

WITHIN-CLASS SCATTER

BETWEEN-CLASS SCATTER

We call S_W the *within-class scatter matrix*. It is proportional to the sample covariance matrix for the pooled d -dimensional data. It is symmetric and positive semidefinite, and it is usually nonsingular if $n > d$. Likewise, S_B is called the *between-class scatter matrix*. It is also symmetric and positive semidefinite, but because it is the outer product of two vectors, its rank is at most one. In particular, for any w , $S_B w$ is in the direction of $m_1 - m_2$, and S_B is quite singular.

In terms of S_B and S_W , the criterion function $J(\cdot)$ can be written as

$$J(w) = \frac{w^T S_B w}{w^T S_W w}. \tag{103}$$

This expression is well known in mathematical physics as the generalized Rayleigh quotient. It is easy to show that a vector w that maximizes $J(\cdot)$ must satisfy

$$S_B w = \lambda S_W w, \tag{104}$$

for some constant λ , which is a generalized eigenvalue problem (Problem 42). This can also be seen informally by noting that at an extremum of $J(w)$ a small change in w in Eq. 103 should leave unchanged the ratio of the numerator to the denominator. If S_W is nonsingular we can obtain a conventional eigenvalue problem by writing

$$S_W^{-1} S_B w = \lambda w. \tag{105}$$

In our particular case, it is unnecessary to solve for the eigenvalues and eigenvectors of $S_W^{-1} S_B$ due to the fact that $S_B w$ is always in the direction of $m_1 - m_2$. Because the scale factor for w is immaterial, we can immediately write the solution for the w that optimizes $J(\cdot)$:

$$w = S_W^{-1} (m_1 - m_2). \tag{106}$$

Thus, we have obtained w for Fisher's linear discriminant—the linear function yielding the maximum ratio of between-class scatter to within-class scatter. (The solution w given by Eq. 106 is sometimes called the *canonical variate*.) Thus the classification has been converted from a d -dimensional problem to a hopefully more manageable one-dimensional one. This mapping is many-to-one, and in theory it cannot possibly reduce the minimum achievable error rate if we have a very large training set. In general, one is willing to sacrifice some of the theoretically attainable performance for the advantages of working in one dimension. All that remains is to find the threshold, that is, the point along the one-dimensional subspace separating the projected points.

When the conditional densities $p(x|\omega_i)$ are multivariate normal with equal covariance matrices Σ , we can calculate the threshold directly. In that case we recall from Chapter 2 that the optimal decision boundary has the equation

$$w^T x + w_0 = 0 \tag{107}$$

where

$$w = \Sigma^{-1} (\mu_1 - \mu_2), \tag{108}$$

and where w_0 is a constant involving w and the prior probabilities. If we use sample means and the sample covariance matrix to estimate μ_i and Σ , we obtain a vector in

3.8.3 Multiple Discrimination

TOTAL MEAN VECTOR
TOTAL SCATTER MATRIX

the s
equa
linea
erall
shou
T
crim
ter a

For
invo
spac
gene

wher

and

T
a to

and

The