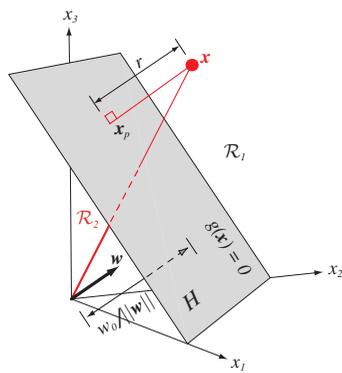
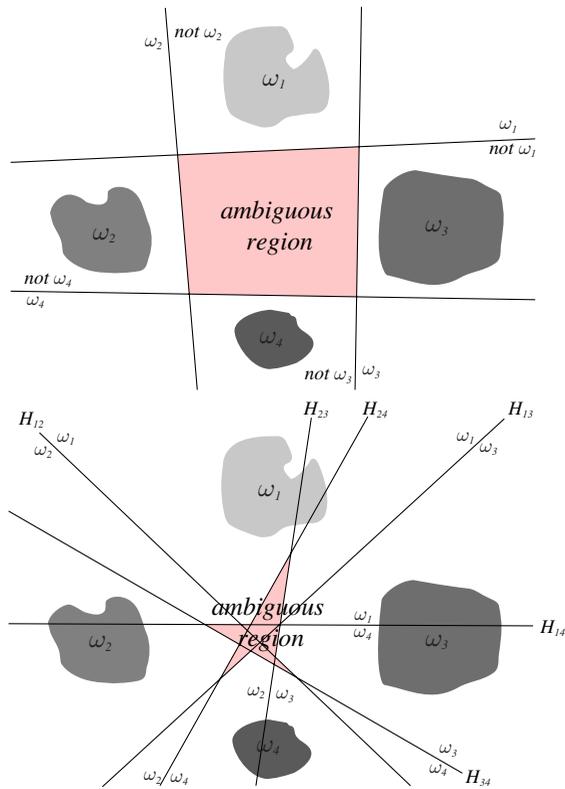


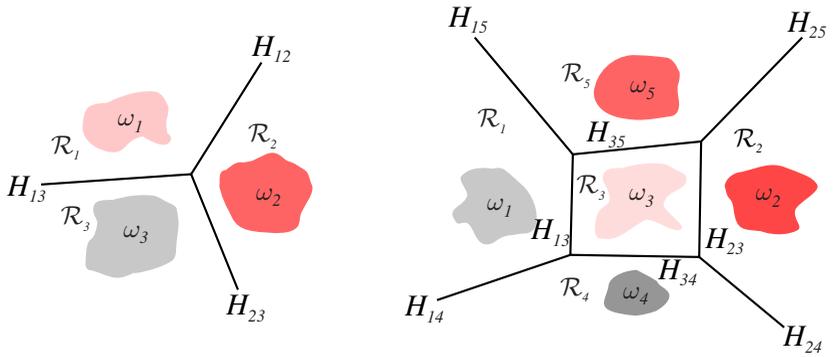
**FIGURE 5.1.** A simple linear classifier having  $d$  input units, each corresponding to the values of the components of an input vector. Each input feature value  $x_i$  is multiplied by its corresponding weight  $w_i$ ; the effective input at the output unit is the sum all these products,  $\sum w_i x_i$ . We show in each unit its effective input-output function. Thus each of the  $d$  input units is linear, emitting exactly the value of its corresponding feature value. The single bias unit unit always emits the constant value 1.0. The single output unit emits a +1 if  $\mathbf{w}'\mathbf{x} + w_0 > 0$  or a -1 otherwise. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



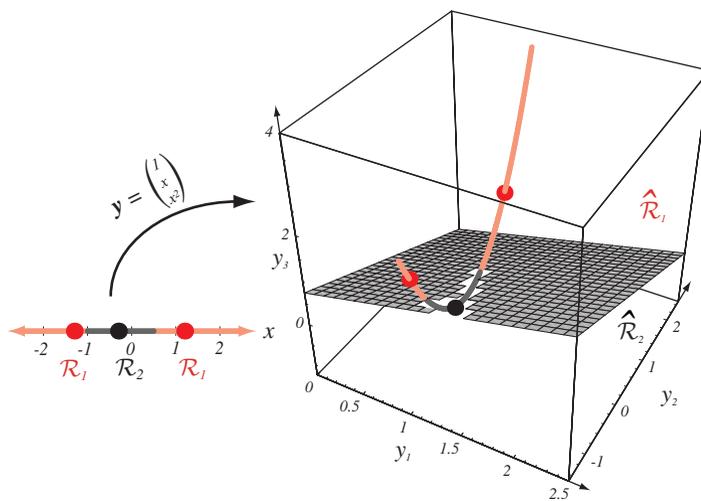
**FIGURE 5.2.** The linear decision boundary  $H$ , where  $g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0 = 0$ , separates the feature space into two half-spaces  $\mathcal{R}_1$  (where  $g(\mathbf{x}) > 0$ ) and  $\mathcal{R}_2$  (where  $g(\mathbf{x}) < 0$ ). From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



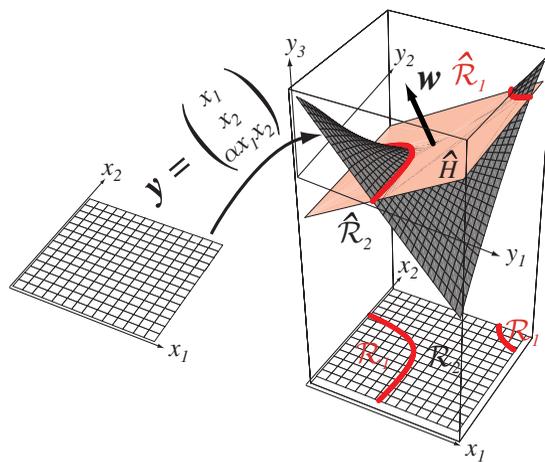
**FIGURE 5.3.** Linear decision boundaries for a four-class problem. The top figure shows  $\omega_i/\text{not } \omega_i$  dichotomies while the bottom figure shows  $\omega_i/\omega_j$  dichotomies and the corresponding decision boundaries  $H_{ij}$ . The pink regions have ambiguous category assignments. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



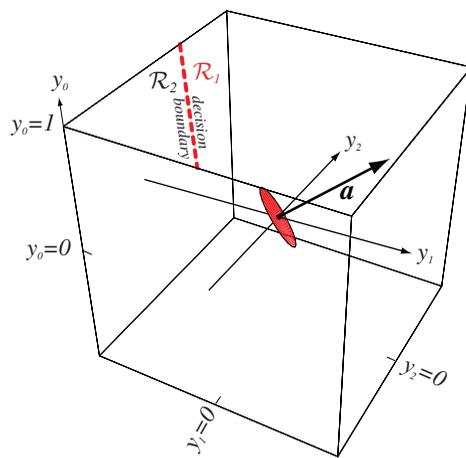
**FIGURE 5.4.** Decision boundaries produced by a linear machine for a three-class problem and a five-class problem. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



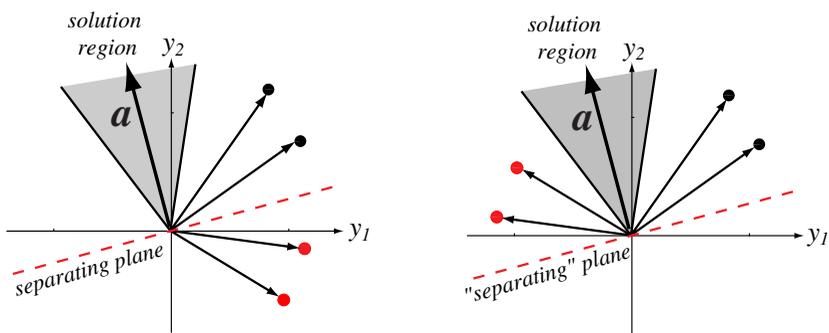
**FIGURE 5.5.** The mapping  $y = (1, x, x^2)^t$  takes a line and transforms it to a parabola in three dimensions. A plane splits the resulting y-space into regions corresponding to two categories, and this in turn gives a nonsimply connected decision region in the one-dimensional x-space. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



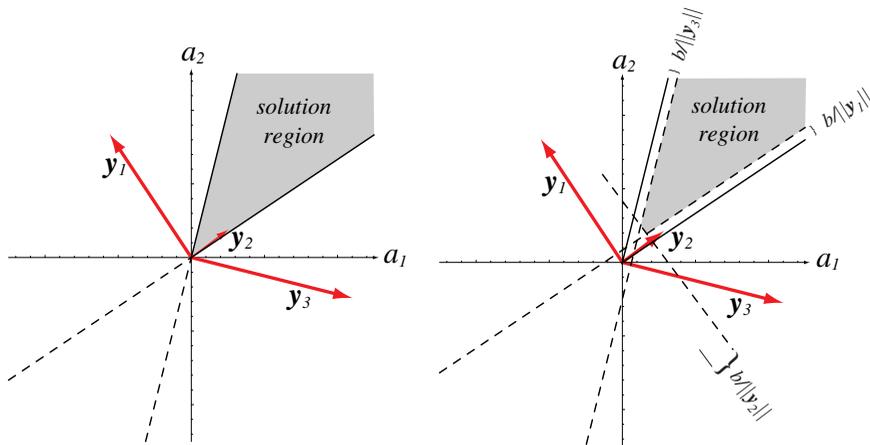
**FIGURE 5.6.** The two-dimensional input space  $\mathbf{x}$  is mapped through a polynomial function  $f$  to  $\mathbf{y}$ . Here the mapping is  $y_1 = x_1$ ,  $y_2 = x_2$  and  $y_3 \propto x_1 x_2$ . A linear discriminant in this transformed space is a hyperplane, which cuts the surface. Points to the positive side of the hyperplane  $\hat{H}$  correspond to category  $\omega_1$ , and those beneath it correspond to category  $\omega_2$ . Here, in terms of the  $\mathbf{x}$  space,  $\mathcal{R}_1$  is a not simply connected. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



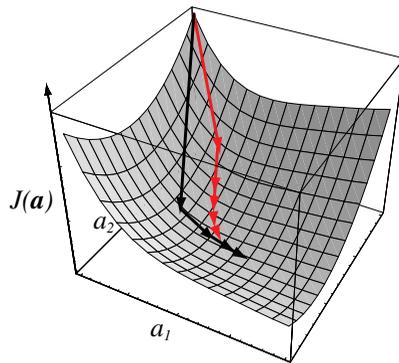
**FIGURE 5.7.** A three-dimensional augmented feature space  $\mathbf{y}$  and augmented weight vector  $\mathbf{a}$  (at the origin). The set of points for which  $\mathbf{a}^t \mathbf{y} = 0$  is a plane (or more generally, a hyperplane) perpendicular to  $\mathbf{a}$  and passing through the origin of  $\mathbf{y}$ -space, as indicated by the red disk. Such a plane need not pass through the origin of the two-dimensional feature space of the problem, as illustrated by the dashed decision boundary shown at the top of the box. Thus there exists an augmented weight vector  $\mathbf{a}$  that will lead to any straight decision line in  $\mathbf{x}$ -space. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



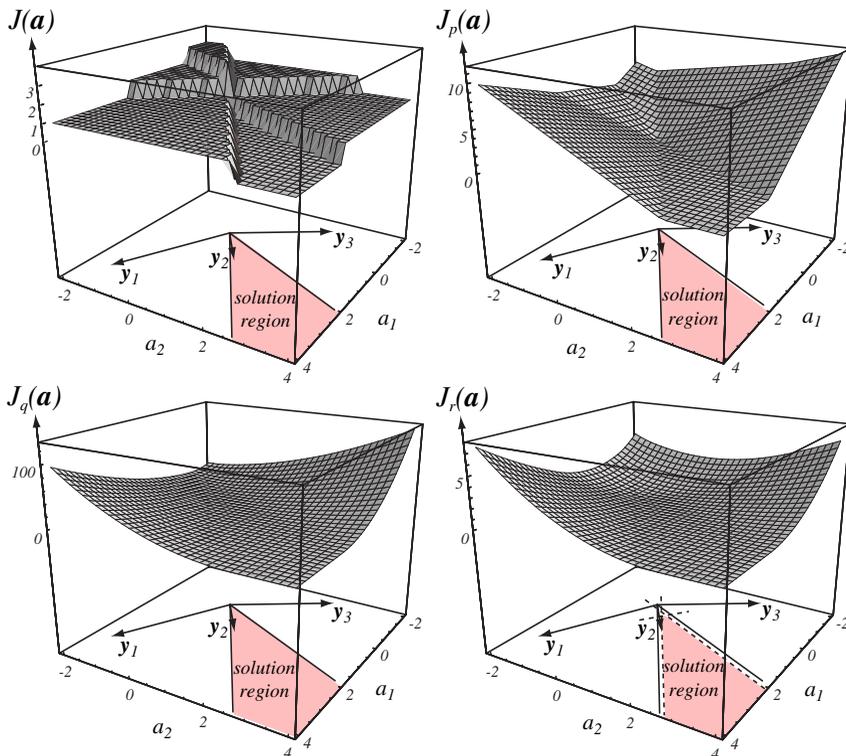
**FIGURE 5.8.** Four training samples (black for  $\omega_1$ , red for  $\omega_2$ ) and the solution region in feature space. The figure on the left shows the raw data; the solution vectors leads to a plane that separates the patterns from the two categories. In the figure on the right, the red points have been “normalized”—that is, changed in sign. Now the solution vector leads to a plane that places all “normalized” points on the same side. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



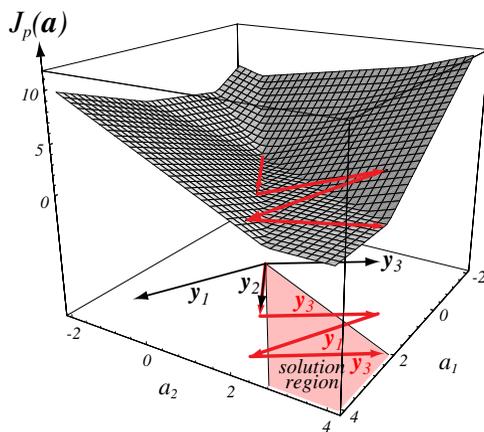
**FIGURE 5.9.** The effect of the margin on the solution region. At the left is the case of no margin ( $b = 0$ ) equivalent to a case such as shown at the left in Fig. 5.8. At the right is the case  $b > 0$ , shrinking the solution region by margins  $b/\|y_i\|$ . From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



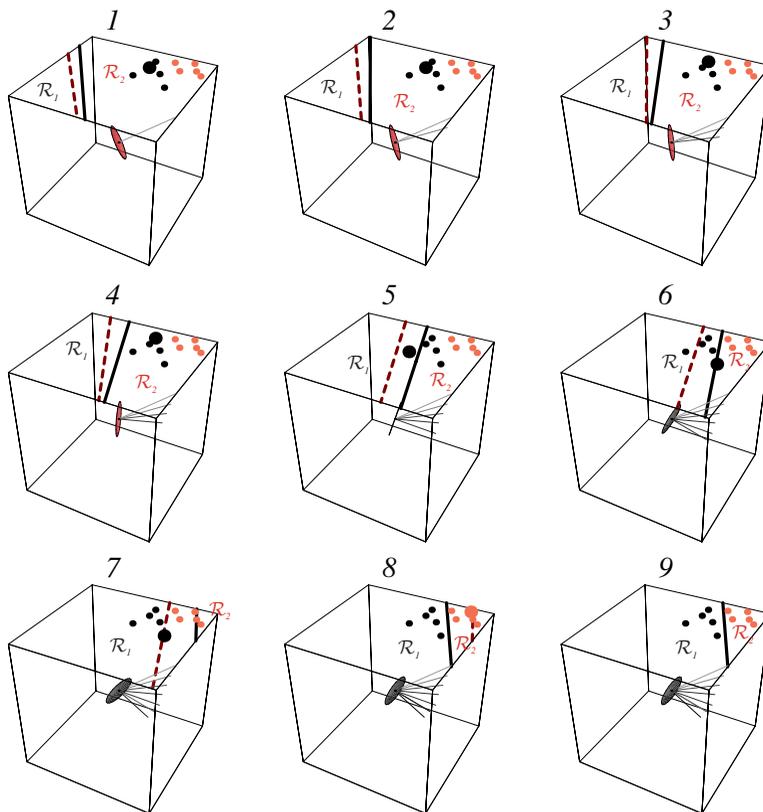
**FIGURE 5.10.** The sequence of weight vectors given by a simple gradient descent method (red) and by Newton's (second order) algorithm (black). Newton's method typically leads to greater improvement per step, even when using optimal learning rates for both methods. However the added computational burden of inverting the Hessian matrix used in Newton's method is not always justified, and simple gradient descent may suffice. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



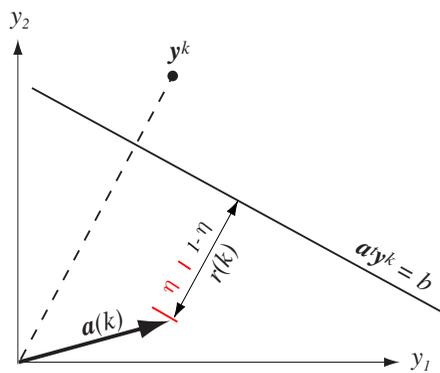
**FIGURE 5.11.** Four learning criteria as a function of weights in a linear classifier. At the upper left is the total number of patterns misclassified, which is piecewise constant and hence unacceptable for gradient descent procedures. At the upper right is the Perceptron criterion (Eq. 16), which is piecewise linear and acceptable for gradient descent. The lower left is squared error (Eq. 32), which has nice analytic properties and is useful even when the patterns are not linearly separable. The lower right is the square error with margin (Eq. 33). A designer may adjust the margin  $b$  in order to force the solution vector to lie toward the middle of the  $b = 0$  solution region in hopes of improving generalization of the resulting classifier. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



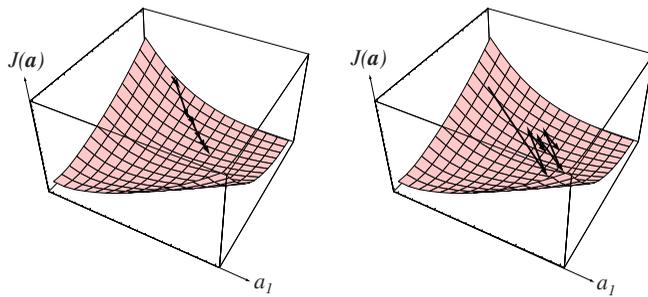
**FIGURE 5.12.** The Perceptron criterion,  $J_p(\mathbf{a})$ , is plotted as a function of the weights  $a_1$  and  $a_2$  for a three-pattern problem. The weight vector begins at  $\mathbf{0}$ , and the algorithm sequentially adds to it vectors equal to the “normalized” misclassified patterns themselves. In the example shown, this sequence is  $\mathbf{y}_2, \mathbf{y}_3, \mathbf{y}_1, \mathbf{y}_3$ , at which time the vector lies in the solution region and iteration terminates. Note that the second update (by  $\mathbf{y}_3$ ) takes the candidate vector *farther* from the solution region than after the first update (cf. Theorem 5.1). From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



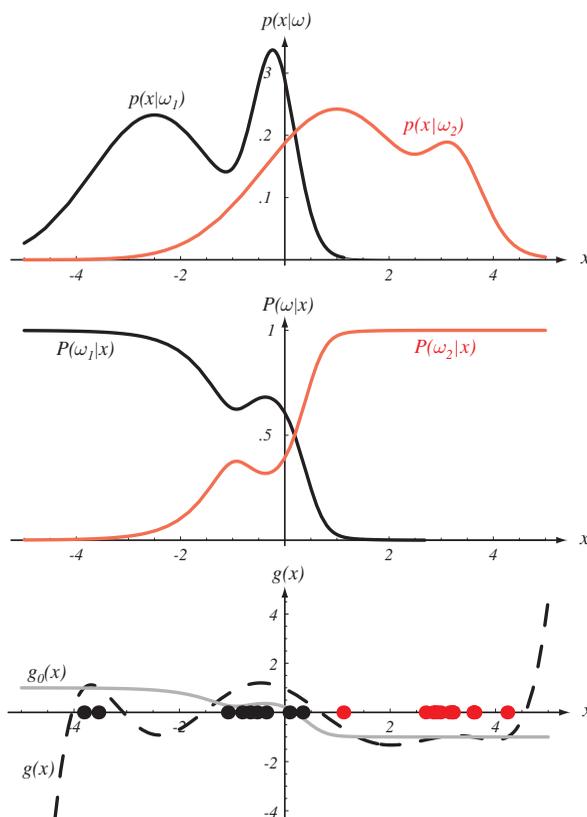
**FIGURE 5.13.** Samples from two categories,  $\omega_1$  (black) and  $\omega_2$  (red) are shown in augmented feature space, along with an augmented weight vector  $\mathbf{a}$ . At each step in a fixed-increment rule, one of the misclassified patterns,  $\mathbf{y}^k$ , is shown by the large dot. A correction  $\Delta\mathbf{a}$  (proportional to the pattern vector  $\mathbf{y}^k$ ) is added to the weight vector—toward an  $\omega_1$  point or away from an  $\omega_2$  point. This changes the decision boundary from the dashed position (from the previous update) to the solid position. The sequence of resulting  $\mathbf{a}$  vectors is shown, where later values are shown darker. In this example, by step 9 a solution vector has been found and the categories are successfully separated by the decision boundary shown. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



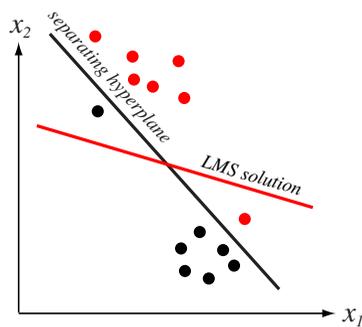
**FIGURE 5.14.** In each step of a basic relaxation algorithm, the weight vector is moved a proportion  $\eta$  of the way toward the hyperplane defined by  $\mathbf{a}'\mathbf{y}^k = b$ . From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



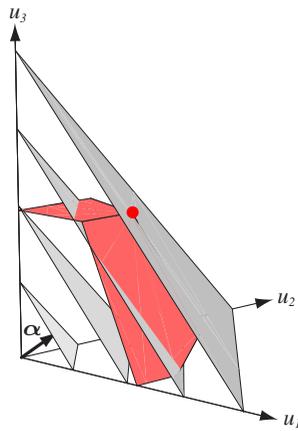
**FIGURE 5.15.** At the left, underrelaxation ( $\eta < 1$ ) leads to needlessly slow descent, or even failure to converge. Overrelaxation ( $1 < \eta < 2$ , shown at the right) describes overshooting; nevertheless, convergence will ultimately be achieved. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



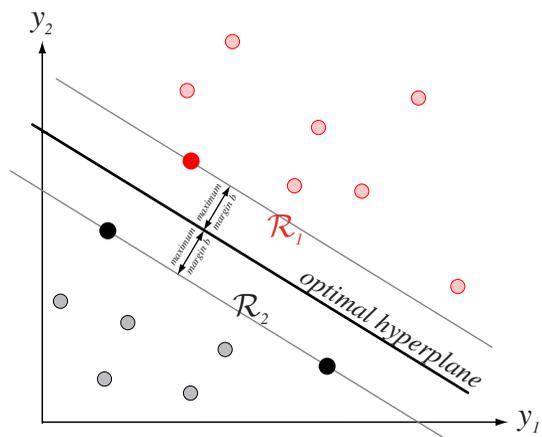
**FIGURE 5.16.** The top figure shows two class-conditional densities, and the middle figure the posteriors, assuming equal priors. Minimizing the MSE error also minimizes the mean-squared-error between  $\mathbf{a}'\mathbf{y}$  and the discriminant function  $g(\mathbf{x})$  (here a seventh-order polynomial) measured over the data distribution, as shown at the bottom. Note that the resulting  $g(x)$  best approximates  $g_0(x)$  in the regions where the data points lie. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



**FIGURE 5.17.** The LMS algorithm need not converge to a separating hyperplane, even if one exists. Because the LMS solution minimizes the sum of the squares of the distances of the training points to the hyperplane, for this example the plane is rotated clockwise compared to a separating hyperplane. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



**FIGURE 5.18.** Surfaces of constant  $z = \alpha^t \mathbf{u}$  are shown in gray, while constraints of the form  $\mathbf{A}\mathbf{u} \leq \mathbf{b}$  are shown in red. The simplex algorithm finds an extremum of  $z$  given the constraints, that is, where the gray plane intersects the red at a single point. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



**FIGURE 5.19.** Training a support vector machine consists of finding the optimal hyperplane, that is, the one with the maximum distance from the nearest training patterns. The support vectors are those (nearest) patterns, a distance  $b$  from the hyperplane. The three support vectors are shown as solid dots. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.