

Pattern Classification

All materials in these slides were taken from Pattern Classification (2nd ed) by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000 with the permission of the authors and the publisher Chapter 2 (part 3) Bayesian Decision Theory (Sections 2-6,2-9)

Discriminant Functions for the Normal Density

Bayes Decision Theory – Discrete Features

Discriminant Functions for the Normal Density

 We saw that the minimum error-rate classification can be achieved by the discriminant function

 $g_i(x) = \ln P(x \mid \omega_i) + \ln P(\omega_i)$

Case of multivariate normal

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^t \sum_{i=1}^{n-1} (x - \mu_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

Pattern Classification, Chapter 2 (Part 3)

Case $\Sigma_i = \sigma^2 I$ (I stands for the identity matrix)

• What does " $\Sigma_i = \sigma^2 I$ " say about the dimensions?

• What about the variance of each dimension?

Note : both $|\Sigma_i|$ and (d/2) $\ln \pi$ are independent of *i* in

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^t \sum_{i=1}^{n-1} (x - \mu_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

Thus we can simplify to :

$$g_i(x) = -\frac{\left\|\chi - \mu_i\right\|^2}{2\sigma^2} + \ln P(\omega_i)$$

where $\|\cdot\|$ denotes the Euclidean norm

We can further simplify by recognizing that the quadratic term x^tx implicit in the Euclidean norm is the same for all *i*.

 $g_i(x) = \mathbf{w}_i^t \mathbf{x} + w_{i0}$ (linear discriminant function) where:

$$\mathbf{w}_i = \frac{\mathbf{\mu}_i}{\sigma^2}; \ w_{i0} = -\frac{1}{2\sigma^2} \mathbf{\mu}_i^t \mathbf{\mu}_i + \ln P(\omega_i)$$

(ω_{i0} is called the threshold for the *i*th category!)

 A classifier that uses linear discriminant functions is called "a linear machine"

The decision surfaces for a linear machine are pieces of hyperplanes defined by:

 $g_i(x) = g_j(x)$

The equation can be written as: $w^{t}(x-x_{0})=0$

• The hyperplane separating \mathcal{R}_i and \mathcal{R}_i

$$\mathbf{x}_{0} = \frac{1}{2} (\boldsymbol{\mu}_{i} + \boldsymbol{\mu}_{j}) - \frac{\sigma^{2}}{\left\|\boldsymbol{\mu}_{i} - \boldsymbol{\mu}_{j}\right\|^{2}} \ln \frac{P(\omega_{i})}{P(\omega_{j})} (\boldsymbol{\mu}_{i} - \boldsymbol{\mu}_{j})$$

always orthogonal to the line linking the means!

if
$$P(\omega_i) = P(\omega_j)$$
 then $\mathbf{x}_0 = \frac{1}{2}(\boldsymbol{\mu}_i + \boldsymbol{\mu}_j)$



FIGURE 2.10. If the covariance matrices for two distributions are equal and proportional to the identity matrix, then the distributions are spherical in *d* dimensions, and the boundary is a generalized hyperplane of d - 1 dimensions, perpendicular to the line separating the means. In these one-, two-, and three-dimensional examples, we indicate $p(\mathbf{x}|\omega_i)$ and the boundaries for the case $P(\omega_1) = P(\omega_2)$. In the three-dimensional case, the grid plane separates \mathcal{R}_1 from \mathcal{R}_2 . From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

7









Pattern Classification, Chapter 2 (Part 3)



FIGURE 2.11. As the priors are changed, the decision boundary shifts; for sufficiently disparate priors the boundary will not lie between the means of these one-, two- and three-dimensional spherical Gaussian distributions. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

• Case $\Sigma_i = \Sigma$ (covariance of all classes are identical but arbitrary!) Hyperplane separating R_i and R_i Has the equation $|\mathbf{w}^{t}(\mathbf{x}-\mathbf{x}_{0})|=0$ Where $\mathbf{w} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_i)$ and $\mathbf{x}_{0} = \frac{1}{2} (\boldsymbol{\mu}_{i} + \boldsymbol{\mu}_{j}) - \frac{\ln \left[P(\boldsymbol{\omega}_{i}) / P(\boldsymbol{\omega}_{j}) \right]}{(\boldsymbol{\mu}_{i} - \boldsymbol{\mu}_{j})^{t} \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_{i} - \boldsymbol{\mu}_{j})} . (\boldsymbol{\mu}_{i} - \boldsymbol{\mu}_{j})$ (the hyperplane separating R_i and R_i is generally not orthogonal to the line between the means!)

Pattern Classification, Chapter 2 (Part 3)



Pattern Classification, Chapter 2 (Part 3)



FIGURE 2.12. Probability densities (indicated by the surfaces in two dimensions and ellipsoidal surfaces in three dimensions) and decision regions for equal but asymmetric Gaussian distributions. The decision hyperplanes need not be perpendicular to the line connecting the means. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

• Case Σ_i = arbitrary

The covariance matrices are different for each category

 $g_{i}(x) = x^{t}W_{i}x + w_{i}^{t}x = w_{i0}$ where : $W_{i} = -\frac{1}{2}\Sigma_{i}^{-1}$ $w_{i} = \Sigma_{i}^{-1}\mu_{i}$ $w_{i0} = -\frac{1}{2}\mu_{i}^{t}\Sigma_{i}^{-1}\mu_{i} - \frac{1}{2}\ln|\Sigma_{i}| + \ln P(\omega_{i})$

The decision surfaces are hyperquadratics (Hyperquadrics are: hyperplanes, pairs of hyperplanes, hyperspheres, hyperellipsoids, hyperparaboloids, hyperhyperboloids)



14



FIGURE 2.14. Arbitrary Gaussian distributions lead to Bayes decision boundaries that are general hyperquadrics. Conversely, given any hyperquadric, one can find two Gaussian distributions whose Bayes decision boundary is that hyperquadric. These variances are indicated by the contours of constant probability density. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

 Components of x are binary or integer valued, x can take only one of m discrete values

concerned with probabilities rather than probability densities in Bayes Formula:

$$P(\omega_{j} | \mathbf{x}) = \frac{P(\mathbf{x} | \omega_{j})P(\omega_{j})}{P(\mathbf{x})}$$

where

$$P(\mathbf{x}) = \sum_{j=1}^{c} P(\mathbf{x} \mid \boldsymbol{\omega}_{j}) P(\boldsymbol{\omega}_{j})$$

• Conditional risk is defined as before: $R(\alpha | \mathbf{x})$

• Approach is still to minimize risk:

$$\alpha^* = \arg\min_i R(\alpha_i \mid \mathbf{x})$$

Case of independent binary features in 2 category problem
Let x = [x₁, x₂, ..., x_d]^t where each x_i is either 0 or 1, with probabilities:

 $p_i = P(x_i = 1 \mid \omega_1)$ $q_i = P(x_i = 1 \mid \omega_2)$

 Assuming conditional independence, P(x|ω_i) can be written as a product of component probabilities:

$$P(\mathbf{x} \mid \omega_1) = \prod_{i=1}^{d} p_i^{x_i} (1 - p_i)^{1 - x_i}$$

and

$$P(\mathbf{x} \mid \omega_2) = \prod_{i=1}^{d} q_i^{x_i} (1 - q_i)^{1 - x_i}$$

yielding a likelihood ratio given by :

$$\frac{P(\mathbf{x} \mid \omega_1)}{P(\mathbf{x} \mid \omega_2)} = \prod_{i=1}^d \left(\frac{p_i}{q_i}\right)^{x_i} \left(\frac{1-p_i}{1-q_i}\right)^{1-x_i}$$

Taking our likelihood ratio

$$\frac{P(\mathbf{x} \mid \omega_1)}{P(\mathbf{x} \mid \omega_2)} = \prod_{i=1}^d \left(\frac{p_i}{q_i}\right)^{x_i} \left(\frac{1-p_i}{1-q_i}\right)^{1-x_i}$$

and plugging it into Eq. 31
$$g(\mathbf{x}) = \ln \frac{p(\mathbf{x} \mid \omega_1)}{p(\mathbf{x} \mid \omega_2)} + \ln \frac{p(\omega_1)}{p(\omega_2)}$$

yields:
$$g(\mathbf{x}) = \sum_{i=1}^d \left[x_i \ln \frac{p_i}{q_i} + (1-x_i) \ln \frac{1-p_i}{1-q_i}\right] + \ln \frac{p(\omega_1)}{p(\omega_2)}$$

The discriminant function in this case is:

$$g(x) = \sum_{i=1}^d w_i x_i + w_0$$

where :

$$w_i = ln \frac{p_i(1-q_i)}{q_i(1-p_i)}$$
 $i = 1,...,d$

and :

$$w_{0} = \sum_{i=1}^{d} ln \frac{1-p_{i}}{1-q_{i}} + ln \frac{P(\omega_{1})}{P(\omega_{2})}$$

decide ω_{1} if $g(x) > 0$ and ω_{2} if $g(x) \le 0$