2}} \int d\mathbf{k} e^{-\frac{1}{2}\mathbf{k}^2 + i\sum_{\mu=0}^p k_\mu K_\mu}$$

with $\mathbf{k} = (k_0, \dots, k_p)$. This allows us to write

$$\begin{aligned}
p(t^0, \dots, t^p | \boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p) &= \left[\frac{\alpha_w L}{2\pi} \right]^{\frac{L}{2}} \left[\frac{\alpha_J}{2\pi} \right]^{\frac{NL}{2}} \left[\frac{\beta}{(2\pi)^2} \right]^{\frac{p+1}{2}} \int d\mathbf{k} e^{-\frac{1}{2}\mathbf{k}^2 + i\sqrt{\beta} \sum_{\mu=0}^p k_\mu t_\mu} \\
&\quad \times \int d\mathbf{J} e^{-\frac{1}{2}\alpha_J \sum_{i=1}^L \sum_{j=1}^N J_{ij}^2} \prod_{i=1}^L \left[\int dw_i e^{-\frac{1}{2}\alpha_w L w_i^2 - i w_i \sqrt{\beta} \sum_{\mu=0}^p k_\mu g(\sum_{j=1}^N J_{ij} \xi_j^\mu)} \right] \\
&= \left[\frac{\alpha_w L}{2\pi} \right]^{\frac{L}{2}} \left[\frac{\alpha_J}{2\pi} \right]^{\frac{NL}{2}} \left[\frac{\beta}{(2\pi)^2} \right]^{\frac{p+1}{2}} \int d\mathbf{k} e^{-\frac{1}{2}\mathbf{k}^2 + i\sqrt{\beta} \sum_{\mu=0}^p k_\mu t_\mu} \int d\mathbf{J} e^{-\frac{1}{2}\alpha_J \sum_{i=1}^L \sum_{j=1}^N J_{ij}^2} \\
&\quad \times \prod_{i=1}^L \left[\int dw e^{-\frac{1}{2}\alpha_w L [w + \frac{i\sqrt{\beta}}{\alpha_w L} \sum_{\mu=0}^p k_\mu g(\sum_{j=1}^N J_{ij} \xi_j^\mu)]^2 - \frac{\beta}{2\alpha_w L} [\sum_{\mu=0}^p k_\mu g(\sum_{j=1}^N J_{ij} \xi_j^\mu)]^2} \right] \\
&= \left[\frac{\alpha_J}{2\pi} \right]^{\frac{NL}{2}} \left[\frac{\beta}{(2\pi)^2} \right]^{\frac{p+1}{2}} \int d\mathbf{k} e^{-\frac{1}{2}\mathbf{k}^2 + i\sqrt{\beta} \sum_{\mu=0}^p k_\mu t_\mu} \\
&\quad \times \int d\mathbf{J} e^{-\frac{1}{2}\alpha_J \sum_{i=1}^L \sum_{j=1}^N J_{ij}^2} e^{-\frac{\beta}{2\alpha_w} \sum_{\mu\nu=0}^p k_\mu k_\nu \left[\frac{1}{L} \sum_{i=1}^L g(\sum_{j=1}^N J_{ij} \xi_j^\mu) g(\sum_{j=1}^N J_{ij} \xi_j^\nu) \right]}
\end{aligned}$$

We re-organize the weights $\{J_{ij}\}$ into L vectors \mathbf{v}_i of dimension N each, $\mathbf{v}_i = (J_{i1}, J_{i2}, \dots, J_{iN})$:

$$\begin{aligned}
p(t^0, \dots, t^p | \boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p) &= \left[\frac{\alpha_J}{2\pi} \right]^{\frac{NL}{2}} \left[\frac{\beta}{(2\pi)^2} \right]^{\frac{p+1}{2}} \int d\mathbf{k} e^{-\frac{1}{2}\mathbf{k}^2 + i\sqrt{\beta} \sum_{\mu=0}^p k_\mu t_\mu} \\
&\quad \times \int \prod_{i=1}^L [d\mathbf{v}_i e^{-\frac{1}{2}\alpha_J \mathbf{v}_i^2}] e^{-\frac{\beta}{2\alpha_w} \sum_{\mu\nu=0}^p k_\mu k_\nu \left[\frac{1}{L} \sum_{i=1}^L g(\mathbf{v}_i \cdot \boldsymbol{\xi}^\mu) g(\mathbf{v}_i \cdot \boldsymbol{\xi}^\nu) \right]}
\end{aligned}$$

If now we take the limit $L \rightarrow \infty$ it is clear that, due to the fact that all Gaussian random vectors \mathbf{v}_i are distributed identically but also *independently*, we may use the general identity

$$\lim_{L \rightarrow \infty} \frac{1}{L} \sum_{i=1}^L F(\mathbf{v}_i) = \langle F(\mathbf{v}) \rangle_{\mathbf{v}} = \int d\mathbf{v} p(\mathbf{v}) F(\mathbf{v})$$

where $p(\mathbf{v})$ denotes the distribution of an individual \mathbf{v}_i , i.e. $p(\mathbf{v}) = (\alpha_J/2\pi)^{N/2} e^{-\frac{1}{2}\alpha_J \mathbf{v}^2}$. Hence, after a further simplifying transformation $\mathbf{k} \rightarrow \mathbf{k}/\sqrt{\beta}$ we find

$$\begin{aligned}
\lim_{L \rightarrow \infty} p(t^0, \dots, t^p | \boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p) &= \int \frac{d\mathbf{k}}{(2\pi)^{p+1}} e^{i \sum_{\mu=0}^p k_\mu t_\mu - \frac{1}{2} \sum_{\mu\nu=0}^p k_\mu k_\nu \left[\frac{1}{\beta} \delta_{\mu\nu} + \frac{1}{\alpha_w} \langle g(\mathbf{v} \cdot \boldsymbol{\xi}^\mu) g(\mathbf{v} \cdot \boldsymbol{\xi}^\nu) \rangle_{\mathbf{v}} \right]} \\
&= \frac{e^{-\frac{1}{2} \sum_{\mu\nu=0}^p t^\mu (\mathbf{C}[\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p])_{\mu\nu}^{-1} t^\nu}}{\sqrt{(2\pi)^{p+1} \det \mathbf{C}[\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p]}} \tag{4.10}
\end{aligned}$$

Again we find a Gaussian distribution (4.10), here characterized by the covariance matrix

$$C_{\mu\nu}[\boldsymbol{\xi}^0, \dots, \boldsymbol{\xi}^p] = \frac{1}{\beta} \delta_{\mu\nu} + \frac{1}{\alpha_w} \int \frac{d\mathbf{z}}{(2\pi)^{N/2}} e^{-\frac{1}{2}\mathbf{z}^2} g\left(\frac{\mathbf{z} \cdot \boldsymbol{\xi}^\mu}{\sqrt{\alpha_J}}\right) g\left(\frac{\mathbf{z} \cdot \boldsymbol{\xi}^\nu}{\sqrt{\alpha_J}}\right) \tag{4.11}$$

Thus, also two-layer feed-forward networks with sigmoidal transfer functions in the hidden layer and a linear output neuron, and suitably scaled Gaussian priors, are found to become Gaussian processes in the limit of an infinitely large hidden layer.

4.2 Learning and Prediction with Gaussian Processes

Definitions and Rationale. The association of outputs $t(\boldsymbol{\xi}) \in \mathfrak{R}$ to inputs $\boldsymbol{\xi} \in \mathfrak{R}^N$ is said to be governed by a Gaussian process if for any finite selection of p inputs $\{\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p\}$ the conditional joint output distribution $p(t^1, \dots, t^p | \boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p)$ is Gaussian:

$$p(t^1, \dots, t^p | \boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p) = \frac{e^{-\frac{1}{2}(\mathbf{t} - \mathbf{a}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p]) \cdot \mathbf{C}^{-1}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p](\mathbf{t} - \mathbf{a}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p])}}{\sqrt{(2\pi)^p \det \mathbf{C}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p]}} \quad (4.12)$$

with the short-hands $\mathbf{t} = (t^1, \dots, t^p)$, $\mathbf{a}[\dots] = (a_1[\dots], \dots, a_p[\dots])$, and with

$$a_\mu[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p] = \int dt p(t^1, \dots, t^p | \boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p) t_\mu \quad (4.13)$$

$$C_{\mu\nu}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p] = \int dt p(t^1, \dots, t^p | \boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p) (t_\mu - a_\mu[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p])(t_\nu - a_\nu[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p]) \quad (4.14)$$

We note:

- Gaussian processes are fully characterized by the functions in (4.13,4.14). The detailed dependence of these averages $a_\mu[\dots]$ and covariances $C_{\mu\nu}[\dots]$ on the arguments $\{\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p\}$ reflects both structural properties and prior parameter assumptions of the underlying models (see e.g. the explicit derivations in the previous section).
- The rationale behind using Gaussian processes is: (i) it allows us to move away from model specific details (e.g. network architecture and weights) and deal with a reasonably large class of systems in one go, and (ii) we will be able to derive surprisingly simple (and hence CPU efficient) expressions for predicted outputs and associated error bars.

Our formulae hold (by definition) for any value of p . To keep notation simple, it will turn out helpful to refer to the new input $\boldsymbol{\xi}$ and its output t to be predicted (which in the previous section we called $\boldsymbol{\xi}^0$ and t^0) as $\boldsymbol{\xi}^{p+1}$ and t^{p+1} . With the same objective, and since the $\{\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p\}$ occur only as arguments of $\mathbf{C}[\dots]$ and $\mathbf{a}[\dots]$, we will also abbreviate

$$\mathbf{C}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p] = \mathbf{C}[p] \quad (p \times p \text{ matrix}), \quad \mathbf{a}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p] = \mathbf{a}[p] \quad (p - \text{dim vector})$$

Prediction means finding (4.3), which upon inserting (4.12) becomes (in our new notation):

$$p(t^{p+1} | \boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^{p+1}, t^1, \dots, t^p) = \frac{p(t^1, \dots, t^{p+1} | \boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^{p+1})}{p(t^1, \dots, t^p | \boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p)}$$

$$= \frac{e^{-\frac{1}{2} \begin{pmatrix} t^1 - a_1[p+1] \\ \vdots \\ t^{p+1} - a_p[p+1] \end{pmatrix} \cdot \mathbf{C}^{-1}[p+1] \begin{pmatrix} t^1 - a_1[p+1] \\ \vdots \\ t^{p+1} - a_p[p+1] \end{pmatrix} + \frac{1}{2} \begin{pmatrix} t^1 - a_1[p] \\ \vdots \\ t^p - a_p[p] \end{pmatrix} \cdot \mathbf{C}^{-1}[p] \begin{pmatrix} t^1 - a_1[p] \\ \vdots \\ t^p - a_p[p] \end{pmatrix}}{\sqrt{2\pi \det \mathbf{C}[p+1] / \det \mathbf{C}[p]}} \quad (4.15)$$

Explicit Expressions for Predictions and Error Bars. The output probability distribution (4.15) for t^{p+1} , given the new input $\boldsymbol{\xi}^{p+1}$ and the data $D = \{(\boldsymbol{\xi}^1, t^1), \dots, (\boldsymbol{\xi}^p, t^p)\}$, must also be Gaussian, since it is a marginal distribution of a multivariate Gaussian distribution. This property allows us to concentrate in the exponent of (4.15) only on terms involving t^{p+1} ; the others simply provide an appropriate data-dependent normalization constant:

$$p(t^{p+1} | \boldsymbol{\xi}^{p+1}, D) \sim e^{-\frac{1}{2}(t^{p+1} - a_{p+1}[p+1])^2 \mathbf{C}^{-1}[p+1]_{p+1,p+1} - (t^{p+1} - a_{p+1}[p+1]) \sum_{\nu=1}^p \mathbf{C}^{-1}[p+1]_{p+1,\nu} (t^\nu - a_\nu[p+1])}$$

or, equivalently,

$$p(t^{p+1} | \boldsymbol{\xi}^{p+1}, D) \sim e^{-\frac{1}{2} \mathbf{C}^{-1}[p+1]_{p+1,p+1} \left[t^{p+1} - a_{p+1}[p+1] + \frac{\sum_{\nu=1}^p \mathbf{C}^{-1}[p+1]_{p+1,\nu} (t^\nu - a_\nu[p+1])}{\mathbf{C}^{-1}[p+1]_{p+1,p+1}} \right]^2}$$

We can now simply read off average and variance of the predicted output. Upon restoring our original notation the resulting prediction and its associated uncertainty read:

$$t^*(\boldsymbol{\xi}) = a_{p+1}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}] - \frac{\sum_{\nu=1}^p \mathbf{C}^{-1}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}]_{p+1,\nu} (t^\nu - a_\nu[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}])}{\mathbf{C}^{-1}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}]_{p+1,p+1}} \quad (4.16)$$

$$\Delta t^*(\boldsymbol{\xi}) = \frac{1}{\sqrt{\mathbf{C}^{-1}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}]_{p+1,p+1}}} \quad (4.17)$$

A non-zero choice for the $a_\mu[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^{p+1}]$ is equivalent to adding a fixed non-parametrized function of the joint inputs to the joint outputs of the underlying model. This is allowed mathematically, but requires rather specific knowledge of the function underlying the observed data, which is usually lacking. The two examples of Gaussian processes discussed earlier, RBF-type systems and two-layer feedforward networks with an infinitely large hidden layer, both had $a_\mu[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^{p+1}] = 0$ for all μ . In the latter cases one finds (4.16) being simplified to

$$t^*(\boldsymbol{\xi}) = -\frac{\sum_{\nu=1}^p \mathbf{C}^{-1}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}]_{p+1,\nu} t^\nu}{\mathbf{C}^{-1}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}]_{p+1,p+1}} \quad (4.18)$$

Given a data set $D = \{(\boldsymbol{\xi}^1, t^1), \dots, (\boldsymbol{\xi}^p, t^p)\}$ and a choice of model, i.e. of the functions $C_{\mu\nu}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^{p+1}]$ and $a[\boldsymbol{\xi}]$, we can now simply insert a novel input $\boldsymbol{\xi}$ into equations (4.16,4.17). More specifically, there is no (iterative) learning process; the computational effort is reduced to the inversion of a symmetric and positive definite $(p+1) \times (p+1)$ matrix.

Simplifications for On-Line Learning. We will next establish a relation between the matrices $\mathbf{C}^{-1}[p+1]$ and $\mathbf{C}^{-1}[p]$ (in our short-hand notation). In those cases where the data $(\boldsymbol{\xi}^\mu, t^\mu)$ arrive sequentially, it will reduce drastically the computational cost generated by the need to invert $\mathbf{C}[p+1]$. The idea is to exploit the property that, by definition, all entries of the covariance matrix $\mathbf{C}[p]$ also occur as the entries in the top-left $p \times p$ corner of $\mathbf{C}[p+1]$:

$$\mathbf{C}[p+1] = \begin{pmatrix} & & & k_1 \\ & \mathbf{C}[p] & & k_2 \\ & & & \vdots \\ k_1 & k_2 & \cdots & k_{p+1} \end{pmatrix} \quad k_\mu = C_{\mu,p+1}[p+1]$$

We write the inverse $\mathbf{C}^{-1}[p+1]$ to be calculated as

$$\mathbf{C}^{-1}[p+1] = \begin{pmatrix} & & & u_1 \\ & \mathbf{D} & & u_2 \\ & & & \vdots \\ u_1 & u_2 & \cdots & u_{p+1} \end{pmatrix}$$

with as yet unknown numbers $\{u_\mu\}$ and an as yet unknown (symmetric) $p \times p$ matrix \mathbf{D} . Working out the matrix condition $\mathbf{I} = \mathbf{C}^{-1}[p+1] \cdot \mathbf{C}[p+1]$ leads us to the expression

$$\mathbf{I} = \begin{pmatrix} (\mathbf{DC}[p])_{11} + u_1 k_1 & \cdots & (\mathbf{DC})_{1p} + u_1 k_p & \sum_{\mu=1}^p D_{1\mu} k_\mu + u_1 k_{p+1} \\ (\mathbf{DC}[p])_{21} + u_2 k_1 & \cdots & (\mathbf{DC})_{2p} + u_2 k_p & \sum_{\mu=1}^p D_{2\mu} k_\mu + u_2 k_{p+1} \\ \vdots & & & \vdots \\ (\mathbf{DC}[p])_{p1} + u_p k_1 & \cdots & (\mathbf{DC})_{pp} + u_p k_p & \sum_{\mu=1}^p D_{p\mu} k_\mu + u_p k_{p+1} \\ \sum_{\mu=1}^p C_{1\mu} u_\mu + u_{p+1} k_1 & \cdots & \sum_{\mu=1}^p C_{p\mu} u_\mu + u_{p+1} k_p & \sum_{\mu=1}^{p+1} u_\mu k_\mu \end{pmatrix}$$

Left and right-hand side are identical if and only if the following hold:

$$(\mathbf{DC}[p])_{\mu\nu} + u_\mu k_\nu = \delta_{\mu\nu} \quad \text{for all } \mu, \nu = 1, \dots, p$$

$$\sum_{\nu=1}^p D_{\mu\nu} k_\nu + u_\mu k_{p+1} = 0 \quad \sum_{\nu=1}^p C_{\mu\nu}[p] u_\nu + u_{p+1} k_\mu = 0 \quad \text{for all } \mu = 1, \dots, p$$

These equations are easily solved, giving

$$D_{\mu\nu} = (\mathbf{C}^{-1}[p])_{\mu\nu} + \frac{[\sum_{\rho \leq p} (\mathbf{C}^{-1}[p])_{\mu\rho} k_\rho] [\sum_{\rho \leq p} (\mathbf{C}^{-1}[p])_{\nu\rho} k_\rho]}{k_{p+1} - \sum_{\lambda\rho \leq p} k_\lambda (\mathbf{C}^{-1}[p])_{\lambda\rho} k_\rho} \quad (4.19)$$

$$u_\mu = \frac{\sum_{\rho \leq p} (\mathbf{C}^{-1}[p])_{\mu\rho} k_\rho}{k_{p+1} - \sum_{\lambda\rho \leq p} k_\lambda (\mathbf{C}^{-1}[p])_{\lambda\rho} k_\rho} \quad (\mu \leq p) \quad (4.20)$$

$$u_{p+1} = \frac{1}{k_{p+1} - \sum_{\lambda\rho \leq p} k_\lambda (\mathbf{C}^{-1}[p])_{\lambda\rho} k_\rho} \quad (4.21)$$

These results allow us, firstly, to replace expressions (4.16,4.17) by

$$t^*(\boldsymbol{\xi}) = a_{p+1}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}] + \sum_{\mu\nu \leq p} C_{p+1,\mu}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}] \mathbf{C}^{-1}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p]_{\mu\nu} (t^\nu - a_\nu[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}]) \quad (4.22)$$

$$\Delta t^*(\boldsymbol{\xi}) =$$

$$\sqrt{C_{p+1,p+1}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}] - \sum_{\mu\nu \leq p} C_{p+1,\mu}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}] (\mathbf{C}^{-1}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p])_{\mu\nu} C_{\nu,p+1}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p, \boldsymbol{\xi}]} \quad (4.23)$$

More importantly, however (since the new equations (4.22,4.23) are not simpler than the previous (4.16,4.17)), the identities (4.19,4.20,4.21) ensure that the inverse of the $p \times p$ matrix $\mathbf{C}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p]$ can be calculated *iteratively* from the inverse of the $(p-1) \times (p-1)$ matrix $\mathbf{C}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^{p-1}]$, i.e. upon adding the data points to our data set one by one.

4.3 The Choice of Covariance Matrix

Structure of Covariance Matrices. The form chosen for the covariance matrix, equivalent to the selection of data-generating model assumed responsible for the data, is crucial for the type of predictions that will be made by Gaussian processes. The two neural network examples of section 4.1, which reduced to Gaussian processes, were both found to be characterized not only by $a_\mu[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p] = 0$, but also by specific simple forms of the covariance matrices $C_{\mu\nu}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p]$, viz. (4.7) and (4.11):

$$\text{RBF : } C_{\mu\nu}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p] = \frac{1}{\beta} \delta_{\mu\nu} + \frac{1}{\alpha} \sum_{i=1}^L \phi_i(\boldsymbol{\xi}^\mu) \phi_i(\boldsymbol{\xi}^\nu) \quad (4.24)$$

$$\text{MLP : } C_{\mu\nu}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p] = \frac{1}{\beta} \delta_{\mu\nu} + \frac{1}{\alpha_w} \int \frac{dz}{(2\pi)^{N/2}} e^{-\frac{1}{2}z^2} g\left(\frac{z \cdot \boldsymbol{\xi}^\mu}{\sqrt{\alpha_J}}\right) g\left(\frac{z \cdot \boldsymbol{\xi}^\nu}{\sqrt{\alpha_J}}\right) \quad (4.25)$$

(in both cases $\boldsymbol{\xi} \in \mathfrak{R}^N$; the RBF-type system has L basis functions, the MLP has an infinitely large hidden layer). Both are of the following general form:

$$C_{\mu\nu}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p] = \frac{1}{\beta} \delta_{\mu\nu} + C[\boldsymbol{\xi}^\mu, \boldsymbol{\xi}^\nu] \quad (4.26)$$

It is in fact easy to show that *all* data-generating models of the general form (3.19) (not necessarily with Gaussian priors), i.e.

$$p(t|\boldsymbol{\xi}; \mathbf{w}) = e^{-\frac{1}{2}\beta[t-f(\boldsymbol{\xi}; \mathbf{w})]^2} \quad (4.27)$$

have the property that, if they reduce to a Gaussian process, then this process must have a covariance matrix of the form (4.26), and in addition $a_\mu[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p] = a[\boldsymbol{\xi}^\mu]$ for some function $a[\boldsymbol{\xi}]$ (note: the two examples of (4.7) and (4.11) both had $a[\boldsymbol{\xi}] = 0$). The joint output distribution for the models (4.27) is

$$p(t^1, \dots, t^p | \boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p) = \left[\frac{2\pi}{\beta} \right]^{\frac{p}{2}} \int d\mathbf{w} p(\mathbf{w}) e^{-\frac{1}{2}\beta \sum_{\mu=1}^p [t^\mu - f(\boldsymbol{\xi}^\mu; \mathbf{w})]^2}$$

Integration over the outputs then immediately leads to

$$a_\mu[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p] = a(\boldsymbol{\xi}^\mu) = \int d\mathbf{w} p(\mathbf{w}) f(\boldsymbol{\xi}^\mu; \mathbf{w}) \quad (4.28)$$

$$C_{\mu\nu}[\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^p] = \frac{1}{\beta} \delta_{\mu\nu} + \int d\mathbf{w} p(\mathbf{w}) [f(\boldsymbol{\xi}^\mu; \mathbf{w}) - a(\boldsymbol{\xi}^\mu)] [f(\boldsymbol{\xi}^\nu; \mathbf{w}) - a(\boldsymbol{\xi}^\nu)] \quad (4.29)$$

This is indeed of the form (4.26), as claimed.

Examples – Homogeneous Covariance Matrices. Let us now work out further the covariance matrix (4.7) of RBF networks, upon choosing Gaussian basis functions of width σ , i.e. $\phi_k(\boldsymbol{\xi}) = (2\pi\sigma^2)^{-N/2} e^{-\frac{1}{2}[\boldsymbol{\xi} - \mathbf{r}_k]^2/\sigma^2}$. It corresponds to the following function $C[\boldsymbol{\xi}, \boldsymbol{\xi}']$, in the representation (4.26):

$$C[\boldsymbol{\xi}, \boldsymbol{\xi}'] = \frac{1}{\alpha(2\pi\sigma^2)^N} \sum_{k=1}^L e^{-\frac{1}{2}[(\boldsymbol{\xi} - \mathbf{r}_k)^2 + (\boldsymbol{\xi}' - \mathbf{r}_k)^2]/\sigma^2}$$

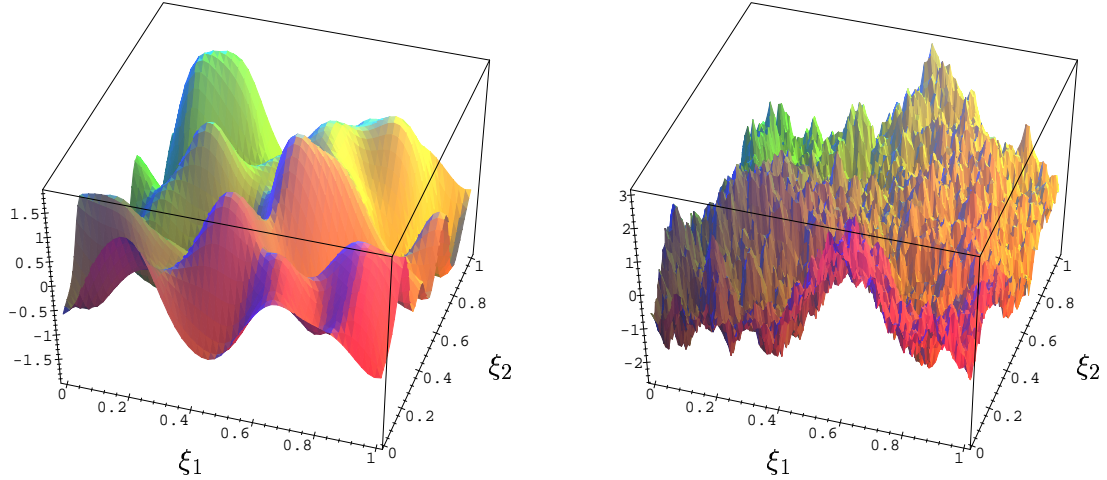


Figure 4.1: Examples of data surfaces, with $\boldsymbol{\xi}, \boldsymbol{\xi}' \in [0, 1]^2$, generated by Gaussian processes with homogeneous covariance matrices and $a[\boldsymbol{\xi}] = 0$. Left: $C[\boldsymbol{\xi}, \boldsymbol{\xi}'] = K_0 e^{-\frac{1}{2}|\boldsymbol{\xi} - \boldsymbol{\xi}'|^2/\sigma^2}$, with $K_0 = 10$ and $\sigma = 0.1$. Right: $C[\boldsymbol{\xi}, \boldsymbol{\xi}'] = K_0 e^{-|\boldsymbol{\xi} - \boldsymbol{\xi}'|/\sigma}$, with $K_0 = 10$ and $\sigma = 0.1$.

We now imagine having a very large number $L \rightarrow \infty$ of basis functions, distributed homogeneously over \mathbb{R}^N , and we re-scale our hyper-parameter α accordingly as $\alpha = \tilde{\alpha}/\Delta_L$. Here Δ_L denotes a small volume element in \mathbb{R}^N such that $\lim_{L \rightarrow \infty} \sum_{k=1}^L \Delta_L f[\mathbf{r}_k] = \int d\mathbf{r} f[\mathbf{r}]$ for any sufficiently smooth function f . We now find

$$\begin{aligned} \lim_{L \rightarrow \infty} C[\boldsymbol{\xi}, \boldsymbol{\xi}'] &= \frac{1}{\tilde{\alpha}(2\pi\sigma^2)^N} \int d\mathbf{r} e^{-\frac{1}{2}[(\boldsymbol{\xi} - \mathbf{r})^2 + (\boldsymbol{\xi}' - \mathbf{r})^2]/\sigma^2} \\ &= \frac{e^{-\frac{1}{2\sigma^2}[\boldsymbol{\xi}^2 + (\boldsymbol{\xi}')^2]}}{\tilde{\alpha}(2\pi\sigma^2)^N} \int d\mathbf{r} e^{-[\mathbf{r}^2 + \mathbf{r} \cdot (\boldsymbol{\xi} + \boldsymbol{\xi}')]/\sigma^2} = \frac{e^{-\frac{1}{4\sigma^2}(\boldsymbol{\xi} - \boldsymbol{\xi}')^2}}{\tilde{\alpha}(2\sigma\sqrt{\pi})^N} \end{aligned} \quad (4.30)$$

Working out the covariance matrix of example (4.11) for arbitrary $g[u]$ is hard. For data vectors $\boldsymbol{\xi}$ close to the origin (alternatively: for large values of the hyper-parameter α_J), however, we may expand typical choices such as $g[u] = \tanh[u]$ as $g[u] = u + \mathcal{O}(u^3)$ and find

$$\begin{aligned} C[\boldsymbol{\xi}, \boldsymbol{\xi}'] &= \frac{1}{\alpha_w \alpha_J} \int \frac{dz}{(2\pi)^{N/2}} e^{-\frac{1}{2}z^2} (z \cdot \boldsymbol{\xi})(z \cdot \boldsymbol{\xi}') + \mathcal{O}(|\boldsymbol{\xi}||\boldsymbol{\xi}'|^3, |\boldsymbol{\xi}|^3|\boldsymbol{\xi}'|) \\ &= \frac{\boldsymbol{\xi} \cdot \boldsymbol{\xi}'}{\alpha_w \alpha_J} + \mathcal{O}(|\boldsymbol{\xi}||\boldsymbol{\xi}'|^3, |\boldsymbol{\xi}|^3|\boldsymbol{\xi}'|) \end{aligned} \quad (4.31)$$

(truncation after the first term is obviously exact only for $g[u] = u$).

Homogeneous covariance matrices are the subclass of (4.26) where the function $C[\boldsymbol{\xi}, \boldsymbol{\xi}']$ is of the form $C[\boldsymbol{\xi}, \boldsymbol{\xi}'] = C[|\boldsymbol{\xi} - \boldsymbol{\xi}'|]$ (like, for instance, the RBF expression (4.30), but unlike the MLP expression (4.31)). Here the output statistics are fully isotropic in input space (i.e. invariant under arbitrary translations and/or rotations of both $\boldsymbol{\xi}$ and $\boldsymbol{\xi}'$). Within this subclass the only remaining freedom we have regarding covariances is the choice of a scalar function $C[u]$ (with $u \geq 0$); however, still much variation is found to be possible in the types

of data that can be generated and predicted. Figure 4.1 shows examples of data surfaces, with $\boldsymbol{\xi}, \boldsymbol{\xi}' \in [0, 1]^2$, generated by Gaussian processes with the following homogeneous covariance matrices:

$$\text{'smooth' functions :} \quad C[\boldsymbol{\xi}, \boldsymbol{\xi}'] = K_0 e^{-\frac{1}{2}|\boldsymbol{\xi} - \boldsymbol{\xi}'|^2 / \sigma^2} \quad (4.32)$$

$$\text{'rough' functions :} \quad C[\boldsymbol{\xi}, \boldsymbol{\xi}'] = K_0 e^{-|\boldsymbol{\xi} - \boldsymbol{\xi}'| / \sigma} \quad (4.33)$$

The qualitative difference between the two can be appreciated upon realizing that smoothness is defined in terms of the behaviour of the kernels $C[\boldsymbol{\xi}, \boldsymbol{\xi}']$ for $\boldsymbol{\xi}' \rightarrow \boldsymbol{\xi}$. Let us choose $\boldsymbol{\xi}' - \boldsymbol{\xi} = \boldsymbol{\epsilon}$, with $|\boldsymbol{\epsilon}| \ll 1$. For the above two examples one finds, respectively:

$$e^{-\frac{1}{2}|\boldsymbol{\xi} - \boldsymbol{\xi}'|^2 / \sigma^2} = 1 - \frac{\boldsymbol{\epsilon}^2}{2\sigma^2} + \mathcal{O}(|\boldsymbol{\epsilon}|^4) \quad e^{-|\boldsymbol{\xi} - \boldsymbol{\xi}'| / \sigma} = 1 - \frac{|\boldsymbol{\epsilon}|}{\sigma} + \mathcal{O}(|\boldsymbol{\epsilon}|^2)$$

Hence the second kernel indeed describes much faster de-correlation of outputs for small differences in the inputs. In a similar fashion one could define a more general family of covariance matrices, $C[\boldsymbol{\xi}, \boldsymbol{\xi}'] = K_0 e^{-|\boldsymbol{\xi} - \boldsymbol{\xi}'|^\gamma / \sigma^\gamma}$, in which the degree of smoothness is controlled monotonically by a parameter $\gamma > 0$.

Chapter 5

Support Vector Machines for Binary Classification

Support vector machines carry out binary classifications. They involve a preprocessing stage where non-linearly separable data are converted (hopefully) into linearly separable ones. First, however, we will have to return to the properties of linearly separable problems (for reasons which will become clear); the pre-processing will be introduced later.

5.1 Optimal Separating Plane for Linearly Separable Tasks

Linear Separations & Stability Parameters. Consider a linearly separable binary classification task, with a set of data of the usual form $D = \{(\boldsymbol{\xi}^1, t^1), \dots, (\boldsymbol{\xi}^p, t^p)\}$, where $\boldsymbol{\xi} \in \mathfrak{R}^N$ and $t^\mu \in \{-1, 1\}$. If this problem is linearly separable, we know that

$$(\exists \mathbf{w} \in \mathfrak{R}^N)(\exists w_0 \in \mathfrak{R}) : \quad t^\mu = \text{sgn}[\mathbf{w} \cdot \boldsymbol{\xi}^\mu + w_0] \quad \text{for all } (\boldsymbol{\xi}^\mu, t^\mu) \in D \quad (5.1)$$

We even have algorithms which are guaranteed to converge towards a solution (\mathbf{w}, w_0) (a ‘separating plane’) of the type (5.1), e.g. the so-called perceptron or AdaTron rules. However, it will be clear that in general there will exist an infinite number of separating planes (\mathbf{w}, w_0) that will meet the requirements (5.1). See e.g. figure 5.1 (left panel) for $N = 2$; here a separating plane is any line separating all points $\boldsymbol{\xi}^\mu$ with $t^\mu = 1$ (marked +) from those with $t^\mu = -1$ (marked \times).

A natural question to ask is whether one can quantify the quality of the various solutions of (5.1). This can be done on the basis of generalization performance: a good plane (\mathbf{w}, w_0) is one which keeps a safe distance (if possible) from the data points, so that detrimental effects of noise (which might change the locations of these points slightly) are minimal. Hence we wish to measure not only whether points are correctly classified by a candidate separating plane, but also the *distance* of such a plane to the data points which it aims to separate. Both are given by the so-called *stability parameters* γ_μ (there is one for every data point):

$$\gamma_\mu(\mathbf{w}, w_0) = t^\mu(\mathbf{w} \cdot \boldsymbol{\xi}^\mu + w_0)/|\mathbf{w}| \quad (5.2)$$

They have the properties:

- The plane (\mathbf{w}, w_0) correctly classifies point $\boldsymbol{\xi}^\mu$ if and only if: $\gamma_\mu(\mathbf{w}, w_0) > 0$
- The distance of the plane (\mathbf{w}, w_0) to point $\boldsymbol{\xi}^\mu$ is: $|\gamma_\mu(\mathbf{w}, w_0)|$

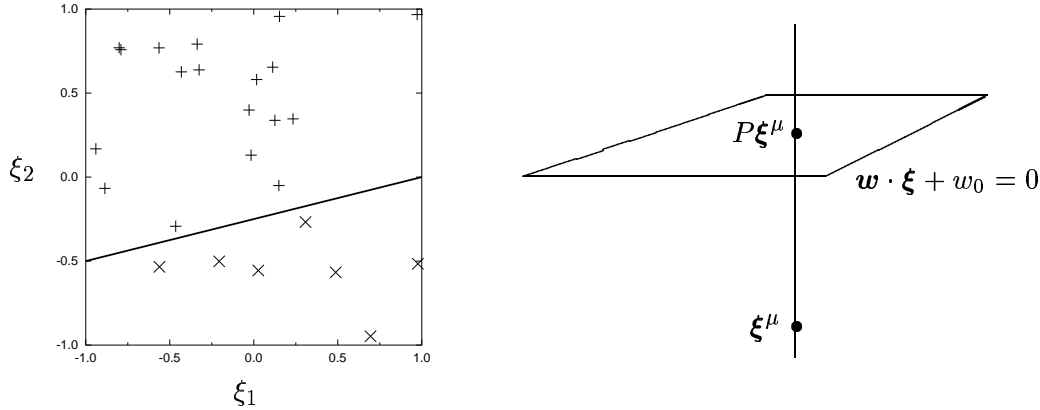


Figure 5.1: Left: illustration of linearly separable tasks for $N = 2$. A separating plane is any line $\mathbf{w} \cdot \boldsymbol{\xi} + w_0 = 0$ separating all points $\boldsymbol{\xi}^\mu$ with $t^\mu = 1$ (marked $+$) from those with $t^\mu = -1$ (marked \times). Note that there will generally be infinitely many such planes. Right: measuring separation quality on the basis of the distance of data points $\boldsymbol{\xi}^\mu$ from the separating plane (as calculated via projection); see the main text.

The first property is immediately obvious from (5.1). Demonstrating the second property requires only simple high-school geometry. We first construct a parametrization of the line in \mathbb{R}^N which goes through $\boldsymbol{\xi}^\mu$ and is orthogonal to our plane (see figure 5.1, right panel): $\boldsymbol{\xi}(\rho) = \boldsymbol{\xi}^\mu + \rho \mathbf{w}$ ($\rho \in \mathbb{R}$). Insertion into the plane's equation $\mathbf{w} \cdot \boldsymbol{\xi} + w_0 = 0$ gives the solution $\rho = -(\mathbf{w} \cdot \boldsymbol{\xi} + w_0)/\mathbf{w}^2$. Hence the projection $P\boldsymbol{\xi}^\mu$ of $\boldsymbol{\xi}^\mu$ onto the plane (\mathbf{w}, w_0) is given by

$$P\boldsymbol{\xi}^\mu = \boldsymbol{\xi}^\mu - \mathbf{w} \left[\frac{\mathbf{w} \cdot \boldsymbol{\xi}^\mu + w_0}{\mathbf{w}^2} \right]$$

And thus the desired distance $|\boldsymbol{\xi}^\mu - P\boldsymbol{\xi}^\mu|$ indeed equals (using $|t^\mu| = 1$ for all μ):

$$|\boldsymbol{\xi}^\mu - P\boldsymbol{\xi}^\mu| = \left| \mathbf{w} \left[\frac{\mathbf{w} \cdot \boldsymbol{\xi}^\mu + w_0}{\mathbf{w}^2} \right] \right| = |\gamma_\mu(\mathbf{w}, w_0)|$$

As claimed. As a consequence of the properties of the stability parameters we may now sharpen our requirements of a separating plane:

$$\begin{aligned} \text{all data points separated correctly \& } \\ \text{all data points at distance } \geq \Delta \text{ from the plane} \end{aligned} \iff \gamma_\mu(\mathbf{w}, w_0) \geq \Delta > 0 \text{ for all } \mu \quad (5.3)$$

There will be an upper limit to the values of the distance $\Delta > 0$ which are compatible with correct separation (and it is a non-trivial task to determine this limit), since clearly

$$\Delta \leq \frac{1}{2} \Delta^\pm \quad \Delta^\pm = \min_{\mu \in I_+, \nu \in I_-} |\boldsymbol{\xi}^\mu - \boldsymbol{\xi}^\nu| \quad I_\pm = \{\mu \in \{1, \dots, p\} \mid t^\mu = \pm 1\}$$

However, if we choose a value of Δ such that a solution (\mathbf{w}, w_0) of (5.3) exists, then there is an algorithm which is guaranteed to achieve the objective (a straightforward generalization of the classic Perceptron learning rule, which is indeed recovered for $\Delta \rightarrow 0$):

Generalized Perceptron Learning Rule:

1. drawn at random a $\mu \in \{1, \dots, p\}$
2. calculate $\gamma_\mu(\mathbf{w}, w_0)$ as given in (5.2)
3. if $\gamma_\mu(\mathbf{w}, w_0) < \Delta$, modify parameters: $\mathbf{w} \rightarrow \mathbf{w} + t^\mu \boldsymbol{\xi}^\mu$
 $w_0 \rightarrow w_0 + t^\mu$
4. return to 1.

Generalized Perceptron Convergence Theorem (E. Gardner):

If $(\exists \mathbf{w}^* \in \mathfrak{R}^N)(\exists w_0^* \in \mathfrak{R}) : \gamma_\mu(\mathbf{w}^*, w_0^*) > \Delta$ for all $\mu \in \{1, \dots, p\}$, then the above algorithm will converge towards a solution (\mathbf{w}, w_0) with $\gamma_\mu(\mathbf{w}, w_0) \geq \Delta$ in a finite number of iteration steps.

The proof will not be given here. For $w_0 = 0$ it is a simple generalization of the original perceptron convergence proof. For $w_0 \neq 0$, however, the situation is more complicated; the reason is that the threshold w_0 does not occur in the denominator of the stability parameters (5.2), so that it cannot simply be absorbed into the classic proof by adding a ‘dummy’ extra component $\xi_0^\mu = 1$ to all input vectors $\boldsymbol{\xi}^\mu$.

The Optimal Separating Plane. We are now in a position to define the optimal separating plane (\mathbf{w}^*, w_0^*) as the plane for which the *smallest* of the $\{\gamma_\mu(\mathbf{w}, w_0)\}$ (i.e. the stability parameter of the data point which is most in danger of being mis-classified when data are subject to noise) is as large as possible:

$$\text{optimal separating plane } (\mathbf{w}^*, w_0^*) = \arg \max_{(\mathbf{w}, w_0)} \left[\min_{\mu} \gamma_\mu(\mathbf{w}, w_0) \right] \quad (5.4)$$

Note that:

- If the problem is linearly separable, definition (5.4) will produce a separating plane with the largest distance to the closest data point. If the problem is not linearly separable, definition (5.4) will produce a plane which does not separate the data (since that is impossible), but will have the minimum distance to the plane of those mis-classified points which are furthest away from the plane.
- By construction, there must be at least two μ such that $\gamma_\mu(\mathbf{w}^*, w_0^*) = \min_{\rho} \gamma_\rho(\mathbf{w}^*, w_0^*)$ for the optimal separating plane (5.4) (see e.g. figure 5.1).

Finding the optimal plane has now been reduced to a relatively simple optimization exercise (5.4). The only problem with it is that its solution (\mathbf{w}^*, w_0^*) is still not unique, since there remains the trivial degree of freedom relating to overall scaling: if we multiply $(\mathbf{w}, w_0) \rightarrow (\rho \mathbf{w}, \rho w_0)$ with $\rho > 0$ we simply get the same plane. To eliminate this freedom we nail down ρ by insisting that

$$\min_{\mu} t^\mu (\mathbf{w} \cdot \boldsymbol{\xi}^\mu + w_0) = 1$$

(given the assumption of linear separability this is always possible). This implies that $\min_{\mu} \gamma_\mu(\mathbf{w}, w_0) = 1/|\mathbf{w}|$, and that $\max_{(\mathbf{w}, w_0)} [\min_{\mu} \gamma_\mu(\mathbf{w}, w_0)] = 1/\min_{(\mathbf{w}, w_0)} |\mathbf{w}|$. Hence we have now converted the initial optimization problem (5.4) into

$$\text{find : } \min_{(\mathbf{w}, w_0)} \frac{1}{2} \mathbf{w}^2 \quad \text{under constraints} \quad (\forall \mu) : t^\mu (\mathbf{w} \cdot \boldsymbol{\xi}^\mu + w_0) \geq 1 \quad (5.5)$$

5.2 Representation in terms of Support Vectors

Solution of the Optimization Problem. We now solve the constrained optimization problem (5.5) using the standard method of Lagrange multipliers. This leads to the following coupled equations:

$$\frac{\partial}{\partial w_i} \frac{1}{2} \mathbf{w}^2 = \sum_{\mu=1}^p \lambda_{\mu}(\mathbf{w}, w_0) \frac{\partial}{\partial w_i} [t^{\mu}(\mathbf{w} \cdot \boldsymbol{\xi}^{\mu} + w_0) - 1] \quad i = 0, \dots, N$$

$$\begin{aligned} \text{with :} \quad \lambda_{\mu}(\mathbf{w}, w_0) &= 0 & \text{if extremum has } t_{\mu}(\mathbf{w} \cdot \boldsymbol{\xi}^{\mu} + w_0) &> 1 \\ \lambda_{\mu}(\mathbf{w}, w_0) &\neq 0 & \text{if extremum has } t_{\mu}(\mathbf{w} \cdot \boldsymbol{\xi}^{\mu} + w_0) &= 1 \end{aligned}$$

We define the index set S (dependent on \mathbf{w} and w_0) as follows:

$$\begin{aligned} \mu \in S : \quad t^{\mu}(\mathbf{w} \cdot \boldsymbol{\xi}^{\mu} + w_0) &= 1 & (\text{so } \lambda_{\mu}(\mathbf{w}, w_0) &\neq 0) \\ \mu \notin S : \quad t^{\mu}(\mathbf{w} \cdot \boldsymbol{\xi}^{\mu} + w_0) &> 1 & (\text{so } \lambda_{\mu}(\mathbf{w}, w_0) &= 0) \end{aligned} \quad (5.6)$$

with $|S| \leq p$ elements. Without danger of confusion we will henceforth drop the arguments of the Lagrange parameters and write them simply as λ_{μ} . Working out the above partial derivatives gives

$$\mathbf{w} = \sum_{\mu \in S} \lambda_{\mu} t^{\mu} \boldsymbol{\xi}^{\mu} \quad \sum_{\mu \in S} \lambda_{\mu} t^{\mu} = 0 \quad (5.7)$$

$$\mu \in S : \quad t_{\mu} \left[\sum_{\nu \in S} \lambda_{\nu} t^{\nu} (\boldsymbol{\xi}^{\mu} \cdot \boldsymbol{\xi}^{\nu}) + w_0 \right] = 1 \quad \mu \notin S : \quad t_{\mu} \left[\sum_{\nu \in S} \lambda_{\nu} t^{\nu} (\boldsymbol{\xi}^{\mu} \cdot \boldsymbol{\xi}^{\nu}) + w_0 \right] > 1 \quad (5.8)$$

We will now show that, given knowledge of the index set S , the solution of (5.7,5.8) follows directly. In doing so we will need the following (symmetric and non-negative definite) $|S| \times |S|$ matrix \mathbf{C} :

$$C_{\mu\nu} = t^{\mu}(\boldsymbol{\xi}^{\mu} \cdot \boldsymbol{\xi}^{\nu}) t^{\nu} \quad \text{for } \mu, \nu \in S \quad (5.9)$$

Firstly, after insertion of the first expression of (5.7) (to eliminate \mathbf{w}) the first equation of (5.8) allows us to write the non-zero Lagrange parameters as

$$\mu \in S : \quad \lambda_{\mu} = \sum_{\nu \in S} (\mathbf{C}^{-1})_{\mu\nu} (1 - w_0 t^{\nu}) \quad (5.10)$$

Inserting (5.10) into the second equation of (5.7) gives us the threshold w_0 :

$$w_0 = \frac{\sum_{\mu\nu \in S} t^{\mu} (\mathbf{C}^{-1})_{\mu\nu}}{\sum_{\mu\nu \in S} t^{\mu} (\mathbf{C}^{-1})_{\mu\nu} t^{\nu}} \quad (5.11)$$

Our problem of finding the solution (\mathbf{w}^*, w_0^*) of (5.7,5.8), i.e. the optimal separating plane, has thus been reduced to finding the index set S^1 . We note in the above equations that only those inputs play a role which correspond to indices $\mu \in S$. This then leads us to the definition of *Support Vectors*:

¹If the matrix \mathbf{C} has zero eigenvalues (it can never have negative ones) we must replace \mathbf{C}^{-1} in the above equations by the pseudo-inverse of \mathbf{C} .

- Support vectors are those input vectors ξ^μ in the training set which, in the construction of the optimal separating plane as given above, correspond to indices $\mu \in S$.
- Support vectors are those input vectors which are closest to the separating plane, since the distance for pattern μ equals $\gamma_\mu = t^\mu(\mathbf{w} \cdot \xi^\mu + w_0)/|\mathbf{w}|$ and (5.8) tells us that

$$\mu \in S: \quad \gamma_\mu = 1/|\mathbf{w}| \qquad \mu \notin S: \quad \gamma_\mu > 1/|\mathbf{w}|$$

- Only support vectors and their correlations occur in the construction of the optimal separating plane. All other inputs are classified correctly once the support vectors are.

Let us at this stage derive two further properties of the solution, for later use. Firstly:

$$\sum_{\mu \in S} \lambda_\mu \geq 0 \tag{5.12}$$

Proof: we use the fact that \mathbf{C} is non-negative definite and insert the first equation of (5.8), i.e. $\sum_{\nu \in S} \lambda_\nu t^\nu (\xi^\mu \cdot \xi^\nu) t^\nu = 1 - w_0 t^\mu$, followed by the second equation of (5.7), giving

$$0 \leq \sum_{\mu\nu \in S} \lambda_\mu C_{\mu\nu} \lambda_\nu = \sum_{\mu \in S} \lambda_\mu (1 - w_0 t^\mu) = \sum_{\mu \in S} \lambda_\mu$$

Secondly, given the parameter vector \mathbf{w} , the threshold w_0 can be written as

$$w_0 = -\frac{1}{2} \mathbf{w} \cdot (\xi^\mu + \xi^\nu) \quad \text{for any } \mu, \nu \in S \quad \text{with } t^\mu = 1, t^\nu = -1 \tag{5.13}$$

Proof: this follows immediately from the definitions (5.6) in combination with $t^\mu \in \{-1, 1\}$, according to which $\mathbf{w} \cdot \xi^\mu = t^\mu - w_0 =$ for $\mu \in S$.

Transformation of the Optimization Problem. Knowing the structure of the solution, we can transform the problem to a set of equations in the space of the Lagrange multipliers $\{\lambda_\mu\}$, which makes sense particularly if $p < N$. Insertion of $\mathbf{w} = \sum_{\mu=1}^p \lambda_\mu t^\mu \xi^\mu$, following (5.7), into the original problem (5.5) gives with $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_p)$:

$$\text{find: } \min_{(\boldsymbol{\lambda}, w_0)} \frac{1}{2} \sum_{\mu\nu=1}^p \lambda_\mu t^\mu (\xi^\mu \cdot \xi^\nu) t^\nu \lambda_\nu \quad \text{under constraints } \forall \mu: t^\mu \sum_{\nu=1}^p (\xi^\mu \cdot \xi^\nu) t^\nu \lambda_\nu \geq 1 - w_0 t^\mu \tag{5.14}$$

Support vectors correspond to those μ for which the last inequality becomes an equality. In preparation of a statement to be proven below, we note that the equations from which to solve the solution of (5.5) and of (5.14) can be summarized as

$$W_\mu = \sum_{\nu} t^\mu (\xi^\mu \cdot \xi^\nu) t^\nu \lambda_\nu + w_0 t^\mu - 1, \tag{5.15}$$

$$\sum_{\mu=1}^p \lambda_\mu t^\mu = 0, \quad \forall \mu: \quad \{\lambda_\mu = 0, W_\mu > 0\} \quad \text{or} \quad \{\lambda_\mu \neq 0, W_\mu = 0\} \tag{5.16}$$

(these are indeed $2p + 1$ equations for the $2p + 1$ unknowns $\{W_\mu, \lambda_\mu\}$ and w_0). In the case of multiple solutions we must select the one for which the quantity $\frac{1}{2} \sum_{\mu\nu=1}^p \lambda_\mu t^\mu (\xi^\mu \cdot \xi^\nu) t^\nu \lambda_\nu$ is minimal.

A final transformation brings our problem into the form most commonly encountered in textbooks. The statement is that the solution of (5.14) is identical to that of the following similar constrained minimization problem

$$\text{find: } \min_{\lambda} \left\{ \frac{1}{2} \sum_{\mu\nu=1}^p \lambda_{\mu} t^{\mu} (\boldsymbol{\xi}^{\mu} \cdot \boldsymbol{\xi}^{\nu}) t^{\nu} \lambda_{\nu} - \sum_{\mu} \lambda_{\mu} \right\} \quad \text{constraints: } \sum_{\mu=1}^p \lambda_{\mu} t^{\mu} = 0, \quad \forall \mu: \lambda_{\mu} \geq 0 \quad (5.17)$$

To appreciate the origin of this statement (we will not give the full elaborate proof here), let us work out the Lagrange equations for (5.17), denoting the Lagrange parameters of the inequality constraints $\lambda_{\mu} \geq 0$ as W_{μ} , and the Lagrange parameter of the constraint $\sum_{\mu=1}^p \lambda_{\mu} t^{\mu} = 0$ as $-W_0$:

$$W_{\mu} = \sum_{\nu} t_{\nu} (\boldsymbol{\xi}^{\mu} \cdot \boldsymbol{\xi}^{\nu}) t^{\nu} \lambda_{\nu} + W_0 t^{\mu} - 1 \quad (5.18)$$

$$\sum_{\mu=1}^p \lambda_{\mu} t^{\mu} = 0, \quad \forall \mu: \quad \{\lambda_{\mu} = 0, W_{\mu} \neq 0\} \quad \text{or} \quad \{\lambda_{\mu} > 0, W_{\mu} = 0\} \quad (5.19)$$

Comparison with (5.15,5.16) shows that the two sets of equations and requirements are identical if the following is true:

$$\text{Solution of (5.14): } \lambda_{\mu} \geq 0 \quad \forall \mu \quad \quad \text{Solution of (5.17): } W_{\mu} \geq 0 \quad \forall \mu$$

Both statements are very similar, and can be shown to be true by virtue of the fact that we are only looking for solutions of the Lagrange equations which are *minima* of the quadratic surface to be extremized (the sign of a Lagrange parameter refers to the direction of the associated constraining force, relative to the gradient of the surface). Since the solution of a quadratic minimization problem with linear inequality constraints can also be shown to have a *unique* solution, the two minimization problems are indeed found to be equivalent².

A Simple Example. Let us assume, for simplicity, that the input vectors in our data set D are orthogonal and normalized, i.e. $\boldsymbol{\xi}^{\mu} \cdot \boldsymbol{\xi}^{\nu} = \delta_{\mu\nu}$ (which implies also that $p \leq N$). We aim to calculate the optimal separating plane (\mathbf{w}^*, w_0^*) from

$$\mathbf{w} = \sum_{\mu \in S} \lambda_{\mu} t^{\mu} \boldsymbol{\xi}^{\mu} \quad w_0 = \frac{\sum_{\mu\nu \in S} t^{\mu} (\mathbf{C}^{-1})_{\mu\nu}}{\sum_{\mu\nu \in S} t^{\mu} (\mathbf{C}^{-1})_{\mu\nu} t^{\nu}} \quad (5.20)$$

with the matrix \mathbf{C} as given in (5.9), and with the index set S giving the labels μ of the support vectors (i.e. those with $\lambda_{\mu} > 0$). We note that in the present example $C_{\mu\nu} = \delta_{\mu\nu}$, so that the second equation of (5.20) becomes

$$w_0 = |S|^{-1} \sum_{\mu \in S} t^{\mu} \quad -1 \leq w_0 \leq 1$$

We calculate the $\{\lambda_{\mu}\}$ via the original equations (5.10), which reduces to

$$\mu \in S: \quad \lambda_{\mu} = 1 - w_0 t^{\mu}$$

²See the specialized literature on convex optimization for full and further details.

Insertion into (5.8) subsequently gives

$$\mu \in S: \quad 1 = 1 \qquad \mu \notin S: \quad t_\mu w_0 > 1$$

We conclude that, in view of $|w_0| \leq 1$, the only possible solution is the one where $S = \{1, \dots, p\}$, i.e. all input vectors are support vectors. Hence

$$\mathbf{w}^* = \sum_{\mu=1}^p [t^\mu - w_0] \boldsymbol{\xi}^\mu \qquad w_0^* = \frac{1}{p} \sum_{\mu=1}^p t^\mu \qquad (5.21)$$

(which is known as the ‘biased Hebb rule’). The corresponding stability parameters $\gamma_\mu = |\mathbf{w}|^{-1}$ (since $\mu \in S$ for all μ) are

$$\gamma_\mu = \frac{1}{\sqrt{\sum_{\mu=1}^p [t^\mu - w_0]^2}}$$

5.3 Pre-Processing, SVM Kernels

Pre-Processing of Non-Linearly-Separable Data. Let us now turn to problems which are not linearly separable, but could possibly be made so via suitable pre-processing. This implies that we will no longer view the $\boldsymbol{\xi}$ as the input vectors, but see the $\boldsymbol{\xi}$ as the outcome of the pre-processing of true input vectors $\mathbf{x} \in \mathfrak{R}^M$. If we write the combined action of our classification machine as $S(\mathbf{x}) \in \{-1, 1\}$ we have

$$\boldsymbol{\xi} = \boldsymbol{\phi}(\mathbf{x}) : \qquad S(\mathbf{x}) = \text{sgn}[\mathbf{w} \cdot \boldsymbol{\phi}(\mathbf{x}) + w_0] \qquad (5.22)$$

with $\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_N(\mathbf{x}))$ and usually with $N > M$. The data set is now $D = \{(\mathbf{x}^1, t^1), \dots, (\mathbf{x}^p, t^p)\}$. The set of pre-processing operations $\phi_i(\mathbf{x})$ is typically chosen to contain non-linear functions (otherwise there would be no point in pre-processing), such as

$$\begin{aligned} \phi_i(\mathbf{x}) &= x_i && \text{for } i = 1, \dots, M \\ \phi_i(\mathbf{x}) &= x_i^2 && \text{for } i = M + 1, \dots, 2M \\ \phi_{2M+1}(\mathbf{x}) &= x_1 x_2 && \phi_{2M+2}(\mathbf{x}) = x_1 x_3 \quad \phi_{2M+3}(\mathbf{x}) = x_2 x_3 \quad \text{etc} \end{aligned}$$

Our previous results regarding the optimal separating plane (\mathbf{w}^*, w_0^*) can be immediately applied via the replacement $\boldsymbol{\xi}^\mu \rightarrow \boldsymbol{\phi}(\mathbf{x}^\mu)$, which leads to

$$\mathbf{w}^* = \sum_{\mu=1}^p \lambda_\mu t^\mu \boldsymbol{\phi}(\mathbf{x}^\mu) \qquad w_0^* = -\frac{1}{2} \mathbf{w}^* \cdot [\boldsymbol{\phi}(\mathbf{x}^+) + \boldsymbol{\phi}(\mathbf{x}^-)] \qquad (5.23)$$

Here the support vectors are those \mathbf{x}^μ for which $\lambda_\mu = 0$, \mathbf{x}^\pm are (any) support vectors with $t^\mu = \pm 1$, and the $\{\lambda_\mu\}$ are determined by solving:

$$\min_{\boldsymbol{\lambda}} \left\{ \frac{1}{2} \sum_{\mu\nu=1}^p \lambda_\mu \lambda_\nu t^\mu t^\nu \boldsymbol{\phi}(\mathbf{x}^\mu) \cdot \boldsymbol{\phi}(\mathbf{x}^\nu) - \sum_{\mu} \lambda_\mu \right\} \quad \text{constraints:} \quad \sum_{\mu=1}^p \lambda_\mu t^\mu = 0, \quad \forall \mu: \lambda_\mu \geq 0 \qquad (5.24)$$

Insertion of the above expressions (5.23) into the machine’s operation rule (5.22) gives the final classification recipe

$$S(\mathbf{x}) = \operatorname{sgn} \left[\sum_{\mu=1}^p \lambda_{\mu} t^{\mu} \left\{ \phi(\mathbf{x}^{\mu}) \cdot \phi(\mathbf{x}) - \frac{1}{2} \phi(\mathbf{x}^{\mu}) \cdot \phi(\mathbf{x}^{+}) - \frac{1}{2} \phi(\mathbf{x}^{\mu}) \cdot \phi(\mathbf{x}^{-}) \right\} \right] \quad (5.25)$$

At this stage one makes the important observation that neither in (5.25), nor in the equations (5.24) from which to extract the $\{\lambda_{\mu}\}$ do we retain any ‘bare’ occurrence of the pre-processing functions $\phi_i(\mathbf{x})$. Instead we consistently encounter only the symmetric form

$$K(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^N \phi_i(\mathbf{x}) \phi_i(\mathbf{y}) \quad (5.26)$$

More specifically, in terms of (5.26) our full equations become

$$S(\mathbf{x}) = \operatorname{sgn} \left[\sum_{\mu=1}^p \lambda_{\mu} t^{\mu} \left\{ K(\mathbf{x}^{\mu}, \mathbf{x}) - \frac{1}{2} K(\mathbf{x}^{\mu}, \mathbf{x}^{+}) - \frac{1}{2} K(\mathbf{x}^{\mu}, \mathbf{x}^{-}) \right\} \right] \quad (5.27)$$

$$\min_{\lambda} \left\{ \frac{1}{2} \sum_{\mu\nu=1}^p \lambda_{\mu} \lambda_{\nu} t^{\mu} t^{\nu} K(\mathbf{x}^{\mu}, \mathbf{x}^{\nu}) - \sum_{\mu} \lambda_{\mu} \right\} \quad \text{constraints:} \quad \sum_{\mu=1}^p \lambda_{\mu} t^{\mu} = 0, \quad \forall \mu: \lambda_{\mu} \geq 0 \quad (5.28)$$

We can henceforth work directly with (5.27,5.28), and call the function $K(\mathbf{x}, \mathbf{y})$ in these equations the SVM Kernel. The decision boundary in input space of our machine is defined as the set of all $\mathbf{x} \in \mathfrak{R}^M$ for which

$$\sum_{\mu=1}^p \lambda_{\mu} t^{\mu} \left\{ K(\mathbf{x}^{\mu}, \mathbf{x}) - \frac{1}{2} K(\mathbf{x}^{\mu}, \mathbf{x}^{+}) - \frac{1}{2} K(\mathbf{x}^{\mu}, \mathbf{x}^{-}) \right\} = 0 \quad (5.29)$$

We can even generalize our classification machine to more general non-negative symmetric kernels $K(\mathbf{x}, \mathbf{y})$ than those of the form (5.26).

SVM Kernels – Properties and Examples. We have seen that, at least for SVM kernels of the form (5.26), making a choice for the $K(\mathbf{x}, \mathbf{y})$ implies choosing a specific type of pre-processing³. Not all non-negative and symmetric kernels $K(\mathbf{x}, \mathbf{y})$ can be written in the form (5.26), however; the conditions for this to be possible are given in Mercer’s Theorem (which we will not prove here):

A symmetric kernel $K(\mathbf{x}, \mathbf{y})$ has an expansion of the form $K(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{\infty} \phi_i(\mathbf{x}) \phi_i(\mathbf{y})$ in the function space $L^2(\mathfrak{R}^M)$ if and only if

$$0 < \int d\mathbf{x} d\mathbf{y} g(\mathbf{x}) K(\mathbf{x}, \mathbf{y}) g(\mathbf{y}) < \infty \quad \text{for all functions } g \in L^2(\mathfrak{R}^M)$$

Apparently, choosing an acceptable kernel directly, as opposed to via (5.26), will generally be equivalent to working with an *infinite* number of pre-processing functions. To see this more explicitly, we inspect a couple of kernel examples:

³Note: choosing bi-linear kernels $K(\mathbf{x}, \mathbf{y}) = \sum_{ij} x_i A_{ij} y_j$ is pointless: we would then assume linear separability of the original data, and would be better of classifying directly via the optimal plane in input space.

- Kernels generating *pre-processing via polynomials*: $K(\mathbf{x}, \mathbf{y}) = [\mathbf{x} \cdot \mathbf{y} + 1]^d$

Expansion of this type of kernel gives:

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= 1 + d(\mathbf{x} \cdot \mathbf{y}) + \frac{1}{2}d(d-1)(\mathbf{x} \cdot \mathbf{y})^2 + \frac{1}{6}d(d-1)(d-2)(\mathbf{x} \cdot \mathbf{y})^3 + \dots \\ &= 1 + d \sum_i x_i y_i + \frac{1}{2}d(d-1) \sum_{ij} (x_i x_j)(y_i y_j) \\ &\quad + \frac{1}{6}d(d-1)(d-2) \sum_{ijk} (x_i x_j x_k)(y_i y_j y_k) + \dots \end{aligned}$$

This expression is indeed seen to be of the form $K(\mathbf{x}, \mathbf{y}) = \sum_{\ell=1}^{\infty} \phi_{\ell}(\mathbf{x})\phi_{\ell}(\mathbf{y})$, describing pre-processing with polynomial functions $\phi_{\ell}(\mathbf{x})$. The series truncates only if and only if d is chosen to be a positive integer.

- Kernels generating *pre-processing via radial basis functions*: $K(\mathbf{x}, \mathbf{y}) = K(|\mathbf{x} - \mathbf{y}|)$

Let us consider non-negative definite kernels of the form $K(\mathbf{x} - \mathbf{y})$, such that the Fourier transform $\hat{K}(\mathbf{k})$ of the function $K(\mathbf{z})$ exists:

$$\hat{K}(\mathbf{k}) = \int d\mathbf{x} e^{i\mathbf{k} \cdot \mathbf{x}} K(\mathbf{x}) \quad K(\mathbf{x}) = \int \frac{d\mathbf{k}}{(2\pi)^M} e^{-i\mathbf{k} \cdot \mathbf{x}} \hat{K}(\mathbf{k})$$

We first show that $\hat{K}(\mathbf{k}) \geq 0$ for all $\mathbf{k} \in \Re^M$, due to K being a non-negative kernel. This follows from working out

$$\int d\mathbf{y} K(\mathbf{x}, \mathbf{y}) e^{i\mathbf{k} \cdot \mathbf{y}} = \int d\mathbf{y} K(|\mathbf{x} - \mathbf{y}|) e^{i\mathbf{k} \cdot \mathbf{y}} = \hat{K}(\mathbf{k}) e^{i\mathbf{k} \cdot \mathbf{x}}$$

Thus the $\hat{K}(\mathbf{k})$ are seen to be eigenvalues of the kernel K , and therefore non-negative. We can now construct the following symmetric kernel:

$$L(\mathbf{x}, \mathbf{y}) = L(\mathbf{x} - \mathbf{y}) \quad L(\mathbf{x}) = \int \frac{d\mathbf{k}}{(2\pi)^M} e^{-i\mathbf{k} \cdot \mathbf{x}} \sqrt{\hat{K}(\mathbf{k})}$$

We work out the product

$$\begin{aligned} \int d\mathbf{z} L(\mathbf{x} - \mathbf{z})L(\mathbf{z} - \mathbf{y}) &= \int \frac{d\mathbf{u}d\mathbf{v}}{(2\pi)^{2M}} \sqrt{\hat{K}(\mathbf{u})\hat{K}(\mathbf{v})} \int d\mathbf{z} e^{-i\mathbf{u} \cdot (\mathbf{x} - \mathbf{z}) - i\mathbf{v} \cdot (\mathbf{z} - \mathbf{y})} \\ &= \int \frac{d\mathbf{u}d\mathbf{v}}{(2\pi)^{2M}} \sqrt{\hat{K}(\mathbf{u})\hat{K}(\mathbf{v})} e^{-i\mathbf{u} \cdot \mathbf{x} + i\mathbf{v} \cdot \mathbf{y}} \int d\mathbf{z} e^{i\mathbf{z} \cdot (\mathbf{u} - \mathbf{v})} \\ &= \int \frac{d\mathbf{u}}{(2\pi)^M} \hat{K}(\mathbf{u}) e^{-i\mathbf{u} \cdot (\mathbf{x} - \mathbf{y})} = K(\mathbf{x} - \mathbf{y}) \end{aligned}$$

(where we used the identity $\delta(\mathbf{z}) = (2\pi)^{-M} \int d\mathbf{u} e^{i\mathbf{u} \cdot \mathbf{z}}$, see also appendix C). If $K(\mathbf{x})$ is of the form $K(|\mathbf{x}|)$, then also $\hat{K}(\mathbf{k}) = \hat{K}(|\mathbf{k}|)$ and $L(\mathbf{x}) = L(|\mathbf{x}|)$, so that

$$K(|\mathbf{x} - \mathbf{y}|) = \int d\mathbf{z} \phi_{\mathbf{z}}(\mathbf{x})\phi_{\mathbf{z}}(\mathbf{y}), \quad \phi_{\mathbf{z}}(\mathbf{x}) = L(|\mathbf{x} - \mathbf{z}|)$$

This describes pre-processing with radial basis functions.

- Kernels generating *pre-processing via sigmoidal neurons*: $K(\mathbf{x}, \mathbf{y}) = g[\mathbf{x} \cdot \mathbf{y}]$

Here we choose $g[u]$ to be sigmoidal function, such as $g[u] = \tanh[au+b]$, with constants $a, b \in \mathfrak{R}$. This kernel can not always be written in the separable form $K(\mathbf{x}, \mathbf{y}) = \sum_{\ell \geq 1} \phi_{\ell}(\mathbf{x})\phi_{\ell}(\mathbf{y})$, but by simply working out equation (5.27) it becomes clear that we will be doing pre-processing via a ‘hidden’ layer of sigmoidal neurons:

$$S(\mathbf{x}) = \operatorname{sgn} \left[\sum_{\mu=1}^p \lambda_{\mu} t^{\mu} g[\mathbf{x}^{\mu} \cdot \mathbf{x}] + w_0 \right] \quad w_0 = -\frac{1}{2} \sum_{\mu=1}^p \lambda_{\mu} t^{\mu} \{g[\mathbf{x}^{\mu} \cdot \mathbf{x}^+] + g[\mathbf{x}^{\mu} \cdot \mathbf{x}^-]\}$$

The number of support vectors will be the number of hidden neurons, and the input-to-hidden weights are given by the support vectors $\boldsymbol{\xi}^{\mu}$ themselves.

Application Example. We close this section with an example of the application of Support Vector Machines to a binary classification task which is not linearly separable. We have generated a data set D with $p = 30$ and with data points $\boldsymbol{\xi}^{\mu} \in \mathfrak{R}^2$; half of these points have $t^{\mu} = 1$, these are located close to the origin, the other half have $t^{\mu} = -1$ and are located further away from the origin. We try to separate the two classes with an SVM using the polynomial kernel $K(\mathbf{x}, \mathbf{y}) = [\mathbf{x} \cdot \mathbf{y} + 1]^d$. The resulting decision boundaries (5.29) in the input plane \mathfrak{R}^2 are shown, together with the input vectors, in figure 5.2 (after solving numerically the constrained optimization problem (5.28) for the $\{\lambda_{\mu}\}$) for different values of the parameter $d \in \{2, 3, 4, 5\}$. Note that for the present kernel, the parameter d controls the number of pre-processing functions used. As d is increased (and thus the complexity of the SVM) we observe decision boundaries with increasingly detailed class separation, and also a decrease in the number of support vectors. This example illustrates also one of the unresolved issues with Support Vector Machines (which is in fact being studied intensively at the moment): the method appears to be outside the province of Bayesian analysis, and hence does not yet allow for dealing properly with noisy data, overfitting and decision reliability.

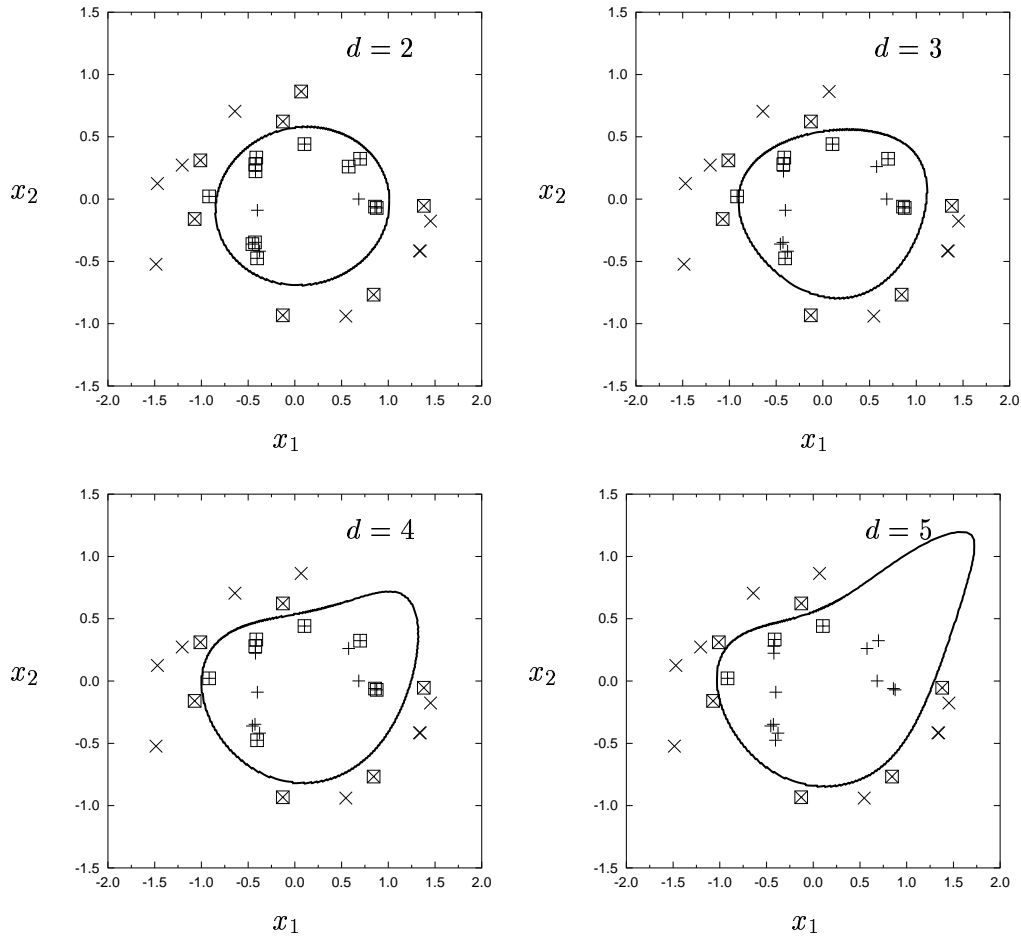


Figure 5.2: Binary classification of $p = 30$ non-linearly-separable input vectors $\mathbf{x}^\mu = (x_1^\mu, x_2^\mu)$ (half with $t^\mu = 1$, close to the origin, indicated by '+'; the other half with $t^\mu = -1$, away from the origin, indicated by 'x'), by an SVM systems with kernel $K(\mathbf{x}, \mathbf{y}) = [1 + \mathbf{x} \cdot \mathbf{y}]^d$. From top left to bottom right: $d \in \{2, 3, 4, 5\}$. Thick curve: the decision boundary as given by (5.29). The support vectors of the solution are those marked with '□'.

Appendix A

Probability in a Nutshell

We define ‘events’ \mathbf{x} as n -dimensional vectors, drawn from some event set $A \subseteq \mathfrak{R}^n$. We associate with each event a real-valued and non-negative probability measure $p(\mathbf{x}) \geq 0$.

Discrete Event Sets

Definitions & Conventions. If A is discrete and countable, each component x_i of \mathbf{x} can only assume values from a discrete set A_i so $A \subseteq A_1 \otimes A_2 \otimes \dots \otimes A_n$. We drop explicit mentioning of sets where possible; e.g. \sum_{x_i} will mean $\sum_{x_i \in A_i}$, and $\sum_{\mathbf{x}}$ will mean $\sum_{\mathbf{x} \in A}$, etc. No problems arise as long as the arguments of $p(\dots)$ are symbols; only once we evaluate probabilities for explicit values of the arguments we need to indicate to which components of \mathbf{x} such values are assigned. The probabilities are normalized according to $\sum_{\mathbf{x}} p(\mathbf{x}) = 1$.

Interpretation of Probability. Imagine a system which generates events $\mathbf{x} \in A$ sequentially, giving the infinite series $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots$. We choose an arbitrary one-to-one index mapping $\pi : \{1, 2, \dots\} \rightarrow \{1, 2, \dots\}$, and one particular event $\mathbf{x} \in A$ (in that order), and calculate for the first M sequence elements $\{\mathbf{x}_{\pi(1)}, \dots, \mathbf{x}_{\pi(M)}\}$ the frequency $f_M(\mathbf{x})$ with which \mathbf{x} occurred:

$$f_M(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \delta_{\mathbf{x}, \mathbf{x}_{\pi(m)}} \quad \begin{cases} \delta_{\mathbf{x}, \mathbf{y}} = 1 & \text{if } \mathbf{x} = \mathbf{y} \\ \delta_{\mathbf{x}, \mathbf{y}} = 0 & \text{if } \mathbf{x} \neq \mathbf{y} \end{cases}$$

We define *random events* as those generated by a system as above with the property that for each one-to-one index map π , for each event $\mathbf{x} \in A$ the frequency of occurrence $f_M(\mathbf{x})$ tends to a limit as $M \rightarrow \infty$. This limit is then defined as the ‘probability’ associated with \mathbf{x} :

$$\forall \mathbf{x} \in A : \quad p(\mathbf{x}) = \lim_{M \rightarrow \infty} f_M(\mathbf{x})$$

Since $f_M(\mathbf{x}) \geq 0$ for each \mathbf{x} , and since $\sum_{\mathbf{x}} f_M(\mathbf{x}) = 1$ for any M , it follows that $p(\mathbf{x}) \geq 0$ and that $\sum_{\mathbf{x}} p(\mathbf{x}) = 1$ (as it should).

Marginal & Conditional Probabilities, Statistical Independence. The so-called ‘marginal probabilities’ are obtained upon summing over individual components of $\mathbf{x} = (x_1, \dots, x_n)$:

$$p(x_1, \dots, x_{\ell-1}, x_{\ell+1}, \dots, x_n) = \sum_{x_{\ell}} p(x_1, \dots, x_n) \quad (\text{A.1})$$

In particular we obtain (upon repeating this procedure):

$$p(x_i) = \sum_{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n} p(x_1, \dots, x_n) \quad (\text{A.2})$$

Marginal probabilities are normalized. This follows upon combining their definition with the basic normalization $\sum_{\mathbf{x}} p(\mathbf{x}) = 1$, e.g.

$$\sum_{x_1, \dots, x_{\ell-1}, x_{\ell+1}, \dots, x_n} p(x_1, \dots, x_{\ell-1}, x_{\ell+1}, \dots, x_n) = 1, \quad \sum_{x_i} p(x_i) = 1$$

For any two disjoint subsets $\{i_1, \dots, i_k\}$ and $\{j_1, \dots, j_\ell\}$ of the index set $\{1, \dots, n\}$ (with necessarily $k + \ell \leq n$) we next define the ‘conditional probability’

$$p(x_{i_1}, \dots, x_{i_k} | x_{j_1}, \dots, x_{j_\ell}) = \frac{p(x_{i_1}, \dots, x_{i_k}, x_{j_1}, \dots, x_{j_\ell})}{p(x_{j_1}, \dots, x_{j_\ell})} \quad (\text{A.3})$$

(A.3) gives the probability that the k components $\{i_1, \dots, i_k\}$ of \mathbf{x} take the values $\{x_{i_1}, \dots, x_{i_k}\}$, given the knowledge that the ℓ components $\{j_1, \dots, j_\ell\}$ take the values $\{x_{j_1}, \dots, x_{j_\ell}\}$.

The concept of statistical independence now follows naturally. Loosely speaking: statistical independence means that conditioning in the sense defined above does not affect any of the marginal probabilities. Thus the n events $\{x_1, \dots, x_n\}$ are said to be statistically independent if for *any* two disjoint subsets $\{i_1, \dots, i_k\}$ and $\{j_1, \dots, j_\ell\}$ of $\{1, \dots, n\}$ we have

$$p(x_{i_1}, \dots, x_{i_k} | x_{j_1}, \dots, x_{j_\ell}) = p(x_{i_1}, \dots, x_{i_k})$$

This can be shown to be equivalent to saying

$$\{x_1, \dots, x_n\} \text{ are independent : } \begin{array}{l} p(x_{i_1}, \dots, x_{i_k}) = p(x_{i_1})p(x_{i_2}) \dots p(x_{i_k}) \\ \text{for every subset } \{i_1, \dots, i_k\} \subseteq \{1, \dots, n\} \end{array} \quad (\text{A.4})$$

Continuous Event Sets

Definitions & Conventions. Here A is no longer countable. Each component x_i of \mathbf{x} can assume values from a continuous set $A_i = \{x \in \mathfrak{R} \mid \exists \mathbf{x} \in A \text{ with } x_i = x\}$, and $A \subseteq A_1 \otimes A_2 \otimes \dots \otimes A_n$. As before we drop explicit mentioning of sets where possible; e.g. $\int dx_i$ will mean $\int_{A_i} dx_i$ and $\int d\mathbf{x}$ will mean $\int_A d\mathbf{x}$, etc. The function $p(\mathbf{x})$ is now to be seen as a *probability density*, which is accordingly normalized via integration: $\int d\mathbf{x} p(\mathbf{x}) = 1$.

Interpretation of Probability. Again we imagine a system which generates events $\mathbf{x} \in A$ sequentially, giving $\mathbf{x}_1, \mathbf{x}_2, \dots$. We define ‘boxes’ (hypercubes) $B(\mathbf{x}, \Delta)$ in \mathfrak{R}^n as follows:

$$B(\mathbf{x}, \Delta) = \{\mathbf{y} \in \mathfrak{R}^n \mid x_i \leq y_i < x_i + \Delta_i \text{ for all } i\}$$

in which all $\Delta_i > 0$. The volume of such a box is simply $\prod_{i=1}^n \Delta_i$. We choose an arbitrary one-to-one index mapping $\pi : \{1, 2, \dots\} \rightarrow \{1, 2, \dots\}$, and one particular event $\mathbf{x} \in A$ (in that order), and calculate for the first M sequence elements $\{\mathbf{x}_{\pi(1)}, \dots, \mathbf{x}_{\pi(M)}\}$ the frequency $f_M(\mathbf{x}, \Delta)$ with which \mathbf{x} happened to lie in box $B(\mathbf{x}, \Delta)$:

$$f_M(\mathbf{x}, \Delta) = \frac{1}{M} \sum_{m=1}^M I(\mathbf{x}_{\pi(m)}; \mathbf{x}, \Delta) \quad \begin{cases} I(\mathbf{y}; \mathbf{x}, \Delta) = 1 & \text{if } \mathbf{y} \in B(\mathbf{x}, \Delta) \\ I(\mathbf{y}; \mathbf{x}, \Delta) = 0 & \text{if } \mathbf{y} \notin B(\mathbf{x}, \Delta) \end{cases}$$

We define *random events* as those generated by a system as above with the property that for each one-to-one index map π , for each box $B(\mathbf{x}, \Delta)$ the frequency of occurrence $f_M(\mathbf{x}, \Delta)$ tends to a limit as $M \rightarrow \infty$. This limit is subsequently used to define the ‘probability density’ $p(\mathbf{x})$ associated with \mathbf{x} :

$$\forall \mathbf{x} \in A : \quad p(\mathbf{x}) = \lim_{\Delta \rightarrow \mathbf{0}} \frac{\lim_{M \rightarrow \infty} f_M(\mathbf{x}, \Delta)}{\prod_{i=1}^n \Delta_i}$$

(provided this limit exists). Since $f_M(\mathbf{x}, \Delta) \geq 0$ for each \mathbf{x} , and since $\int d\mathbf{x} f_M(\mathbf{x}, \Delta) = \prod_{i=1}^n \Delta_i$ for any M , it follows that $p(\mathbf{x}) \geq 0$ and that $\int d\mathbf{x} p(\mathbf{x}) = 1$ (as it should).

Marginal & Conditional Probability Densities, Statistical Independence. The marginal probabilities are now obtained upon integrating over individual components of \mathbf{x} :

$$p(x_1, \dots, x_{\ell-1}, x_{\ell+1}, \dots, x_n) = \int dx_{\ell} p(x_1, \dots, x_n) \quad (\text{A.5})$$

In particular we obtain (upon repeating this procedure):

$$p(x_i) = \int dx_1 \dots dx_{i-1} dx_{i+1} \dots dx_n p(x_1, \dots, x_n) \quad (\text{A.6})$$

Again marginal probabilities are normalized (in integration sense). This follows as before upon combining their definition with the basic normalization $\int d\mathbf{x} p(\mathbf{x}) = 1$, e.g.

$$\int dx_1 \dots dx_{\ell-1} dx_{\ell+1} \dots dx_n p(x_1, \dots, x_{\ell-1}, x_{\ell+1}, \dots, x_n) = 1, \quad \int dx_i p(x_i) = 1$$

For any two disjunct subsets $\{i_1, \dots, i_k\}$ and $\{j_1, \dots, j_{\ell}\}$ of the index set $\{1, \dots, n\}$ we define the ‘conditional probability density’

$$p(x_{i_1}, \dots, x_{i_k} | x_{j_1}, \dots, x_{j_{\ell}}) = \frac{p(x_{i_1}, \dots, x_{i_k}, x_{j_1}, \dots, x_{j_{\ell}})}{p(x_{j_1}, \dots, x_{j_{\ell}})} \quad (\text{A.7})$$

It gives the probability density for the k components $\{i_1, \dots, i_k\}$, given the knowledge that the ℓ components $\{j_1, \dots, j_{\ell}\}$ take the values $\{x_{j_1}, \dots, x_{j_{\ell}}\}$.

Statistical independence can as before be defined in the following way:

$$\{x_1, \dots, x_n\} \text{ are independent :} \quad \begin{aligned} p(x_{i_1}, \dots, x_{i_k}) &= p(x_{i_1})p(x_{i_2}) \dots p(x_{i_k}) \\ \text{for every subset } \{i_1, \dots, i_k\} &\subseteq \{1, \dots, n\} \end{aligned} \quad (\text{A.8})$$

Averages of Specific Random Variables

We define ‘random variables’ as arbitrary functions $F(\mathbf{x})$ of random events $\mathbf{x} \in A$. We define *averages* or *expectation values* or *mean values* $\langle F(\mathbf{x}) \rangle$ of random variables $F(\mathbf{x})$ as follows:

$$\begin{aligned} \text{discrete random variables :} & \quad \langle F(\mathbf{x}) \rangle = \sum_{\mathbf{x}} p(\mathbf{x}) F(\mathbf{x}) \\ \text{continuous random variables :} & \quad \langle F(\mathbf{x}) \rangle = \int d\mathbf{x} p(\mathbf{x}) F(\mathbf{x}) \end{aligned} \quad (\text{A.9})$$

Let us now turn to the definitions and properties of a couple of relevant random variables (discrete or continuous) and averages. Note that in general one cannot be sure beforehand (without explicit proof) that the following averages will actually exist (i.e. are finite):

average: $\mu_i = \langle x_i \rangle$

variance: $\sigma_i^2 = \langle x_i^2 \rangle - \langle x_i \rangle^2$

σ_i^2 is non-negative, since $\langle x_i^2 \rangle - \langle x_i \rangle^2 = \langle [x_i - \langle x_i \rangle]^2 \rangle$. This also shows that $\sigma_i^2 = 0$ implies that $x_i = x'_i$ for any two events $\mathbf{x} \in A$ and $\mathbf{x}' \in A$ with nonzero probabilities.

covariance matrix: $C_{ij} = \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle$

Note that $C_{ii} = \sigma_i^2$. The covariance matrix is symmetric so all eigenvalues are real. It is non-negative definite (so all eigenvalues are non-negative), since it can be written as $\langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle = \langle (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle) \rangle$, from which it follows that

$$\text{for any } \mathbf{z} \in \Re^n : \quad \mathbf{z} \cdot \mathbf{C} \mathbf{z} = \langle \left\{ \sum_{i=1}^n z_i [x_i - \langle x_i \rangle] \right\}^2 \rangle \geq 0$$

moments: $\langle x_{i_1}^{m_{i_1}} x_{i_2}^{m_{i_2}} \dots x_{i_k}^{m_{i_k}} \rangle$, with $m_{i_l} \in \{0, 1, 2, 3, \dots\}$

Gaussian (or ‘Normal’) Distribution:

$$p(x) = \frac{e^{-\frac{1}{2}(x-\mu)^2/\sigma^2}}{\sigma\sqrt{2\pi}}, \quad x \in \Re \tag{A.10}$$

Normalization is built-in (see also appendix B):

$$\int dx p(x) = \int \frac{dx}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu)^2/\sigma^2} = \int \frac{dz}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} = 1$$

For the Gaussian distribution we also obtain:

$$\langle (x-\mu)^n \rangle = \int \frac{dx}{\sigma\sqrt{2\pi}} (x-\mu)^n e^{-\frac{1}{2}(x-\mu)^2/\sigma^2} = \sigma^n \int \frac{dz}{\sqrt{2\pi}} z^n e^{-\frac{1}{2}z^2}$$

For n odd the result is zero. For n even we get:

$$\begin{aligned} \langle (x-\mu)^{2m} \rangle &= \sigma^{2m} (-1)^m \lim_{y \rightarrow \frac{1}{2}} \frac{d^m}{dy^m} \int \frac{dz}{\sqrt{2\pi}} e^{-yz^2} = \frac{\sigma^{2m}}{\sqrt{2}} (-1)^m \lim_{y \rightarrow \frac{1}{2}} \frac{d^m}{dy^m} y^{-\frac{1}{2}} \\ &= \frac{\sigma^{2m}}{\sqrt{2}} \lim_{y \rightarrow \frac{1}{2}} \left[\frac{1}{2} \cdot \frac{3}{2} \cdot \dots \cdot \frac{2m-1}{2} \right] y^{-\frac{1}{2}-m} = \sigma^{2m} [1.3 \dots (2m-1)] \end{aligned}$$

From these relations, in turn, follow the moments. For example:

$$\begin{aligned} \langle (x-\mu) \rangle = 0 &\Rightarrow \langle x \rangle = \mu \\ \langle (x-\mu)^2 \rangle = \sigma^2 &\Rightarrow \langle x^2 \rangle = \mu^2 + \sigma^2 \\ \langle (x-\mu)^3 \rangle = 0 &\Rightarrow \langle x^3 \rangle = \mu^3 + 3\mu\sigma^2 \\ \langle (x-\mu)^4 \rangle = 0 &\Rightarrow \langle x^4 \rangle = 3(\mu^2 + \sigma^2)^2 \end{aligned}$$

Relations such as $\langle x^4 \rangle = 3(\mu^2 + \sigma^2)^2$ can often serve as a quick test to see whether an unknown distribution could be Gaussian.

Appendix B

Gaussian Integrals

Real, Symmetric, Positive Definite Matrices. The symmetric $N \times N$ matrix \mathbf{A} is assumed to be positive definite, i.e. $\mathbf{x} \cdot \mathbf{A}\mathbf{x} > 0$ for all $\mathbf{x} \in \mathfrak{R}^N$ with $|\mathbf{x}| \neq 0$. The eigenvalue polynomial $\det[\mathbf{A} - \lambda\mathbf{I}] = 0$ is of order N , so \mathbf{A} will have N (possibly complex) solutions λ (some may coincide) of the eigenvalue problem $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ ($\mathbf{x} \neq 0$). We denote complex conjugation of complex numbers z as: $z = a + ib$, $z^* = a - ib$ ($a, b \in \mathfrak{R}$), and $|z|^2 = z^*z \in \mathfrak{R}$. We define the unit matrix $\mathbf{I}_{ij} = \delta_{ij}$. We will use the following facts, proofs of which are found in any elementary linear algebra course:

Fact 1: All eigenvalues of the matrix \mathbf{A} are real.

Fact 2: All eigenvectors can be chosen real-valued.

Fact 3: All eigenvalues λ are positive.

Fact 4: We can construct a complete orthogonal basis in \mathfrak{R}^N of \mathbf{A} -eigenvectors.

Thus there exist a set of vectors $\{\hat{\mathbf{e}}^i\}$ ($\lambda = 1, \dots, N$) with the properties:

$$\mathbf{A}\hat{\mathbf{e}}^i = \lambda_i\hat{\mathbf{e}}^i \quad \lambda_i \in \mathfrak{R}, \lambda_i > 0 \quad \hat{\mathbf{e}}^i \in \mathfrak{R}^N, \hat{\mathbf{e}}^i \cdot \hat{\mathbf{e}}^j = \delta_{ij}$$

We can now bring \mathbf{A} onto diagonal form by a unitary transformation \mathbf{U} , which we construct from the components of the normalised eigenvectors $\hat{\mathbf{e}}: U_{ij} = \hat{\mathbf{e}}_i^j$. We denote its transpose by \mathbf{U}^\dagger , $U_{ij}^\dagger = U_{ji}$, and show that \mathbf{U} is indeed unitary, i.e. $\mathbf{U}^\dagger\mathbf{U} = \mathbf{U}\mathbf{U}^\dagger = \mathbf{I}$:

$$\begin{aligned} \sum_j (\mathbf{U}^\dagger\mathbf{U})_{ij} x_j &= \sum_{jk} U^{ki} U_{kj} x_j = \sum_{jk} \hat{\mathbf{e}}_k^i \hat{\mathbf{e}}_k^j x_j = \sum_j \delta_{ij} x_j = x_i \\ \sum_j (\mathbf{U}\mathbf{U}^\dagger)_{ij} x_j &= \sum_{jk} U^{ik} U_{jk} x_j = \sum_{jk} \hat{\mathbf{e}}_i^k \hat{\mathbf{e}}_j^k x_j = \sum_k \hat{\mathbf{e}}_i^k (\hat{\mathbf{e}} \cdot \mathbf{x}) = x_i \end{aligned}$$

(since $\{\hat{\mathbf{e}}^\ell\}$ forms a complete orthogonal basis). From \mathbf{U} being unitary it follows that \mathbf{U} and \mathbf{U}^\dagger leave inner products invariant:

$$\mathbf{U}\mathbf{x} \cdot \mathbf{U}\mathbf{y} = \mathbf{x} \cdot \mathbf{U}^\dagger\mathbf{U}\mathbf{y} = \mathbf{x} \cdot \mathbf{y} \quad \mathbf{U}^\dagger\mathbf{x} \cdot \mathbf{U}^\dagger\mathbf{y} = \mathbf{x} \cdot \mathbf{U}\mathbf{U}^\dagger\mathbf{y} = \mathbf{x} \cdot \mathbf{y}$$

We can see explicitly that \mathbf{U} indeed brings \mathbf{A} onto diagonal form:

$$(\mathbf{U}^\dagger\mathbf{A}\mathbf{U})_{ij} = \sum_{kl=1}^N U_{ik}^\dagger A_{kl} U_{lj} = \sum_{kl=1}^N \hat{\mathbf{e}}_k^i A_{kl} \hat{\mathbf{e}}_l^j = \lambda_j \sum_{k=1}^N \hat{\mathbf{e}}_k^i \hat{\mathbf{e}}_k^j = \lambda_j \delta_{ij} \quad (\text{B.1})$$

The inverse \mathbf{A}^{-1} of the matrix \mathbf{A} exists, and can be written as $(\mathbf{A}^{-1})_{ij} = \sum_{k=1}^N \lambda_k^{-1} \hat{\mathbf{e}}_i^k \hat{\mathbf{e}}_j^k$.

Gaussian Integrals. We now turn to the associated Gaussian integrals

$$I = \int d\mathbf{x} f(\mathbf{x}) e^{-\frac{1}{2}\mathbf{x}\cdot\mathbf{A}\mathbf{x}} \quad (\text{B.2})$$

The simplest such integral is $I = \int d\mathbf{x} e^{-\frac{1}{2}\mathbf{x}^2}$. It is calculated by writing its square as an integral in \mathfrak{R}^2 , and switching to polar coordinates, ($z_1 = r \cos \phi$, $z_2 = r \sin \phi$). The Jacobian of this transformation is r . Hence

$$I^2 = \int dz_1 dz_2 e^{-\frac{1}{2}\mathbf{z}^2} = \int_0^{2\pi} d\phi \int_0^\infty dr r e^{-\frac{1}{2}r^2} = 2\pi \left[-e^{-\frac{1}{2}r^2} \right]_0^\infty = 2\pi \quad \text{so} \quad I = \sqrt{2\pi}$$

For $f(\mathbf{x}) = 1$ we can do the integral (B.2) by using the previous results on the diagonalisability of the matrix \mathbf{A} . We put $\mathbf{x} = \mathbf{U}\mathbf{z}$ (since \mathbf{U} leaves inner products invariant: $d\mathbf{x} = d\mathbf{z}$):

$$\begin{aligned} \int d\mathbf{x} e^{-\frac{1}{2}\mathbf{x}\cdot\mathbf{A}\mathbf{x}} &= \int d\mathbf{z} e^{-\frac{1}{2}\mathbf{z}\cdot\mathbf{U}^\dagger\mathbf{A}\mathbf{U}\mathbf{z}} = \prod_{\ell=1}^N \left[\int dz e^{-\frac{1}{2}\lambda_\ell z^2} \right] \\ &= \left[\prod_{\ell=1}^N \frac{1}{\sqrt{\lambda_\ell}} \right] \left[\int dz e^{-\frac{1}{2}z^2} \right]^N = \frac{(2\pi)^{N/2}}{\sqrt{\det \mathbf{A}}} \end{aligned} \quad (\text{B.3})$$

(note: the determinant of \mathbf{A} is unvariant under rotations, so it can be evaluated with \mathbf{A} on diagonal form, which gives the product of the N eigenvalues).

Due to the symmetry of the integrand in (B.2) under reflection $\mathbf{x} \rightarrow -\mathbf{x}$, the integral reduces to zero for $f(\mathbf{x}) = x_i$. For $f(\mathbf{x}) = x_i x_j$ we find:

$$\begin{aligned} \int d\mathbf{x} x_i x_j e^{-\frac{1}{2}\mathbf{x}\cdot\mathbf{A}\mathbf{x}} &= \lim_{\mathbf{b} \rightarrow 0} \frac{\partial^2}{\partial b_i \partial b_j} \int d\mathbf{x} e^{-\frac{1}{2}\mathbf{x}\cdot\mathbf{A}\mathbf{x} + \mathbf{b}\cdot\mathbf{x}} \\ &= \lim_{\mathbf{b} \rightarrow 0} \frac{\partial^2}{\partial b_i \partial b_j} \int d\mathbf{z} e^{-\frac{1}{2}\mathbf{z}\cdot\mathbf{U}^\dagger\mathbf{A}\mathbf{U}\mathbf{z} + \mathbf{z}\cdot\mathbf{U}^\dagger\mathbf{b}} = \lim_{\mathbf{b} \rightarrow 0} \frac{\partial^2}{\partial b_i \partial b_j} \prod_{\ell=1}^N \left[\int dz e^{-\frac{1}{2}\lambda_\ell z^2 + z(\mathbf{U}^\dagger\mathbf{b})_\ell} \right] \\ &= \lim_{\mathbf{b} \rightarrow 0} \frac{\partial^2}{\partial b_i \partial b_j} \prod_{\ell=1}^N \left[\int dz e^{-\frac{1}{2}\lambda_\ell [z - (\mathbf{U}^\dagger\mathbf{b})_\ell \lambda_\ell^{-1}]^2 + \frac{1}{2}\lambda_\ell^{-1} (\mathbf{U}^\dagger\mathbf{b})_\ell^2} \right] \\ &= \lim_{\mathbf{b} \rightarrow 0} \frac{\partial^2}{\partial b_i \partial b_j} e^{\frac{1}{2} \sum_{ij\ell=1}^N \lambda_\ell^{-1} U_{i\ell} b_i U_{j\ell} b_j} \prod_{\ell=1}^N \left[\int dz e^{-\frac{1}{2}\lambda_\ell z^2} \right] \\ &= \frac{(2\pi)^{N/2}}{\sqrt{\det \mathbf{A}}} \lim_{\mathbf{b} \rightarrow 0} \frac{\partial^2}{\partial b_i \partial b_j} e^{\frac{1}{2} \sum_{ij=1}^N b_i b_j \sum_{\ell=1}^N \lambda_\ell^{-1} \hat{e}_i^\ell \hat{e}_j^\ell} = (A^{-1})_{ij} \frac{(2\pi)^{N/2}}{\sqrt{\det \mathbf{A}}} \end{aligned} \quad (\text{B.4})$$

In particular, by combining the last two results, we find a powerful (and completely general) relation for Gaussian probability distributions with zero mean $\mathbf{x} = 0$ (the latter one can always achieve by a simple translation):

$$\frac{\int d\mathbf{x} x_i x_j e^{-\frac{1}{2}\mathbf{x}\cdot\mathbf{A}\mathbf{x}}}{\int d\mathbf{x} e^{-\frac{1}{2}\mathbf{x}\cdot\mathbf{A}\mathbf{x}}} = (A^{-1})_{ij} \quad (\text{B.5})$$

If we know that a given distribution is Gaussian, with zero average, we apparently only need to calculate the correlations $\langle x_i x_j \rangle$ to know the full distribution:

$$P(\mathbf{x}) \text{ Gaussian, with } \langle \mathbf{x} \rangle = 0 \quad \Rightarrow \quad P(\mathbf{x}) = \frac{e^{-\frac{1}{2}\mathbf{x}\cdot\mathbf{A}\mathbf{x}}}{(2\pi)^{N/2} \det^{-\frac{1}{2}} \mathbf{A}}, \quad \text{with } (A^{-1})_{ij} = \langle x_i x_j \rangle \quad (\text{B.6})$$

Appendix C

The δ -Distribution

Definition. There are several ways of introducing the δ -distribution. Here we will go for an intuitive definition first, and a formal one later. We define the δ -distribution as the probability distribution $\delta(x)$ corresponding to a random variable in the limit where the randomness in the variable vanishes. If x is ‘distributed’ around zero, this implies

$$\int dx f(x)\delta(x) = f(0) \quad \text{for any function } f$$

The problem arises when we want to actually write down an expression for $\delta(x)$. Intuitively one could think of writing something like

$$\delta(x) = \lim_{\Delta \rightarrow 0} G_{\Delta}(x) \quad G_{\Delta}(x) = \frac{1}{\Delta\sqrt{2\pi}} e^{-\frac{1}{2}x^2/\Delta^2} \quad (\text{C.1})$$

This is not a true function in a mathematical sense; $\delta(x)$ is zero for $x \neq 0$ and $\delta(0) = \infty$. The way to interpret and use expressions like (C.1) is to realise that $\delta(x)$ only has a meaning when appearing inside an integration. One then takes the limit $\Delta \rightarrow 0$ *after* performing the integration. Upon adopting this convention, we can use (C.1) to derive the following properties (for sufficiently well-behaved and differentiable functions f^1):

$$\int dx \delta(x)f(x) = \lim_{\Delta \rightarrow 0} \int dx G_{\Delta}(x)f(x) = \lim_{\Delta \rightarrow 0} \int \frac{dx}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} f(\Delta x) = f(0)$$

$$\begin{aligned} \int dx \delta'(x)f(x) &= \lim_{\Delta \rightarrow 0} \int dx \left\{ \frac{d}{dx} [G_{\Delta}(x)f(x)] - G_{\Delta}(x)f'(x) \right\} \\ &= \lim_{\Delta \rightarrow 0} [G_{\Delta}(x)f(x)]_{-\infty}^{\infty} - f'(0) = -f'(0) \end{aligned}$$

both can be summarised in and generalised to the single expression:

$$\int dx f(x) \frac{d^n}{dx^n} \delta(x) = (-1)^n \lim_{x \rightarrow 0} \frac{d^n}{dx^n} f(x) \quad (n = 0, 1, 2, \dots) \quad (\text{C.2})$$

¹The conditions on the so-called ‘test-functions’ f can be properly formalised; this being not a course on distribution theory, here we just concentrate on the basic ideas and properties

Equivalently we can take the result (C.2) as our definition of the δ -distribution.

Representations, Relations, Generalisations. We can use the definitions of Fourier transforms and inverse Fourier transforms to obtain an integral representation of the δ -distribution:

$$\begin{aligned}\mathcal{F} : f(x) &\rightarrow \hat{f}(k) & \hat{f}(k) &= \int dx e^{-2\pi i k x} f(x) \\ \mathcal{F}^{-1} : \hat{f}(k) &\rightarrow f(x) & f(x) &= \int dk e^{2\pi i k x} \hat{f}(k)\end{aligned}$$

In combination these relations give the identity:

$$f(x) = \int dk e^{2\pi i k x} \int dy e^{-2\pi i k y} f(y)$$

Application to $f(x) = \delta(x)$ gives:

$$\delta(x) = \int dk e^{2\pi i k x} = \int \frac{dk}{2\pi} e^{i k x} \quad (\text{C.3})$$

A second useful relation is the following one, which relates the δ -distribution to the step-function:

$$\delta(x) = \frac{d}{dx} \theta(x) \quad (\text{C.4})$$

This we prove by showing that both have the same effect inside an integration (with an arbitrary test-function):

$$\begin{aligned}\int dx \left[\delta(x) - \frac{d}{dx} \theta(x) \right] \phi(x) &= \phi(0) - \lim_{\epsilon \rightarrow 0} \int_{-\epsilon}^{\epsilon} dx \left\{ \frac{d}{dx} [\theta(x)\phi(x)] - \phi'(x)\theta(x) \right\} \\ &= \phi(0) - \lim_{\epsilon \rightarrow 0} [\phi(\epsilon) - 0] + \lim_{\epsilon \rightarrow 0} \int_0^{\epsilon} dx \phi'(x) = 0\end{aligned}$$

Thirdly we can inspect the effect of performing a continuously differentiable and invertible transformation f on a variable that occurs inside a δ -distribution, giving rise to the following identity:

$$\delta[f(x) - f(a)] = \frac{\delta(x - a)}{f'(a)} \quad (\text{C.5})$$

Again this is proved by showing that both sides have the same effect inside an integration (with an arbitrary test-function):

$$\begin{aligned}\int dx \phi(x) \left\{ \delta[f(x) - f(a)] - \frac{\delta(x - a)}{f'(a)} \right\} &= \int dx f'(x) \left[\frac{\phi(x)}{f'(x)} \right] \delta[f(x) - f(a)] - \frac{\phi(a)}{f'(a)} \\ &= \int dk \frac{\phi(f^{-1}(k))}{f'(f^{-1}(k))} \delta[k - f(a)] - \frac{\phi(a)}{f'(a)} = \frac{\phi(f^{-1}(f(a)))}{f'(f^{-1}(f(a)))} - \frac{\phi(a)}{f'(a)} = 0\end{aligned}$$

Finally, the following generalisation is straightforward:

$$\mathbf{x} \in \mathcal{R}^N : \quad \delta(\mathbf{x}) = \prod_{i=1}^N \delta(x_i) \quad (\text{C.6})$$