

Perceptron Learning

4.1 Learning algorithms for neural networks

In the two preceding chapters we discussed two closely related models, McCulloch–Pitts units and perceptrons, but the question of how to find the parameters adequate for a given task was left open. If two sets of points have to be separated linearly with a perceptron, adequate weights for the computing unit must be found. The operators that we used in the preceding chapter, for example for edge detection, used hand customized weights. Now we would like to find those parameters automatically. The *perceptron learning algorithm* deals with this problem.

A learning algorithm is an adaptive method by which a network of computing units self-organizes to implement the desired behavior. This is done in some learning algorithms by presenting some examples of the desired input-output mapping to the network. A correction step is executed iteratively until the network learns to produce the desired response. The learning algorithm is a closed loop of presentation of examples and of corrections to the network parameters, as shown in Figure 4.1.

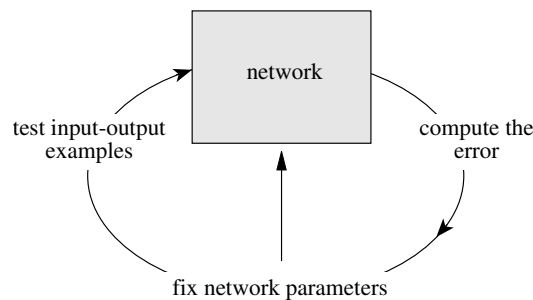


Fig. 4.1. Learning process in a parametric system

In some simple cases the weights for the computing units can be found through a sequential test of stochastically generated numerical combinations. However, such algorithms which look blindly for a solution do not qualify as “learning”. A learning algorithm must adapt the network parameters according to previous experience until a solution is found, if it exists.

4.1.1 Classes of learning algorithms

Learning algorithms can be divided into *supervised* and *unsupervised* methods. Supervised learning denotes a method in which some input vectors are collected and presented to the network. The output computed by the network is observed and the deviation from the expected answer is measured. The weights are corrected according to the magnitude of the error in the way defined by the learning algorithm. This kind of learning is also called *learning with a teacher*, since a control process knows the correct answer for the set of selected input vectors.

Unsupervised learning is used when, for a given input, the exact numerical output a network should produce is unknown. Assume, for example, that some points in two-dimensional space are to be classified into three clusters. For this task we can use a classifier network with three output lines, one for each class (Figure 4.2). Each of the three computing units at the output must specialize by firing only for inputs corresponding to elements of each cluster. If one unit fires, the others must keep silent. In this case we do not know a priori which unit is going to specialize on which cluster. Generally we do not even know how many well-defined clusters are present. Since no “teacher” is available, the network must organize itself in order to be able to associate clusters with units.

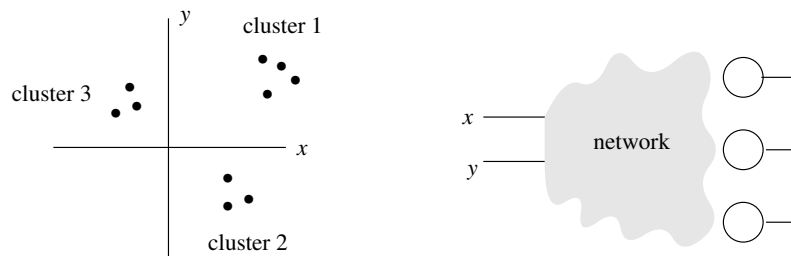


Fig. 4.2. Three clusters and a classifier network

Supervised learning is further divided into methods which use reinforcement or error correction. Reinforcement learning is used when after each presentation of an input-output example we only know whether the network produces the desired result or not. The weights are updated based on this information (that is, the Boolean values *true* or *false*) so that only the input

vector can be used for weight correction. In learning with error correction, the magnitude of the error, together with the input vector, determines the magnitude of the corrections to the weights, and in many cases we try to eliminate the error in a single correction step.

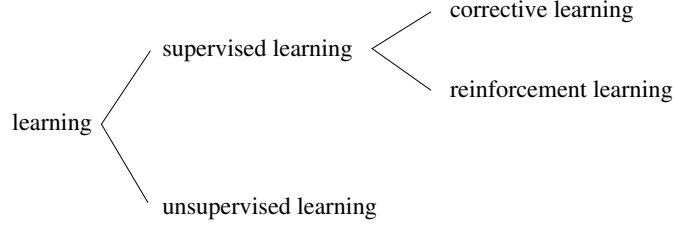


Fig. 4.3. Classes of learning algorithms

The perceptron learning algorithm is an example of supervised learning with reinforcement. Some of its variants use supervised learning with error correction (corrective learning).

4.1.2 Vector notation

In the following sections we deal with learning methods for perceptrons. To simplify the notation we adopt the following conventions. The input (x_1, x_2, \dots, x_n) to the perceptron is called the input vector. If the weights of the perceptron are the real numbers w_1, w_2, \dots, w_n and the threshold is θ , we call $\mathbf{w} = (w_1, w_2, \dots, w_n, w_{n+1})$ with $w_{n+1} = -\theta$ the *extended weight vector* of the perceptron and $(x_1, x_2, \dots, x_n, 1)$ the *extended input vector*. The threshold computation of a perceptron will be expressed using scalar products. The arithmetic test computed by the perceptron is thus

$$\mathbf{w} \cdot \mathbf{x} \geq \theta,$$

if \mathbf{w} and \mathbf{x} are the weight and input vectors, and

$$\mathbf{w} \cdot \mathbf{x} \geq 0$$

if \mathbf{w} and \mathbf{x} are the extended weight and input vectors. It will always be clear from the context whether normal or extended vectors are being used.

If, for example, we are looking for the weights and threshold needed to implement the AND function with a perceptron, the input vectors and their associated outputs are

$$\begin{aligned} (0, 0) &\mapsto 0, \\ (0, 1) &\mapsto 0, \\ (1, 0) &\mapsto 0, \\ (1, 1) &\mapsto 1. \end{aligned}$$

If a perceptron with threshold zero is used, the input vectors must be extended and the desired mappings are

$$\begin{aligned}(0, 0, 1) &\mapsto 0, \\ (0, 1, 1) &\mapsto 0, \\ (1, 0, 1) &\mapsto 0, \\ (1, 1, 1) &\mapsto 1.\end{aligned}$$

A perceptron with three still unknown weights (w_1, w_2, w_3) can carry out this task.

4.1.3 Absolute linear separability

The proof of convergence of the perceptron learning algorithm assumes that each perceptron performs the test $\mathbf{w} \cdot \mathbf{x} > 0$. So far we have been working with perceptrons which perform the test $\mathbf{w} \cdot \mathbf{x} \geq 0$. We must just show that both classes of computing units are equivalent when the training set is finite, as is always the case in learning problems. We need the following definition.

Definition 3. *Two sets A and B of points in an n -dimensional space are called absolutely linearly separable if $n + 1$ real numbers w_1, \dots, w_{n+1} exist such that every point $(x_1, x_2, \dots, x_n) \in A$ satisfies $\sum_{i=1}^n w_i x_i > w_{n+1}$ and every point $(x_1, x_2, \dots, x_n) \in B$ satisfies $\sum_{i=1}^n w_i x_i < w_{n+1}$*

If a perceptron with threshold zero can linearly separate two finite sets of input vectors, then only a small adjustment to its weights is needed to obtain an absolute linear separation. This is a direct corollary of the following proposition.

Proposition 7. *Two finite sets of points, A and B , in n -dimensional space which are linearly separable are also absolutely linearly separable.*

Proof. Since the two sets are linearly separable, weights w_1, \dots, w_{n+1} exist for which, without loss of generality, the inequality

$$\sum_{i=1}^n w_i a_i \geq w_{n+1}$$

holds for all points $(a_1, \dots, a_n) \in A$ and

$$\sum_{i=1}^n w_i b_i < w_{n+1}$$

for all points $(b_1, \dots, b_n) \in B$. Let $\varepsilon = \max\{\sum_{i=1}^n w_i b_i - w_{n+1} \mid (b_1, \dots, b_n) \in B\}$. It is clear that $\varepsilon < \varepsilon/2 < 0$. Let $w' = w_{n+1} + \varepsilon/2$. For all points in A it holds that

$$\sum_{i=1}^n w_i a_i - (w' - \frac{1}{2}\varepsilon) \geq 0.$$

This means that

$$\sum_{i=1}^n w_i a_i - w' \geq -\frac{1}{2}\varepsilon > 0 \Rightarrow \sum_{i=1}^n w_i a_i > w'. \quad (4.1)$$

For all points in B a similar argument holds since

$$\sum_{i=1}^n w_i b_i - w_{n+1} = \sum_{i=1}^n w_i b_i - (w' - \frac{1}{2}\varepsilon) \leq \varepsilon,$$

and from this we deduce

$$\sum_{i=1}^n w_i b_i - w' \leq \frac{1}{2}\varepsilon < 0. \quad (4.2)$$

Equations (4.1) and (4.2) show that the sets A and B are absolutely linearly separable. If two sets are linearly separable in the absolute sense, then they are, of course, linearly separable in the conventional sense. \square

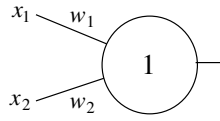
4.1.4 The error surface and the search method

A usual approach for starting the learning algorithm is to initialize the network weights randomly and to improve these initial parameters, looking at each step to see whether a better separation of the training set can be achieved. In this section we identify points (x_1, x_2, \dots, x_n) in n -dimensional space with the vector \mathbf{x} with the same coordinates.

Definition 4. *The open (closed) positive half-space associated with the n -dimensional weight vector \mathbf{w} is the set of all points $\mathbf{x} \in \mathbb{R}^n$ for which $\mathbf{w} \cdot \mathbf{x} > 0$ ($\mathbf{w} \cdot \mathbf{x} \geq 0$). The open (closed) negative half-space associated with \mathbf{w} is the set of all points $\mathbf{x} \in \mathbb{R}^n$ for which $\mathbf{w} \cdot \mathbf{x} < 0$ ($\mathbf{w} \cdot \mathbf{x} \leq 0$).*

We omit the adjectives “closed” or “open” whenever it is clear from the context which kind of linear separation is being used.

Let P and N stand for two finite sets of points in \mathbb{R}^n which we want to separate linearly. A weight vector is sought so that the points in P belong to its associated positive half-space and the points in N to the negative half-space. The *error* of a perceptron with weight vector \mathbf{w} is the number of incorrectly classified points. The learning algorithm must minimize this error function $E(\mathbf{w})$. One possible strategy is to use a local greedy algorithm which works by computing the error of the perceptron for a given weight vector, looking then for a direction in weight space in which to move, and updating the weight vector by selecting new weights in the selected search direction. We can visualize this strategy by looking at its effect in weight space.

**Fig. 4.4.** Perceptron with constant threshold

Let us take as an example a perceptron with constant threshold $\theta = 1$ (Figure 4.4). We are looking for two weights, w_1 and w_2 , which transform the perceptron into a binary AND gate. We can show graphically the error function for all combinations of the two variable weights. This has been done in Figure 4.5 for values of the weights between -0.5 and 1.5 . The solution region is the triangular area in the middle. The learning algorithm should reach this region starting from any other point in weight space. In this case, it is possible to descend from one surface to the next using purely local decisions at each point.

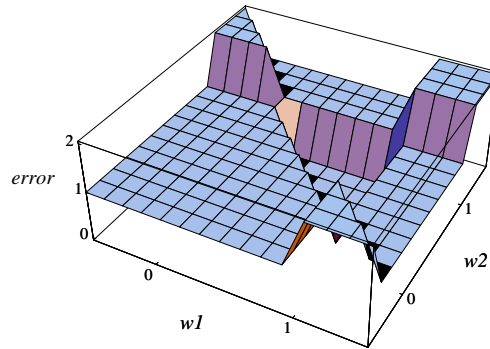
**Fig. 4.5.** Error function for the AND function

Figure 4.6 shows the different regions of the error function as seen from “above”. The solution region is a triangle with error level 0. For the other regions the diagram shows their corresponding error count. The figure illustrates an iterative search process that starts at \mathbf{w}_0 , goes through \mathbf{w}_1 , \mathbf{w}_2 , and finally reaches the solution region at \mathbf{w}^* . Later on, this visualization will help us to understand the computational complexity of the perceptron learning algorithm.

The optimization problem we are trying to solve can be understood as descent on the error surface but also as a search for an inner point of the solution region. Let $N = \{(0, 0), (1, 0), (0, 1)\}$ and $P = \{(1, 1)\}$ be two sets of points to be separated absolutely. The set P must be classified in the positive

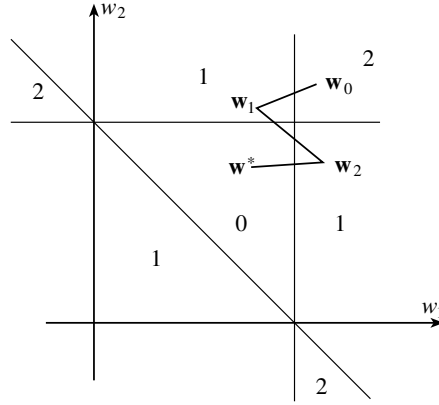


Fig. 4.6. Iteration steps to the region of minimal error

and the set N in the negative half-space. This is the separation corresponding to the AND function.

Three weights w_1, w_2 , and $w_3 = -\theta$ are needed to implement the desired separation with a generic perceptron. The first step is to extend the input vectors with a third coordinate $x_3 = 1$ and to write down the four inequalities that must be fulfilled:

$$(0, 0, 1) \cdot (w_1, w_2, w_3) < 0 \quad (4.3)$$

$$(1, 0, 1) \cdot (w_1, w_2, w_3) < 0 \quad (4.4)$$

$$(0, 1, 1) \cdot (w_1, w_2, w_3) < 0 \quad (4.5)$$

$$(1, 1, 1) \cdot (w_1, w_2, w_3) > 0 \quad (4.6)$$

These equations can be written in the following simpler matrix form:

$$\begin{pmatrix} 0 & 0 & -1 \\ -1 & 0 & -1 \\ 0 & -1 & -1 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \\ w_3 \end{pmatrix} > \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}. \quad (4.7)$$

This can be written as

$$\mathbf{A}\mathbf{w} > \mathbf{0}.$$

where \mathbf{A} is the 4×3 matrix of Equation (4.7) and \mathbf{w} the weight vector (written as a column vector). The learning problem is to find the appropriate weight vector \mathbf{w} .

Equation (4.7) describes all points in the interior of a convex polytope. The sides of the polytope are delimited by the planes defined by each of the inequalities (4.3)–(4.6). Any point in the interior of the polytope represents a solution for the learning problem.

We saw that the solution region for the AND function has a triangular shape when the threshold is fixed at 1. In Figure 4.7 we have a three-dimensional view of the whole solution region when the threshold (i.e., w_3) is allowed to change. The solution region of Figure 4.5 is just a cut of the solution polytope of Figure 4.7 at $w_3 = -1$. The shaded surface represents the present cut, which is similar to any other cut we could make to the polytope for different values of the threshold.

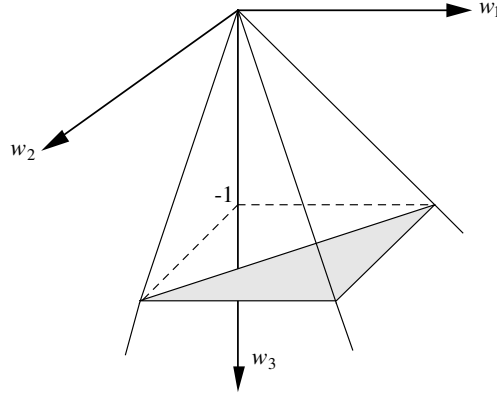


Fig. 4.7. Solution polytope for the AND function in weight space

We can see that the polytope is unbounded in the direction of negative w_3 . This means that the absolute value of the threshold can become as large as we want and we will still find appropriate combinations of w_1 and w_2 to compute the AND function. The fact that we have to look for interior points of polytopes for the solution of the learning problem, is an indication that linear programming methods could be used. We will elaborate on this idea later on.

4.2 Algorithmic learning

We are now in a position to introduce the perceptron learning algorithm. The training set consists of two sets, P and N , in n -dimensional extended input space. We look for a vector \mathbf{w} capable of absolutely separating both sets, so that all vectors in P belong to the open positive half-space and all vectors in N to the open negative half-space of the linear separation.

Algorithm 4.2.1 *Perceptron learning*

start: The weight vector \mathbf{w}_0 is generated randomly,
 set $t := 0$

test: A vector $\mathbf{x} \in P \cup N$ is selected randomly,
 if $\mathbf{x} \in P$ and $\mathbf{w}_t \cdot \mathbf{x} > 0$ go to *test*,
 if $\mathbf{x} \in P$ and $\mathbf{w}_t \cdot \mathbf{x} \leq 0$ go to *add*,
 if $\mathbf{x} \in N$ and $\mathbf{w}_t \cdot \mathbf{x} < 0$ go to *test*,
 if $\mathbf{x} \in N$ and $\mathbf{w}_t \cdot \mathbf{x} \geq 0$ go to *subtract*.

add: set $\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{x}$ and $t := t + 1$, goto *test*

subtract: set $\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{x}$ and $t := t + 1$, goto *test*

This algorithm [312] makes a correction to the weight vector whenever one of the selected vectors in P or N has not been classified correctly. The perceptron convergence theorem guarantees that if the two sets P and N are linearly separable the vector \mathbf{w} is updated only a finite number of times. The routine can be stopped when all vectors are classified correctly. The corresponding test must be introduced in the above pseudocode to make it stop and to transform it into a fully-fledged algorithm.

4.2.1 Geometric visualization

There are two alternative ways to visualize perceptron learning, one more effective than the other. Given the two sets of points $P \in \mathbb{R}^2$ and $N \in \mathbb{R}^2$ to be separated, we can visualize the linear separation in input space, as in Figure 4.8, or in extended input space. In the latter case we extend the input vectors and look for a linear separation through the origin, that is, a plane with equation $w_1x_1 + w_2x_2 + w_3 = 0$. The vector normal to this plane is the weight vector $\mathbf{w} = (w_1, w_2, w_3)$. Figure 4.9 illustrates this approach.

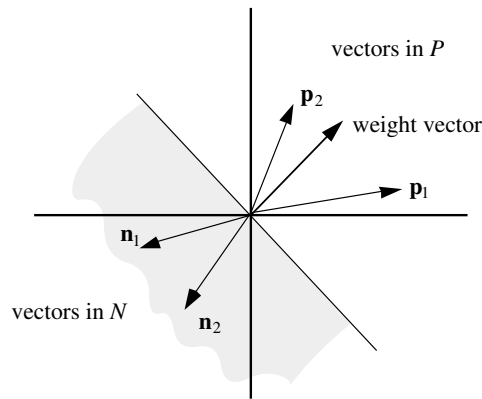


Fig. 4.8. Visualization in input space

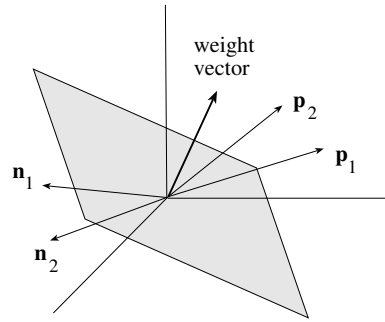


Fig. 4.9. Visualization in extended input space

We are thus looking for a weight vector \mathbf{w} with a positive scalar product with all the extended vectors represented by the points in P and with a negative product with the extended vectors represented by the points in N . Actually, we will deal with this problem in a more general way. Assume that P and N are sets of n -dimensional vectors and a weight vector \mathbf{w} must be found, such that $\mathbf{w} \cdot \mathbf{x} > 0$ holds for all $x \in P$ and $\mathbf{w} \cdot \mathbf{x} < 0$ holds for all $x \in N$.

The perceptron learning algorithm starts with a randomly chosen vector \mathbf{w}_0 . If a vector $\mathbf{x} \in P$ is found such that $\mathbf{w} \cdot \mathbf{x} < 0$, this means that the angle between the two vectors is greater than 90 degrees. The weight vector must be rotated in the direction of \mathbf{x} to bring this vector into the positive half-space defined by \mathbf{w} . This can be done by adding \mathbf{w} and \mathbf{x} , as the perceptron learning algorithm does. If $\mathbf{x} \in N$ and $\mathbf{w} \cdot \mathbf{x} > 0$, then the angle between \mathbf{x} and \mathbf{w} is less than 90 degrees. The weight vector must be rotated away from \mathbf{x} . This is done by subtracting \mathbf{x} from \mathbf{w} . The vectors in P rotate the weight vector in one direction, the vectors in N rotate the negative weight vector in another. If a solution exists it can be found after a finite number of steps. A good initial heuristic is to start with the average of the positive input vectors minus the average of the negative input vectors. In many cases this yields an initial vector near the solution region.

In perceptron learning we are not necessarily dealing with normalized vectors, so that every update of the weight vector of the form $\mathbf{w} \pm \mathbf{x}$ rotates the weight vector by a different angle. If $\mathbf{x} \in P$ and $\|\mathbf{x}\| \gg \|\mathbf{w}\|$ the new weight vector $\mathbf{w} + \mathbf{x}$ is almost equal to \mathbf{x} . This effect and the way perceptron learning works can be seen in Figure 4.10. The initial weight vector is updated by adding \mathbf{x}_1 , \mathbf{x}_3 , and \mathbf{x}_1 again to it. After each correction the weight vector is rotated in one or the other direction. It can be seen that the vector \mathbf{w} becomes larger after each correction in this example. Each correction rotates the weight vector by a smaller angle until the correct linear separation has been found. After the initial updates, successive corrections become smaller and the algorithm “fine tunes” the position of the weight vector. The learning rate, the rate of change of the vector \mathbf{w} , becomes smaller in time and if we

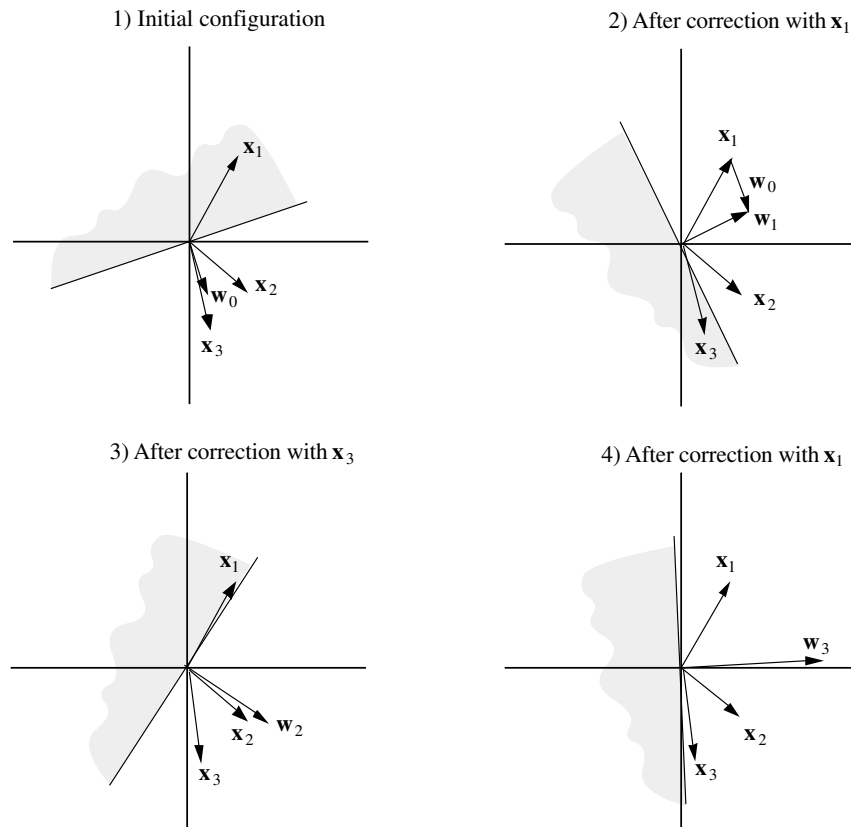


Fig. 4.10. Convergence behavior of the learning algorithm

keep on training, even after the vectors have already been correctly separated, it approaches zero. Intuitively we can think that the learned vectors are increasing the “inertia” of the weight vector. Vectors lying just outside of the positive region are brought into it by rotating the weight vector just enough to correct the error.

This is a typical feature of many learning algorithms for neural networks. They make use of a so-called *learning constant*, which is brought to zero during the learning process to consolidate what has been learned. The perceptron learning algorithm provides a kind of automatic learning constant which determines the degree of adaptivity (the so-called *plasticity* of the network) of the weights.

4.2.2 Convergence of the algorithm

The convergence proof of the perceptron learning algorithm is easier to follow by keeping in mind the visualization discussed in the previous section.

Proposition 8. *If the sets P and N are finite and linearly separable, the perceptron learning algorithm updates the weight vector \mathbf{w}_t a finite number of times. In other words: if the vectors in P and N are tested cyclically one after the other, a weight vector \mathbf{w}_t is found after a finite number of steps t which can separate the two sets.*

Proof. We can make three simplifications without losing generality:

- i) The sets P and N can be joined in a single set $P' = P \cup N^-$, where the set N^- consists of the negated elements of N .
- ii) The vectors in P' can be normalized, because if a weight vector \mathbf{w} is found so that $\mathbf{w} \cdot \mathbf{x} > 0$ this is also valid for any other vector $\eta\mathbf{x}$, where η is a constant.
- iii) The weight vector can also be normalized. Since we assume that a solution for the linear separation problem exists, we call \mathbf{w}^* a normalized solution vector.

Assume that after $t + 1$ steps the weight vector \mathbf{w}_{t+1} has been computed. This means that at time t a vector \mathbf{p}_i was incorrectly classified by the weight vector \mathbf{w}_t and so $\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{p}_i$.

The cosine of the angle ρ between \mathbf{w}_{t+1} and \mathbf{w}^* is

$$\cos \rho = \frac{\mathbf{w}^* \cdot \mathbf{w}_{t+1}}{\|\mathbf{w}_{t+1}\|} \quad (4.8)$$

For the expression in the numerator we know that

$$\begin{aligned} \mathbf{w}^* \cdot \mathbf{w}_{t+1} &= \mathbf{w}^* \cdot (\mathbf{w}_t + \mathbf{p}_i) \\ &= \mathbf{w}^* \cdot \mathbf{w}_t + \mathbf{w}^* \cdot \mathbf{p}_i \\ &\geq \mathbf{w}^* \cdot \mathbf{w}_t + \delta \end{aligned}$$

with $\delta = \min\{\mathbf{w}^* \cdot \mathbf{p} \mid \forall \mathbf{p} \in P'\}$. Since the weight vector \mathbf{w}^* defines an absolute linear separation of P and N we know that $\delta > 0$. By induction we obtain

$$\mathbf{w}^* \cdot \mathbf{w}_{t+1} \geq \mathbf{w}^* \cdot \mathbf{w}_0 + (t + 1)\delta. \quad (4.9)$$

On the other hand for the term in the denominator of (4.8) we know that

$$\begin{aligned} \|\mathbf{w}_{t+1}\|^2 &= (\mathbf{w}_t + \mathbf{p}_i) \cdot (\mathbf{w}_t + \mathbf{p}_i) \\ &= \|\mathbf{w}_t\|^2 + 2\mathbf{w}_t \cdot \mathbf{p}_i + \|\mathbf{p}_i\|^2 \end{aligned}$$

Since $\mathbf{w}_t \cdot \mathbf{p}_i$ is negative or zero (otherwise we would have not corrected \mathbf{w}_t using \mathbf{p}_i) we can deduce that

$$\begin{aligned} \|\mathbf{w}_{t+1}\|^2 &\leq \|\mathbf{w}_t\|^2 + \|\mathbf{p}_i\|^2 \\ &\leq \|\mathbf{w}_t\|^2 + 1 \end{aligned}$$

because all vectors in P have been normalized. Induction then gives us

$$\|\mathbf{w}_{t+1}\|^2 \leq \|\mathbf{w}_0\|^2 + (t+1). \quad (4.10)$$

From (4.9) and (4.10) and Equation (4.8) we get the inequality

$$\cos \rho \geq \frac{\mathbf{w}^* \cdot \mathbf{w}_0 + (t+1)\delta}{\sqrt{\|\mathbf{w}_0\|^2 + (t+1)}}$$

The right term grows proportionally to \sqrt{t} and, since δ is positive, it can become arbitrarily large. However, since $\cos \rho \leq 1$, t must be bounded by a maximum value. Therefore, the number of corrections to the weight vector must be finite. \square

The proof shows that perceptron learning works by bringing the initial vector \mathbf{w}_0 sufficiently close to \mathbf{w}^* (since $\cos \rho$ becomes larger and ρ proportionately smaller).

4.2.3 Accelerating convergence

Although the perceptron learning algorithm converges to a solution, the number of iterations can be very large if the input vectors are not normalized and are arranged in an unfavorable way.

There are faster methods to find the weight vector appropriate for a given problem. When the perceptron learning algorithm makes a correction, an input vector \mathbf{x} is added or subtracted from the weight vector \mathbf{w} . The search direction is given by the vector \mathbf{x} . Each input vector corresponds to the border of one region of the error function defined on weight space. The direction of \mathbf{x} is orthogonal to the step defined by \mathbf{x} on the error surface. The weight vector is displaced in the direction of \mathbf{x} until it “falls” into a region with smaller error.

We can illustrate the dynamics of perceptron learning using the error surface for the OR function as an example. The input $(1, 1)$ must produce the output 1 (for simplicity we fix the threshold of the perceptron to 1). The two weights w_1 and w_2 must fulfill the inequality $w_1 + w_2 \geq 1$. Any other combination produces an error. The contribution to the total error is shown in Figure 4.11 as a step in the error surface. If the initial weight vector lies in the triangular region with error 1, it must be brought up to the verge of the region with error 0. This can be done by adding the vector $(1, 1)$ to \mathbf{w} . However, if the input vector is, for example, $(0.1, 0.1)$, it should be added a few times before the weight combination (w_1, w_2) falls to the region of error 0. In this case we would like to make the correction in a single iteration.

These considerations provide an improvement for the perceptron learning algorithm: if at iteration t the input vector $\mathbf{x} \in P$ is classified erroneously, then we have $\mathbf{w}_t \cdot \mathbf{x} \leq 0$. The error δ can be defined as

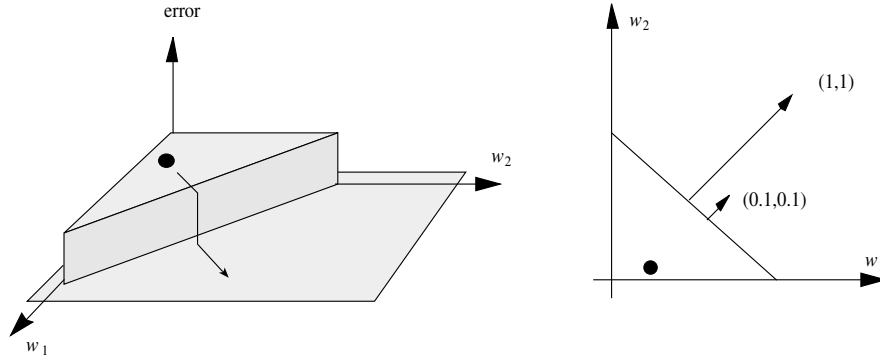


Fig. 4.11. A step on the error surface

$$\delta = -\mathbf{w}_t \cdot \mathbf{x}.$$

The new weight vector \mathbf{w}_{t+1} is calculated as follows:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \frac{\delta + \varepsilon}{\|\mathbf{x}\|^2} \mathbf{x},$$

where ε denotes a small positive real number. The classification of \mathbf{x} has been corrected in one step because

$$\begin{aligned} \mathbf{w}_{t+1} \cdot \mathbf{x} &= (\mathbf{w}_t + \frac{\delta + \varepsilon}{\|\mathbf{x}\|^2} \mathbf{x}) \cdot \mathbf{x} \\ &= \mathbf{w}_t \cdot \mathbf{x} + (\delta + \varepsilon) \\ &= -\delta + \delta + \varepsilon \\ &= \varepsilon > 0 \end{aligned}$$

The number ε guarantees that the new weight vector just barely skips over the border of the region with a higher error. The constant ε should be made small enough to avoid skipping to another region whose error is higher than the current one. When $\mathbf{x} \in N$ the correction step is made similarly, but using the factor $\delta - \varepsilon$ instead of $\delta + \varepsilon$.

The accelerated algorithm is an example of *corrective learning*: We do not just “reinforce” the weight vector, but completely correct the error that has been made. A variant of this rule is correction of the weight vector using a proportionality constant γ as the learning factor, in such a way that at each update the vector $\gamma(\delta + \varepsilon)\mathbf{x}$ is added to \mathbf{w} . The learning constant falls to zero as learning progresses.

4.2.4 The pocket algorithm

If the learning set is not linearly separable the perceptron learning algorithm does not terminate. However, in many cases in which there is no perfect linear

separation, we would like to compute the linear separation which correctly classifies the largest number of vectors in the positive set P and the negative set N . Gallant proposed a very simple variant of the perceptron learning algorithm capable of computing a good approximation to this ideal linear separation. The main idea of the algorithm is to store the best weight vector found so far by perceptron learning (in a “pocket”) while continuing to update the weight vector itself. If a better weight vector is found, it supersedes the one currently stored and the algorithm continues to run [152].

Algorithm 4.2.2 *Pocket algorithm*

- start:* Initialize the weight vector \mathbf{w} randomly. Define a “stored” weight vector $\mathbf{w}_s = \mathbf{w}$. Set h_s , the history of \mathbf{w}_s , to zero.
- iterate:* Update \mathbf{w} using a single iteration of the perceptron learning algorithm. Keep track of the number h of consecutively successfully tested vectors. If at any moment $h > h_s$, substitute \mathbf{w}_s with \mathbf{w} and h_s with h . Continue iterating.

The algorithm can occasionally change a good stored weight vector for an inferior one, since only information from the last run of selected examples is considered. The probability of this happening, however, becomes smaller and smaller as the number of iterations grows. If the training set is finite and the weights and vectors are rational, it can be shown that this algorithm converges to an optimal solution with probability 1 [152].

4.2.5 Complexity of perceptron learning

The perceptron learning algorithm selects a search direction in weight space according to the incorrect classification of the last tested vector and does not make use of global information about the shape of the error function. It is a greedy, local algorithm. This can lead to an exponential number of updates of the weight vector.

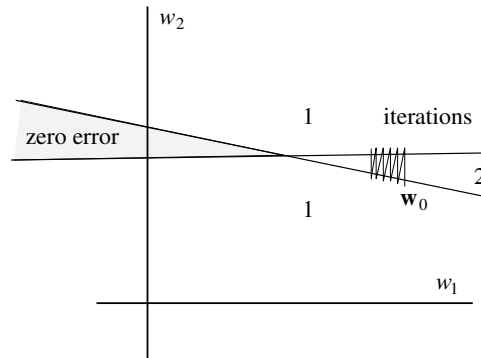


Fig. 4.12. Worst case for perceptron learning (weight space)

Figure 4.12 shows the different error regions in a worst case scenario. The region with error 0 is bounded by two lines which meet at a small angle. Starting the learning algorithm at point \mathbf{w}_0 , the weight updates will lead to a search path similar to the one shown in the figure. In each new iteration a new weight vector is computed, in such a way that one of two vectors is classified

correctly. However, each of these corrections leads to the other vector being incorrectly classified. The iteration jumps from one region with error 1 to the other one. The algorithm converges only after a certain number of iterations, which can be made arbitrarily large by adjusting the angle at which the lines meet.

Figure 4.12 corresponds to the case in which two almost antiparallel vectors are to be classified in the same half-space (Figure 4.13). An algorithm which rotates the separation line in one of the two directions (like perceptron learning) will require more and more time when the angle between the two vectors approaches 180 degrees.

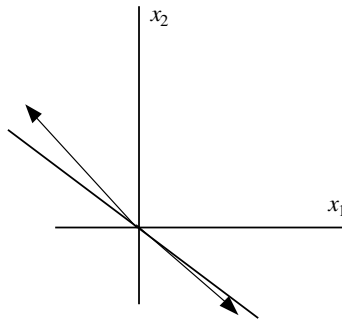


Fig. 4.13. Worst case for perceptron learning (input space)

This example is a good illustration of the advantages of visualizing learning algorithms in both the input space and its dual, weight space. Figure 4.13 shows the concrete problem and Figure 4.12 illustrates why it is difficult to solve.

4.3 Linear programming

A set of input vectors to be separated by a perceptron in a positive and a negative set defines a convex polytope in weight space, whose inner points represent all admissible weight combinations for the perceptron. The perceptron learning algorithm finds a solution when the interior of the polytope is not void. Stated differently: if we want to train perceptrons to classify patterns, we must solve an inner point problem. Linear programming can deal with this kind of task.

4.3.1 Inner points of polytopes

Linear programming was developed to solve the following generic problem: Given a set of n variables x_1, x_2, \dots, x_n a function $c_1x_1 + c_2x_2 + \dots + c_nx_n$ must

be maximized (or minimized). The variables must obey certain constraints given by linear inequalities of the form

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &\leq b_1 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &\leq b_2 \\ &\vdots \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n &\leq b_m \end{aligned}$$

All m linear constraints can be summarized in the matrix inequality $\mathbf{Ax} \leq \mathbf{b}$, in which \mathbf{x} and \mathbf{b} respectively represent n -dimensional and m -dimensional column vectors and \mathbf{A} is a $m \times n$ matrix. It is also necessary that $\mathbf{x} \geq \mathbf{0}$, which can always be guaranteed by introducing additional variables.

As in the case of a perceptron, the m inequalities define a convex polytope of feasible values for the variables x_1, x_2, \dots, x_n . If the optimization problem has a solution, this is found at one of the vertices of the polytope. Figure 4.14 shows a two-dimensional example. The shaded polygon is the feasible region. The function to be optimized is represented by the line normal to the vector \mathbf{c} . Finding the point where this linear function reaches its maximum corresponds to moving the line, without tilting it, up to the farthest position at which it is still in contact with the feasible region, in our case ξ . It is intuitively clear that when one or more solutions exist, one of the vertices of the polytope is one of them.

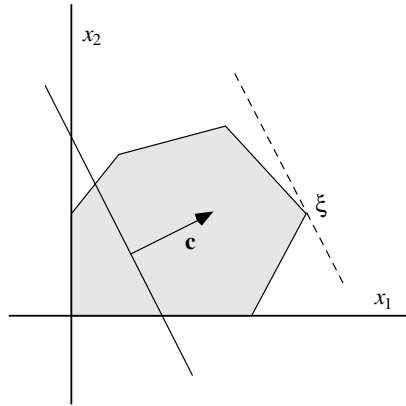


Fig. 4.14. Feasible region for a linear optimization problem

The well-known simplex algorithm of linear programming starts at a vertex of the feasible region and jumps to another neighboring vertex, always moving in a direction in which the function to be optimized increases. In the worst case an exponential number of vertices in the number of inequalities m has to be traversed before a solution is found. On average, however, the

simplex algorithm is not so inefficient. In the case of Figure 4.14 the optimal solution can be found in two steps by starting at the origin and moving to the right. To determine the next vertex to be visited, the simplex algorithm uses as a criterion the length of the projection of the gradient of the function to be optimized on the edges of the polytope. It is in this sense a gradient algorithm. The algorithm can be inefficient because the search for the optimum is constrained to be carried out moving only along the edges of the polytope. If the number of delimiting surfaces is large, a better alternative is to go right through the middle of the polytope.

4.3.2 Linear separability as linear optimization

The simplex algorithm and its variants need to start at a point in the feasible region. In many cases it can be arranged to start at the origin. If the feasible region does not contain the origin as one of its vertices, a feasible point must be found first. This problem can be transformed into a linear program.

Let \mathbf{A} represent the $m \times n$ matrix of coefficients of the linear constraints and \mathbf{b} an m -dimensional column vector. Assume that we are looking for an n -dimensional column vector \mathbf{x} such that $\mathbf{Ax} \leq \mathbf{b}$. This condition is fulfilled only by points in the feasible region. To simplify the problem, assume that $\mathbf{b} \geq \mathbf{0}$ and that we are looking for vectors $\mathbf{x} \geq \mathbf{0}$. Introducing the column vector \mathbf{y} of m additional slack variables (y_1, \dots, y_m) , the inequality $\mathbf{Ax} \leq \mathbf{b}$ can be transformed into the equality $\mathbf{Ax} + \mathbf{Iy} = \mathbf{b}$, where \mathbf{I} denotes the $m \times m$ identity matrix. The linear program to be solved is then

$$\min \left\{ \sum_{i=1}^m y_i \mid \mathbf{Ax} + \mathbf{Iy} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}, \mathbf{y} \geq \mathbf{0} \right\}.$$

An initial feasible solution for the problem is $\mathbf{x} = \mathbf{0}$ and $\mathbf{y} = \mathbf{b}$. Starting from here an iterative algorithm looks for the minimum of the sum of the slack variables. If the minimum is negative the original problem does not have a solution and the feasible region of $\mathbf{Ax} \leq \mathbf{b}$ is void. If the minimum is zero, the value of \mathbf{x} determined during the optimization is an inner point of the feasible region (more exactly, a point at its boundary).

The conditions $\mathbf{x} \geq \mathbf{0}$ and $\mathbf{b} \geq \mathbf{0}$ can be omitted and additional transformations help to transform the more general problem to the canonical form discussed here [157].

Inner points of convex polytopes, defined by separating hyperplanes, can thus be found using linear programming algorithms. Since the computation of the weight vector for a perceptron corresponds to the computation of inner points of convex polytopes, this means that perceptron learning can also be handled in this way. If two sets of vectors are not linearly separable, the linear programming algorithm can detect it. The complexity of linearly separating points in an input space is thus bounded by the complexity of solving linear programming problems.

The perceptron learning algorithm is not the most efficient method for perceptron learning, since the number of steps can grow exponentially in the worst case. In the case of linear programming, theoreticians have succeeded in crafting algorithms which need a polynomial number of iterations and return the optimal solution or an indication that it does not exist.

4.3.3 Karmarkar's algorithm

In 1984 a fast polynomial time algorithm for linear programming was proposed by Karmarkar [236]. His algorithm starts at an inner point of the solution region and proceeds in the direction of steepest ascent (if maximizing), taking care not to step out of the feasible region.

Figure 4.15 schematically shows how the algorithm works. The algorithm starts with a canonical form of the linear programming problem in which the additional constraint $x_1 + x_2 + \dots + x_n = 1$ is added to the basic constraints $\mathbf{Ax} \geq \mathbf{0}$, where x_1, \dots, x_n are the variables in the problem. Some simple transformations can bring the original problem into this form. The point $\mathbf{e} = \frac{1}{n}(1, 1, \dots, 1)$ is considered the middle of the solution polytope and each iteration step tries to transform the original problem in such a way that this is always the starting point.

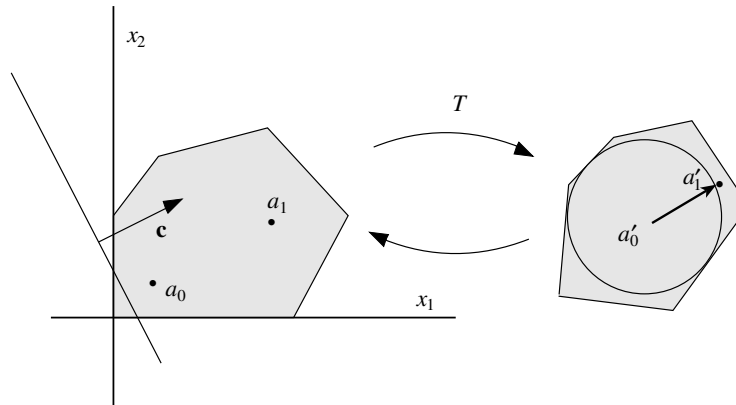


Fig. 4.15. Transformation of the solution polytope

An initial point a_0 is selected in the interior of the solution polytope and then brought into the middle \mathbf{e} of the transformed feasible region using a projective transformation T . A projective transformation maps each point \mathbf{x} in the hyperplane $x_1 + x_2 + \dots + x_n = 1$ to a point \mathbf{x}' in another hyperplane, whereby the line joining \mathbf{x} and \mathbf{x}' goes through a predetermined point p . The transformation is applied on the initial point a_0 , the matrix \mathbf{A} of linear constraints and also to the linear function $\mathbf{c}^T \mathbf{x}$ to be optimized. After the

transformation, the radius of the largest sphere with center a'_0 and inside the transformed feasible region is computed. Starting at the center of the sphere a new point in the direction of the transformed optimizing direction \mathbf{c}' is computed. The step length is made shorter than the computed maximal radius by a small factor, to avoid reaching the surface of the solution polytope. The new point a'_1 computed in this way is a feasible point and is also strictly in the interior of the solution polytope. The point a'_1 is transformed back to the original space using the inverse projective transformation T^{-1} and a new iteration can start again from this point. This basic step is repeated, periodically testing whether a vertex of the polytope is close enough and optimal. At this moment the algorithm stops and takes this vertex as the solution of the problem (Figure 4.16). Additionally, a certain checking must be done in each iteration to confirm that a solution to the optimization problem exists and that the cost function is not unbounded.

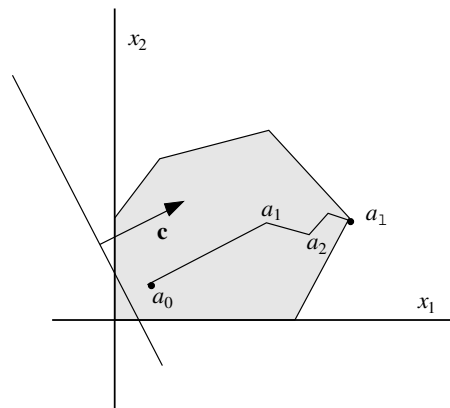


Fig. 4.16. Example of a search path for Karmarkar's algorithm

In the worst case Karmarkar's algorithm executes in a number of iterations proportional to $n^{3.5}$, where n is the number of variables in the problem and other factors are kept constant. Some published modifications of Karmarkar's algorithm are still more efficient but start beating the simplex method in the average case only when the number of variables and constraints becomes relatively large, since the computational overhead for a small number of constraints is not negligible [257].

The existence of a polynomial time algorithm for linear programming and for the solution of interior point problems shows that perceptron learning is not what is called a *hard* computational problem. Given any number of training patterns, the learning algorithm (in this case linear programming) can decide whether the problem has a solution or not. If a solution exists, it finds the appropriate weights in polynomial time at most.

4.4 Historical and bibliographical remarks

The success of the perceptron and the interest it aroused in the 1960s was a direct product of its learning capabilities, different from the hand-design approach of previous models. Later on, research in this field reached an impasse when a learning algorithm for more general networks was still unavailable.

Minsky and Papert [312] analyzed the features and limitations of the perceptron model in a rigorous way. They could show that the perceptron learning algorithm needs an exponential number of learning steps in the worst case. However, perceptron learning is in the average case fairly efficient. Mansfield showed that when the training set is selected randomly from a half-space, the number of iterations of the perceptron learning algorithm is comparable to the number of iterations needed by ellipsoid methods for linear programming (up to dimension 30) [290]. Baum had previously shown that when the learning set is picked by a non-malicious adversary, the complexity of perceptron learning is polynomial [46].

More recently the question has arisen of whether a given set of nonlinearly separable patterns can be decomposed in such a way that the largest linearly separable subset can be detected. Amaldi showed that this is an *NP*-complete problem, that is, a problem for which presumably no polynomial time algorithm exists (compare Chap. 10).

The conditions for perfect perceptron learning can be also relaxed. If the set of patterns is not linearly separable, we can look for the separation that minimizes the average quadratic error, without requiring it to be zero. In this case statistical methods or the backpropagation algorithm (Chap. 7) can be used.

After the invention of the simplex algorithm for linear programming there was a general feeling that it could be proven that one of its variants was of polynomial complexity in the number of constraints and of variables. This was due to the fact that the actual experience with the algorithm showed that in the average case a solution was found in much less than exponential time. However, in 1972 Klee and Minty [247] gave a counterexample which showed that there were situations in which the simplex method visited 2^{n-1} vertices of a feasible region with 2^n vertices. Later it was rigorously proven that the simplex method is polynomial in the average case [64]. The question of the existence of a polynomial time algorithm for linear programming was settled by Khachiyan in 1979, when he showed that a recursive construction of ellipsoids could lead to finding the optimal vertex of the feasible region in polynomial time [244]. His algorithm, however, was very computationally intensive for most of the average-sized problems and could not displace the simplex method. Karmarkar's algorithm, a further development of the ellipsoid method including some very clever transformations, aroused much interest when it was first introduced in 1984. So many variations of the original algorithm have appeared that they are collectively known as Karmarkar-type algorithms. Minimization problems with thousands of constraints can now be

dealt with efficiently by these polynomial time algorithms, but since the simplex method is fast in the average case it continues to be the method of choice in medium-sized problems.

Interesting variations of perceptron learning were investigated by Fontanari and Meir, who coded the different alternatives of weight updates according to the local information available to each weight and let a population of algorithms evolve. With this kind of “evolution strategy” they found competitive algorithms similar to the standard methods [140].

Exercises

1. Implement the perceptron learning algorithm in the computer. Find the weights for an edge detection operator using this program. The input-output examples can be taken from a digitized picture of an object and another one in which only the edges of the object have been kept.
2. Give a numerical example of a training set that leads to many iterations of the perceptron learning algorithm.
3. How many vectors can we pick randomly in an n -dimensional space so that they all belong to the same half-space? Produce a numerical estimate using a computer program.
4. The perceptron learning algorithm is usually fast if the vectors to be linearly separated are chosen randomly. Choose a weight vector \mathbf{w} for a perceptron randomly. Generate p points in input space and classify them in a positive or negative class according to their scalar product with \mathbf{w} . Now train a perceptron using this training set and measure the number of iterations needed. Make a plot of n against p for dimension up to 10 and up to 100 points.



