

Linear Models for Classification

4.1 Discriminant Function

4.1.1 Two classes

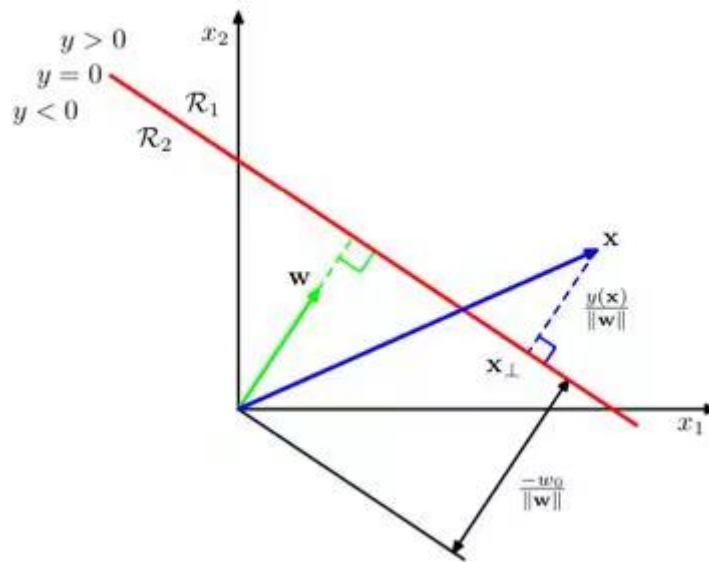
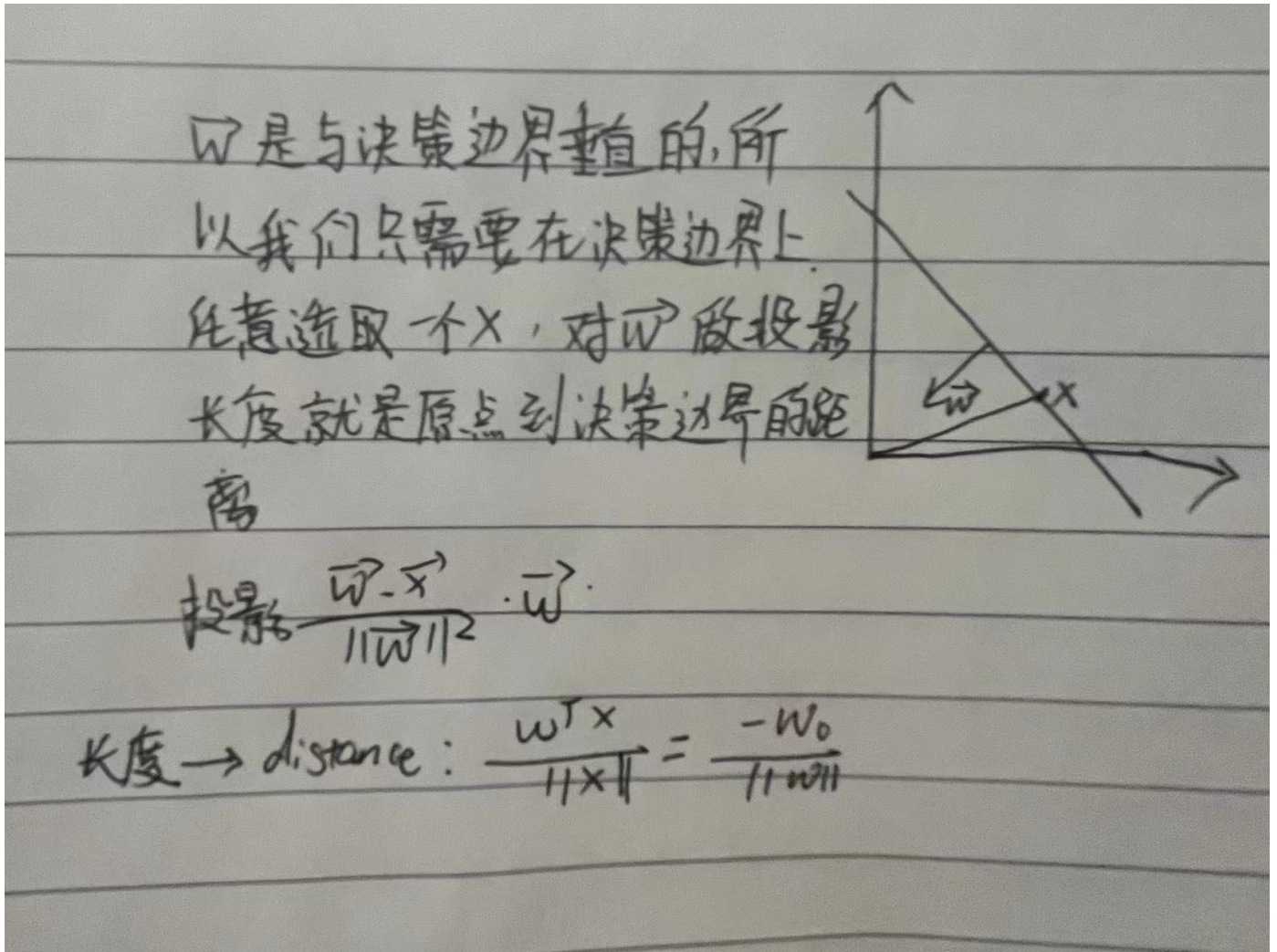
The simplest representation of a linear discriminant function is obtained by taking a linear function of the input vector so that

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

Consider two points \mathbf{x}_A and \mathbf{x}_B both of which lie on the decision surface. Because $y(\mathbf{x}_A) = y(\mathbf{x}_B) = 0$, we have $\mathbf{w}^T(\mathbf{x}_A - \mathbf{x}_B) = 0$ and hence the vector \mathbf{w} is orthogonal to every vector lying within the decision surface, and so \mathbf{w} determines the orientation of the decision surface.

$$\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|}$$

此公式适用于线性分类问题，及决策面可以用线性方程表示的情况



an arbitrary point x and let x_{\perp} be its orthogonal projection onto the decision surface, so that

$$x = x_{\perp} + r \frac{\vec{w}}{\|\vec{w}\|}$$

$$r = \frac{y(x)}{\|\vec{w}\|}$$

两边同时左乘 W^T , 再加上 w_0 .

利用 $y(x) = W^T x + w_0 = 0$

$y(x_L) = W^T x_L + w_0 = 0$

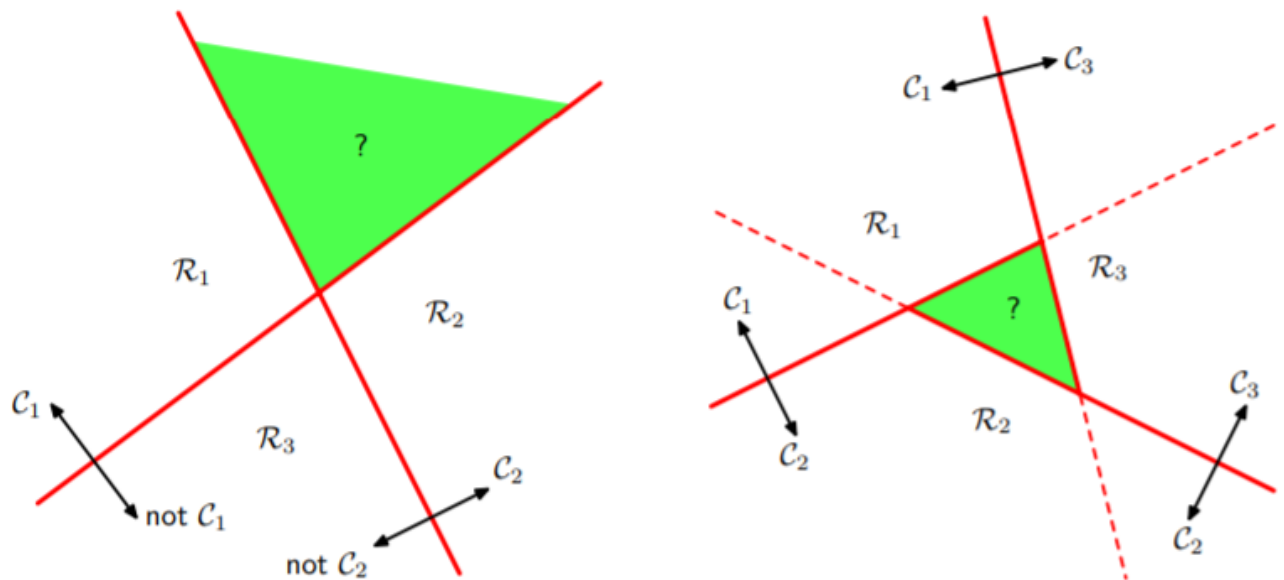
$$\underbrace{W^T x + w_0}_{y(x)} = \underbrace{W^T x_L + w_0}_{y(x_L)} + \gamma \cdot \frac{W^T W}{\|W\|} + w_0$$

$$\gamma = \frac{y(x)}{\|W\|}$$

4.1.2 Multiple classes

One-versus-the-rest Consider the use of $K-1$ classifiers each of which solves a two-class problem of separating points in a particular class C_k from points not in that class.

one-versus-one An alternative is to introduce $K(K-1)/2$ binary discriminant functions, one for every possible pair of classes.

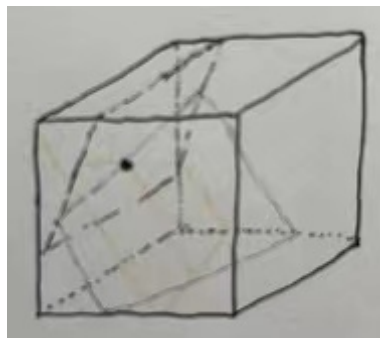


We can avoid these difficulties by considering a single K-class discriminant comprising K linear functions of the form

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

比如此时有三类数据，则判别函数为：

$$\begin{cases} y_1(x) = w_1^T x + w_{10} \\ y_2(x) = w_2^T x + w_{20} \\ y_3(x) = w_3^T x + w_{30} \end{cases}$$

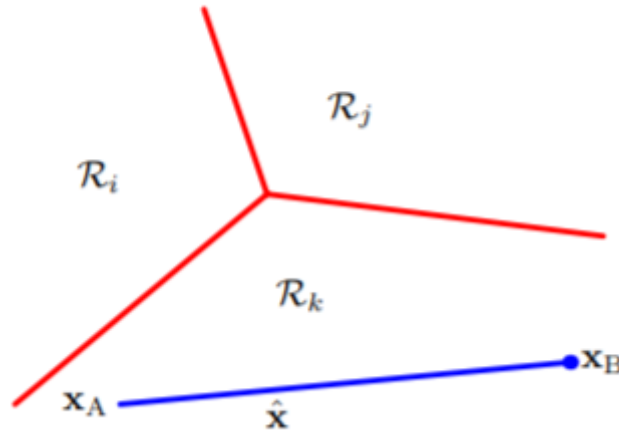


and then assigning a point \mathbf{x} to class C_k if $y_k(\mathbf{x}) > y_j(\mathbf{x})$ for all $j \neq k$. The decision boundary between class C_k and class C_j is therefore given by $y_k(\mathbf{x}) = y_j(\mathbf{x})$ and hence corresponds to a $(D - 1)$ -dimensional hyperplane

defined by

$$(\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0$$

The decision regions of such a discriminant are always singly connected and convex.



$$\hat{\mathbf{x}} = \lambda \mathbf{x}_A + (1 - \lambda) \mathbf{x}_B$$

where $0 \leq \lambda \leq 1$. From the linearity of the discriminant functions, it follows that

$$y_k(\hat{\mathbf{x}}) = \lambda y_k(\mathbf{x}_A) + (1 - \lambda) y_k(\mathbf{x}_B)$$

Because both \mathbf{x}_A and \mathbf{x}_B lie inside \mathcal{R}_k , it follows that $y_k(\mathbf{x}_A) > y_j(\mathbf{x}_A)$, and $y_k(\mathbf{x}_B) > y_j(\mathbf{x}_B)$, for all $j \neq k$, and hence $y_k(\hat{\mathbf{x}}) > y_j(\hat{\mathbf{x}})$, and so $\hat{\mathbf{x}}$ also lies inside \mathcal{R}_k . Thus \mathcal{R}_k is singly connected and convex.

4.1.3 Least squares for classification

Each class C_k is described by its own linear model so that

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

where $k = 1, \dots, K$. We can conveniently group these together using vector notation so that

$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}}$$

where $\widetilde{\mathbf{W}}$ is a matrix whose k^{th} column comprises the $D+1$ -dimensional vector $\widetilde{\mathbf{w}}_k = (w_{k0}, \mathbf{w}_k^T)^T$ and $\widetilde{\mathbf{x}}$ is the corresponding augmented input vector $(1, \mathbf{x}^T)^T$ with a dummy input $x_0 = 1$. This representation was discussed in detail in Section 3.1. A new input \mathbf{x} is then assigned to the class for which the output $y_k = \widetilde{\mathbf{w}}_k^T \widetilde{\mathbf{x}}$ is largest.

We now determine the parameter matrix $\widetilde{\mathbf{W}}$ by minimizing a sum-of-squares error function, as we did for regression in Chapter 3. Consider a training data set $\{\mathbf{x}_n, \mathbf{t}_n\}$ where $n = 1, \dots, N$, and define a matrix \mathbf{T} whose n^{th} row is the vector \mathbf{t}_n^T , together with a matrix $\widetilde{\mathbf{X}}$ whose n^{th} row is $\widetilde{\mathbf{x}}_n^T$. The sum-of-squares error function can then be written as

$$E_D(\widetilde{\mathbf{W}}) = \frac{1}{2} \text{Tr}\{(\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T})^T (\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T})\}$$

$$\tilde{W} = \begin{bmatrix} w_{10} & w_{20} & w_{k0} \\ \vdots & \vdots & \vdots \\ w_{1n} & w_{2n} & w_{kn} \end{bmatrix} \quad \text{每一列都是} \begin{pmatrix} w_{k0} \\ \vdots \\ w_{kn} \end{pmatrix} = \tilde{w}_k$$

$$\tilde{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \quad T = \begin{bmatrix} t_1^T \\ \vdots \\ t_n^T \end{bmatrix}$$

$$\tilde{X} \cdot \tilde{W} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \begin{bmatrix} \tilde{w}_1 & \dots & \tilde{w}_k \end{bmatrix}$$

$$= \begin{bmatrix} P(C_1|x_1) & P(C_k|x_1) \\ \vdots & \vdots \\ P(C_1|x_k) & P(C_k|x_k) \end{bmatrix} \quad x_i \text{ 分到 } C_i \text{ 中的估计 } y_i \text{ 预测值}$$

$$\tilde{X} \cdot \tilde{W} - T = \begin{bmatrix} P(C_1|x_1) - t_{11} & \dots & P(C_k|x_1) - t_{1k} \\ \vdots & \ddots & \vdots \\ P(C_1|x_k) - t_{k1} & \dots & P(C_k|x_k) - t_{kk} \end{bmatrix}$$

$$\text{若 } \tilde{X} \tilde{W} - T = a, \quad a_{ij} = P(C_j|x_i) - t_{ij}$$

$$(\tilde{X} \tilde{W} - T)^T (\tilde{X} \tilde{W} - T) = \begin{bmatrix} a_{11} & \dots & a_{1k} \\ \vdots & \ddots & \vdots \\ a_{k1} & \dots & a_{kk} \end{bmatrix} \begin{bmatrix} a_{11} & \dots & a_{1k} \\ \vdots & \ddots & \vdots \\ a_{k1} & \dots & a_{kk} \end{bmatrix}$$

$$= \begin{bmatrix} a_{11}^2 + a_{21}^2 + \dots + a_{k1}^2 & \vdots & \vdots \\ \vdots & \ddots & \vdots \\ a_{1k}^2 + a_{2k}^2 + \dots + a_{kk}^2 & \vdots & \ddots \end{bmatrix} \quad k \times k$$

$$\text{Tr}\{(\tilde{X} \tilde{W} - T)^T (\tilde{X} \tilde{W} - T)\} = \underbrace{(a_{11}^2 + a_{21}^2 + \dots + a_{k1}^2)}_{C_1 \text{ 分类误差平方和}} + \dots + \underbrace{(a_{1k}^2 + a_{2k}^2 + \dots + a_{kk}^2)}_{C_k \text{ 分类误差平方和}}$$

Setting the derivative with respect to \tilde{W} to zero, and rearranging, we then obtain the solution for \tilde{W} in the form

$$\tilde{W} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T T = \tilde{X}^\dagger T$$

相关推导

因此, 令 $\frac{\partial E_D(\hat{W})}{\partial \hat{W}} = 0$, 得:

$$\hat{W} = \left(\hat{X}^T \hat{X} \right)^{-1} \hat{X}^T T = \hat{X}^\dagger T$$

$\hat{X}^\dagger = \left(\hat{X}^T \hat{X} \right)^{-1} \hat{X}^T$ 为矩阵 \hat{X} 得伪逆矩阵。因此, 判别函数组为:

$$y(x) = \hat{W}^T \hat{x} = T^T \left(\hat{X}^\dagger \right)^T \hat{x}$$

Setting the derivative with respect to \widetilde{W} to zero, and rearranging, we then obtain the solution for \widetilde{W} in the form

$$\widetilde{W} = \left(\widetilde{X}^T \widetilde{X} \right)^{-1} \widetilde{X}^T T = \widetilde{X}^\dagger T$$

where \widetilde{X}^\dagger is the pseudo-inverse of the matrix \widetilde{X} , as discussed in Section 3.1.1. We then obtain the discriminant function in the form

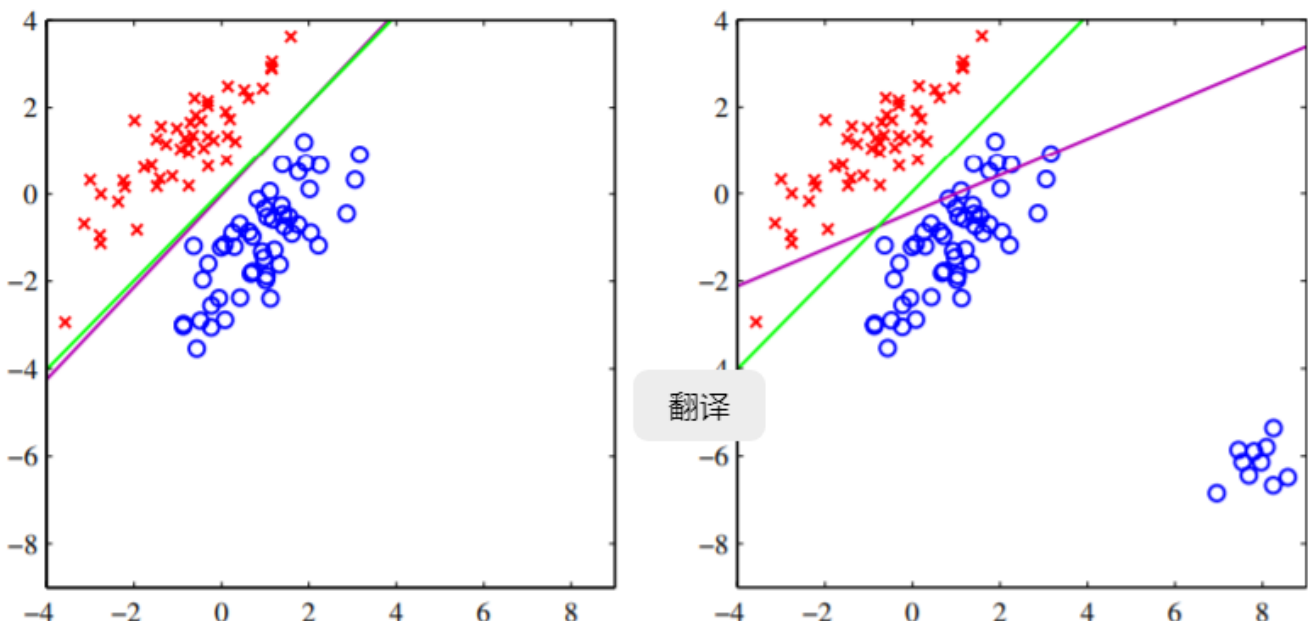
$$y(x) = \widetilde{W}^T \tilde{x} = T^T \left(\widetilde{X}^\dagger \right)^T \tilde{x}$$

An interesting property of least-squares solutions with multiple target variables is that if every target vector in the training set satisfies some linear constraint

$$\mathbf{a}^T \mathbf{t}_n + b = 0$$

for some constants a and b, then the model prediction for any value of x will satisfy the same constraint so that

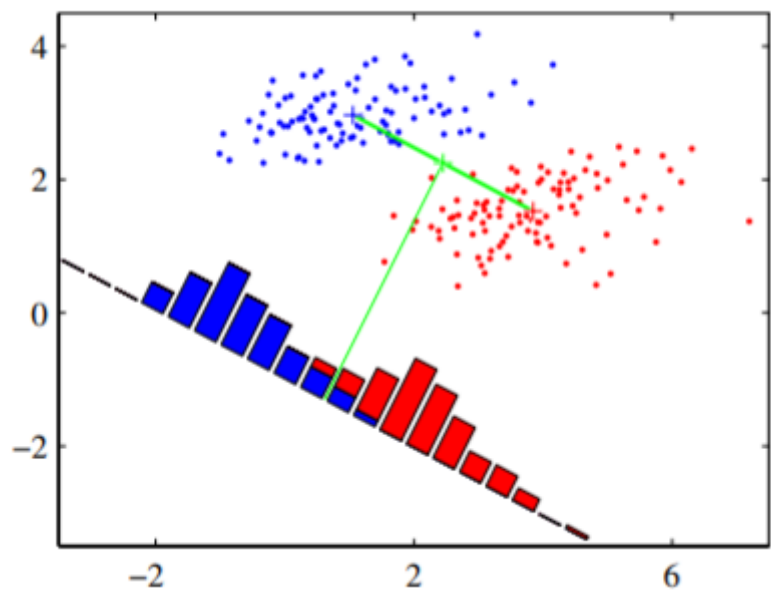
$$\mathbf{a}^T \mathbf{y}(x) + b = 0$$



4.1.4--4.1.6 Fisher's linear discriminant

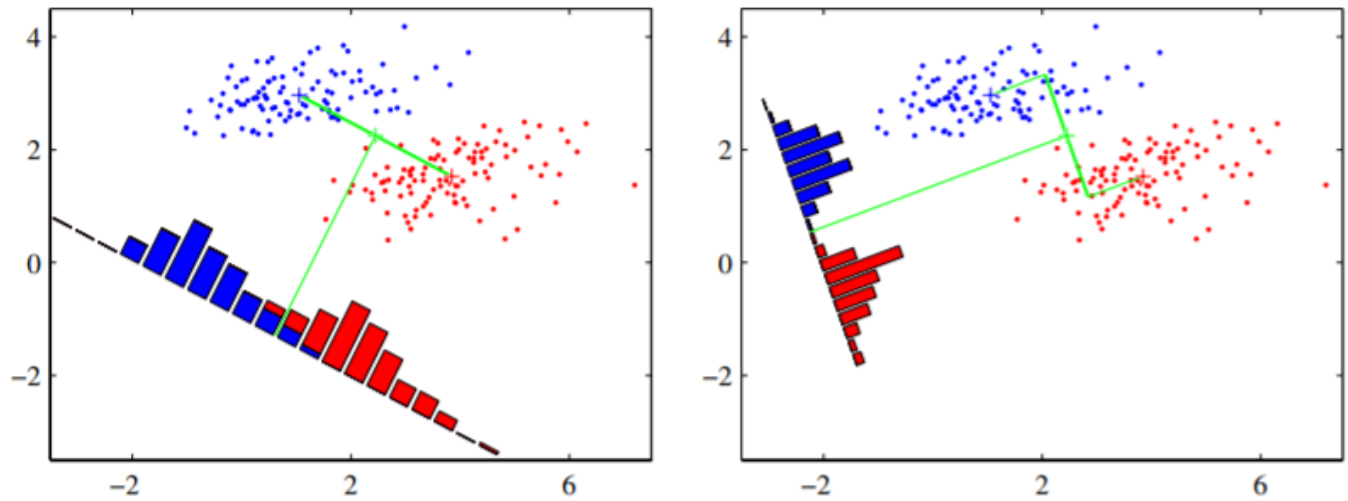
One way to view a linear classification model is in terms of dimensionality reduction. Consider first the case of two classes, and suppose we take the D-dimensional input vector \mathbf{x} and project it down to one dimension using

$$y = \mathbf{w}^T \mathbf{x}$$



If we place a threshold on y and classify $y \geq -w_0$ as class C_1 , and otherwise class C_2 , then we obtain our standard linear classifier discussed in the previous section.

D-dimensional	1-dimensional
mean vectors	
$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in C_2} \mathbf{x}_n$	$m_k = \mathbf{w}^T \mathbf{m}_k$
within-class discrete degree	
$S_K = \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T$ $S_W = S_1 + S_2$	$s_k^2 = \sum_{n \in C_k} (y_n - m_k)^2$ $s_1^2 + s_2^2$
between-class discrete degree	
$S_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$	$(m_2 - m_1)^2$



where $y_n = \mathbf{w}^T \mathbf{x}_n$. We can define the total within-class variance for the whole data set to be simply $s_1^2 + s_2^2$. The Fisher criterion is defined to be the ratio of the between-class variance to the within-class variance and is given by

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

We can make the dependence on \mathbf{w} explicit by using (4.20), (4.23), and (4.24) to rewrite the Fisher criterion in the form

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

$$\begin{aligned}
 (\mathbf{m}_2 - \mathbf{m}_1)^2 &= (\mathbf{W}^T \mathbf{m}_2 - \mathbf{W}^T \mathbf{m}_1)^2 \\
 &= (\mathbf{W}^T \mathbf{m}_2 - \mathbf{W}^T \mathbf{m}_1) (\mathbf{W}^T \mathbf{m}_2 - \mathbf{W}^T \mathbf{m}_1)^T \\
 &= (\mathbf{W}^T \mathbf{m}_2 - \mathbf{W}^T \mathbf{m}_1) [\mathbf{W}^T (\mathbf{m}_2 - \mathbf{m}_1)]^T \\
 &= (\mathbf{W}^T \mathbf{m}_2 - \mathbf{W}^T \mathbf{m}_1) [(\mathbf{W}^T)^T (\mathbf{m}_2 - \mathbf{m}_1)^T] \\
 &= (\mathbf{W}^T \mathbf{m}_2 - \mathbf{W}^T \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{W} \\
 &= \mathbf{W}^T \mathbf{S}_B \mathbf{W}
 \end{aligned}$$

$(AB)^T = A^T \cdot B^T$

$$\begin{aligned}
 S_H^2 &= \sum_{n \in C_1} (y_n - m_k)^2 = \sum_{n \in C_1} (\mathbf{W}^T \mathbf{x}_n - \mathbf{W}^T \mathbf{m}_k)^2 \\
 &= \mathbf{W}^T \left[\sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_k) (\mathbf{x}_n - \mathbf{m}_k)^T \right] \mathbf{W} \\
 &= \mathbf{W}^T \mathbf{S}_W \mathbf{W}
 \end{aligned}$$

where \mathbf{S}_B is the between-class covariance matrix and is given by

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1)^T$$

and \mathbf{S}_W is the total within-class covariance matrix, given by

$$\mathbf{S}_W = \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1) (\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2) (\mathbf{x}_n - \mathbf{m}_2)^T$$

Differentiating (4.26) with respect to \mathbf{w} , we find that $J(\mathbf{w})$ is maximized when

$$(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}.$$

From (4.27), we see that $\mathbf{S}_B \mathbf{w}$ is always in the direction of $(\mathbf{m}_2 - \mathbf{m}_1)$. Furthermore, we do not care about the magnitude of \mathbf{w} , only its direction, and so we can drop the scalar factors $(\mathbf{w}^T \mathbf{S}_B \mathbf{w})$ and $(\mathbf{w}^T \mathbf{S}_W \mathbf{w})$.

Multiplying both sides of (4.29) by \mathbf{S}_W^{-1} we then obtain

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$$

Note that if the within-class covariance is isotropic, so that \mathbf{S}_W is proportional to the unit matrix, we find that \mathbf{w} is proportional to the difference of the class means, as discussed above.

如果类内协方差矩阵是各向同性（各向同性即指随机向量的协方差矩阵为标量乘以单位矩阵，即每个方向的方差相同，对角阵是因为相关性在考察两个不同类之间不重要），则 $\mathbf{S}_W \propto \mathbf{I}$ ，即类内协方差矩阵正比于单位矩

阵，因此其逆矩阵也正比于单位矩阵，则： $w \propto m_2 - m_1$ 此时， w 正比于类内均值之差， w 与类内均值之差的方向相同。

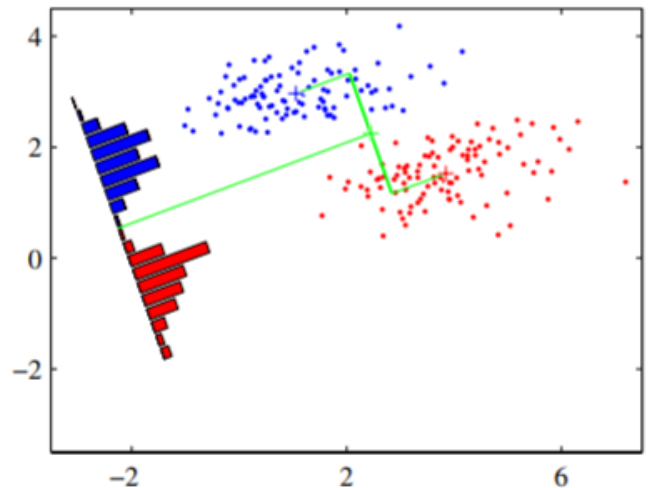
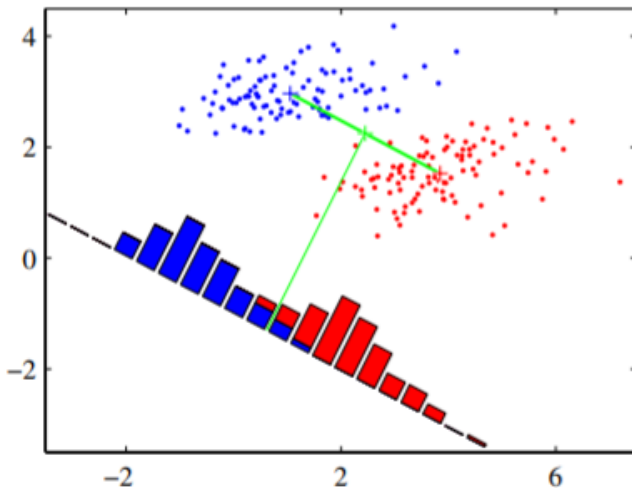
设 $W^T S_W W = a$, $W^T S_B W = b$, 两项均为标量
上式两边左乘 $S^{-1}W$, 得

$$W = \frac{a}{b} S_W^{-1} S_B W = \frac{a}{b} S_W^{-1} (m_2 - m_1)(m_2 - m_1)^T W$$

设 $C = (m_2 - m_1)^T W$ 标量

$$W = \frac{aC}{b} S_W^{-1} (m_2 - m_1)$$

因此 $W \propto S^{-1} W (m_2 - m_1)$



Fisher's discriminant for multiple classes

We now consider the generalization of the Fisher discriminant to $K > 2$ classes, and we shall assume that the dimensionality D of the input space is greater than the number K of classes. Next, we introduce $D' > 1$ linear 'features' $y_k = \mathbf{w}_k^T \mathbf{x}$, where $k=1, \dots, D'$. These feature values can conveniently be grouped together to form a vector \mathbf{y} . Similarly, the weight vectors $\{\mathbf{w}_k\}$ can be considered to be the columns of a matrix \mathbf{W} , so that

$$\mathbf{y} = \mathbf{W}^T \mathbf{x}$$

Note that again we are not including any bias parameters in the definition of \mathbf{y} . The generalization of the within-class covariance matrix to the case of K classes follows from (4.28) to give

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{S}_k$$

where

$$\mathbf{S}_k = \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m}_k) (\mathbf{x}_n - \mathbf{m}_k)^T$$

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}_n$$

and N_k is the number of patterns in class C_k . In order to find a generalization of the between-class covariance matrix, we follow Duda and Hart (1973) and consider first the total covariance matrix

$$\mathbf{S}_T = \sum_{n=1}^N (\mathbf{x}_n - \mathbf{m}) (\mathbf{x}_n - \mathbf{m})^T$$

where \mathbf{m} is the mean of the total data set

$$\mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} \sum_{k=1}^K N_k \mathbf{m}_k$$

and $N = \sum_k N_k$ is the total number of data points. The total covariance matrix can be decomposed into the sum of the within-class covariance matrix, given by (4.40) and (4.41), plus an additional matrix \mathbf{S}_B , which we identify as a measure of the between-class covariance

$$\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B$$

where

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m}) (\mathbf{m}_k - \mathbf{m})^T$$

These covariance matrices have been defined in the original \mathbf{x} -space. We can now define similar matrices in the projected D' - *dimensional* - space

$$\mathbf{s}_W = \sum_{k=1}^K \sum_{n \in C_k} (\mathbf{y}_n - \boldsymbol{\mu}_k) (\mathbf{y}_n - \boldsymbol{\mu}_k)^T$$

and

$$\mathbf{s}_B = \sum_{k=1}^K N_k (\boldsymbol{\mu}_k - \boldsymbol{\mu}) (\boldsymbol{\mu}_k - \boldsymbol{\mu})^T$$

where

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{y}_n, \quad \boldsymbol{\mu} = \frac{1}{N} \sum_{k=1}^K N_k \boldsymbol{\mu}_k$$

Again we wish to construct a scalar that is large when the between-class covariance is large and when the within-class covariance is small. There are now many possible choices of criterion (Fukunaga, 1990). One example is given by

$$J(\mathbf{W}) = \text{Tr}\{\mathbf{S}_W^{-1}\mathbf{S}_B\}$$

This criterion can then be rewritten as an explicit function of the projection matrix \mathbf{W} in the form

$$J(\mathbf{w}) = \text{Tr}\left\{(\mathbf{W}\mathbf{S}_W\mathbf{W}^T)^{-1}(\mathbf{W}\mathbf{S}_B\mathbf{W}^T)\right\}$$

4.1.7 The perceptron algorithm

Another example of a linear discriminant model is the perceptron of Rosenblatt (1962), which occupies an important place in the history of pattern recognition algorithms. It corresponds to a two-class model in which the input vector \mathbf{x} is first transformed using a fixed nonlinear transformation to give a feature vector $\phi(\mathbf{x})$, and this is then used to construct a generalized linear model of the form

$$y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$$

where the nonlinear activation function $f(\cdot)$ is given by a step function of the form

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0. \end{cases} \quad (4.53)$$

We therefore consider an alternative error function known as the perceptron criterion. To derive this, we note that we are seeking a weight vector \mathbf{w} such that patterns \mathbf{x}_n in class C_1 will have $\mathbf{w}^T \phi(\mathbf{x}_n) > 0$, whereas patterns \mathbf{x}_n in class C_2 have $\mathbf{w}^T \phi(\mathbf{x}_n) < 0$. Using the $\{-1, +1\}$ target coding scheme it follows that we would like all patterns to satisfy $\mathbf{w}^T \phi(\mathbf{x}_n) t_n > 0$. The perceptron criterion associates zero error with any pattern that is correctly classified, whereas for a misclassified pattern \mathbf{x}_n it tries to minimize the quantity $-\mathbf{w}^T \phi(\mathbf{x}_n) t_n$. The perceptron criterion is therefore given by

$$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi_n t_n$$

We now apply the stochastic gradient descent algorithm to this error function. The change in the weight vector \mathbf{w} is then given by

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_P(\mathbf{w}) = \mathbf{w}^{(\tau)} + \eta \phi_n t_n$$

If we consider the effect of a single update in the perceptron learning algorithm, we see that the contribution to the error from a misclassified pattern will be reduced because from (4.55) we have

$$-\mathbf{w}^{(\tau+1)T} \phi_n t_n = -\mathbf{w}^{(\tau)T} \phi_n t_n - (\phi_n t_n)^T \phi_n t_n < -\mathbf{w}^{(\tau)T} \phi_n t_n$$

