



Bayesian Decision Theory

Prof. Richard Zanibbi

Bayesian Decision Theory

The Basic Idea

To minimize errors, choose the least risky class, i.e. the class for which the *expected loss* is smallest

Assumptions

Problem posed in probabilistic terms, and all relevant probabilities are known

Probability Mass vs. Probability Density Functions

Probability Mass Function, $P(x)$

Probability for values of discrete random variable x . Each value has its own associated probability

$$\chi = \{v_1, \dots, v_m\}$$

$$P(x) \geq 0, \text{ and } \sum_{x \in \chi} P(x) = 1$$

$$Pr[x \in (a, b)] = \int_a^b p(x) dx$$

Probability Density, $p(x)$

$$p(x) \geq 0 \text{ and } \int_{-\infty}^{\infty} p(x) dx = 1$$

Probability for values of continuous random variable x . Probability returned is for an *interval* within which the value lies (intervals defined by some unit distance)

Prior Probability

Definition ($P(w)$)

The likelihood of a value for a random variable representing the *state of nature* (*true class for the current input*), in the absence of other information

- Informally, “what percentage of the time state X occurs”

Example

The prior probability that an instance taken from two classes is provided as input, in the absence of any features (e.g. $P(\text{cat}) = 0.3$, $P(\text{dog}) = 0.7$)

Class-Conditional Probability Density Function (for Continuous Features)

Definition ($p(x | w)$)

The probability of a value for continuous random variable x , given a state of nature in w

- For each value of x , we have a different class-conditional pdf for each class in w (example next slide)

Example: Class-Conditional Probability Densities

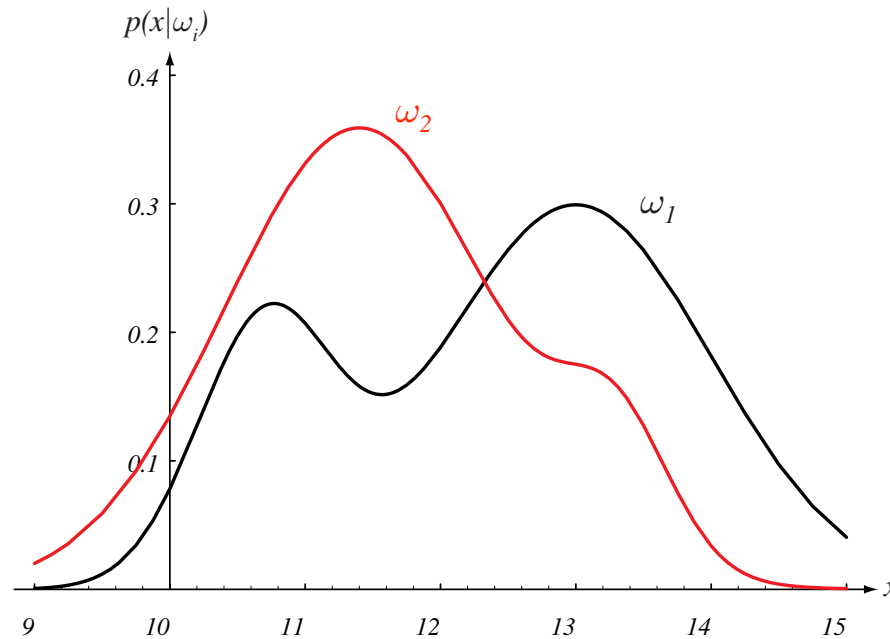


FIGURE 2.1. Hypothetical class-conditional probability density functions show the probability density of measuring a particular feature value x given the pattern is in category ω_i . If x represents the lightness of a fish, the two curves might describe the difference in lightness of populations of two types of fish. Density functions are normalized, and thus the area under each curve is 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Bayes Formula

$$P(\omega_j|x) = \frac{p(x|\omega_j)P(\omega_j)}{p(x)} \quad \text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

where $p(x) = \sum_{j=1}^c p(x|\omega_j)P(\omega_j)$

Purpose

Convert class prior and class-conditional densities to a *posterior probability* for a class: the probability of a class given the input features ('post-observation')

Example: Posterior Probabilities

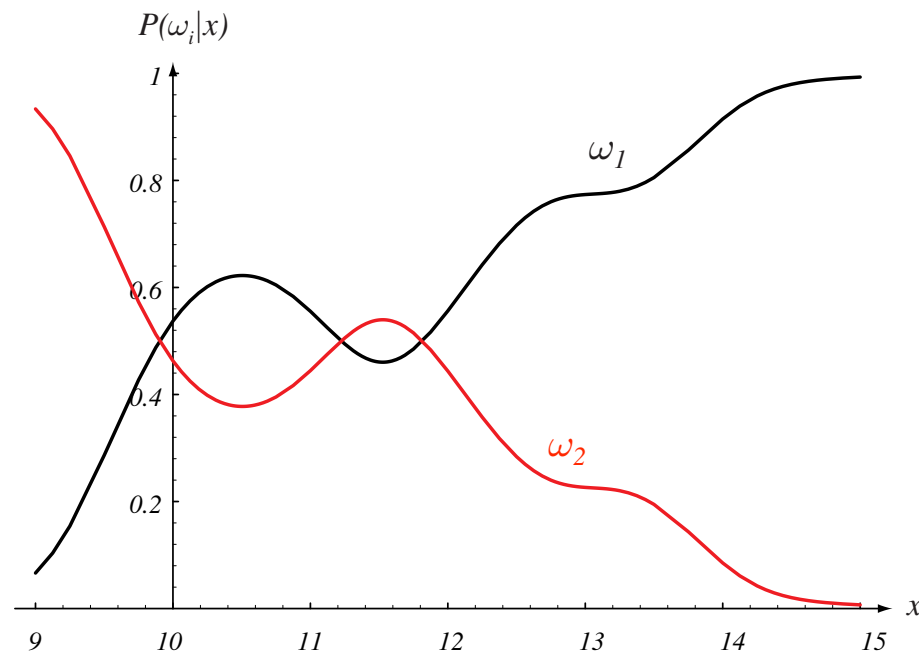


FIGURE 2.2. Posterior probabilities for the particular priors $P(\omega_1) = 2/3$ and $P(\omega_2) = 1/3$ for the class-conditional probability densities shown in Fig. 2.1. Thus in this case, given that a pattern is measured to have feature value $x = 14$, the probability it is in category ω_2 is roughly 0.08, and that it is in ω_1 is 0.92. At every x , the posteriors sum to 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Choosing the Most Likely Class

What happens if we do the following?

Decide ω_1 if $P(\omega_1|x) > P(\omega_2|x)$; otherwise decide ω_2

A. We minimize the average probability of error. Consider the two-class case from previous slide:

$$P(\text{error}|x) = \begin{cases} P(\omega_1|x) & \text{if we choose } \omega_2 \\ P(\omega_2|x) & \text{if we choose } \omega_1 \end{cases}$$

$$P(\text{error}) = \int_{-\infty}^{\infty} P(\text{error}|x)p(x) dx \quad (\text{average error})$$

Expected Loss or *Conditional Risk* of an Action

$$R(\alpha_i|\mathbf{x}) = \sum_{j=1}^c \lambda(\alpha_i|\omega_j)P(\omega_j|\mathbf{x})$$

Explanation

The expected (“average”) loss for taking an action (choosing a class) given an input vector, for a given *conditional loss function* (*lambda*)

Decision Functions and Overall Risk

$$R = \int R(\alpha(x)|x)p(x) dx$$

Decision Function or Decision Rule

($\alpha(x)$): takes on the value of exactly one action for each input vector x

Overall Risk

The expected (average) loss associated with a decision rule

Bayes Decision Rule

Idea

Minimize the overall risk, by choosing the action with the least conditional risk for input vector x

Bayes Risk (R^*)

The resulting overall risk produced using this procedure. **This is the best performance that can be achieved given available information.**

Bayes Decision Rule: Two Category Case

Bayes Decision Rule

For each input, select class with least conditional risk, i.e. choose class one if:

$$R(\alpha_1|\mathbf{x}) < R(\alpha_2|\mathbf{x})$$

where

$$\lambda_{ij} = \lambda(\alpha_i|\omega_j)$$

$$R(\alpha_1|\mathbf{x}) = \lambda_{11}P(\omega_1|\mathbf{x}) + \lambda_{12}P(\omega_2|\mathbf{x})$$

$$R(\alpha_2|\mathbf{x}) = \lambda_{21}P(\omega_1|\mathbf{x}) + \lambda_{22}P(\omega_2|\mathbf{x})$$

Alternate Equivalent Expressions of Bayes Decision Rule (“Choose Class One If..”)

Posterior Class Probabilities

$$(\lambda_{21} - \lambda_{11})P(\omega_1|\mathbf{x}) > (\lambda_{12} - \lambda_{22})P(\omega_2|\mathbf{x})$$

Class Priors and Conditional Densities

Produced by applying Bayes Formula to the above, multiplying both sides by $p(\mathbf{x})$

$$(\lambda_{21} - \lambda_{11})p(\mathbf{x}|\omega_1)P(\omega_1) > (\lambda_{12} - \lambda_{22})p(\mathbf{x}|\omega_2)P(\omega_2)$$

Likelihood Ratio $\frac{p(\mathbf{x}|\omega_1)}{p(\mathbf{x}|\omega_2)} > \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \frac{P(\omega_2)}{P(\omega_1)}$

The Zero-One Loss

Definition

$$\lambda(\alpha_i|\omega_j) = \begin{cases} 0 & i = j \\ 1 & i \neq j \end{cases} \quad i, j = 1, \dots, c$$

Conditional Risk for Zero-One Loss

$$R(\alpha_i|\mathbf{x}) = \sum_{j=1}^c \lambda(\alpha_i|\omega_j)P(\omega_j|\mathbf{x}) = \sum_{j \neq i} P(\omega_j|\mathbf{x}) = 1 - P(\omega_i|\mathbf{x})$$

Bayes Decision Rule (min. error rate)

Decide ω_i if $P(\omega_i|\mathbf{x}) > P(\omega_j|\mathbf{x})$ for all $j \neq i$

Example: Likelihood Ratio

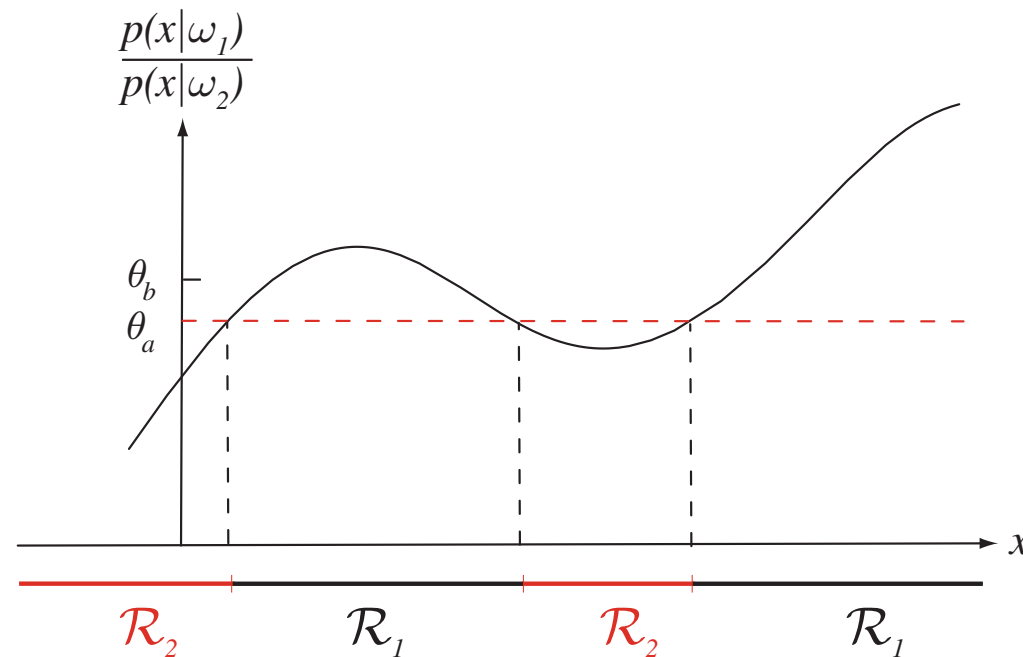


FIGURE 2.3. The likelihood ratio $p(x|\omega_1)/p(x|\omega_2)$ for the distributions shown in Fig. 2.1. If we employ a zero-one or classification loss, our decision boundaries are determined by the threshold θ_a . If our loss function penalizes miscategorizing ω_2 as ω_1 patterns more than the converse, we get the larger threshold θ_b , and hence \mathcal{R}_1 becomes smaller. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Bayes Classifiers

Recall the “Canonical Model”

Decide class i if:

$$g_i(\mathbf{x}) > g_j(\mathbf{x}) \quad \text{for all } j \neq i$$

For Bayes Classifiers

Use the first discriminant def'n below for general case, second for zero-one loss

$$g_i(\mathbf{x}) = -R(\alpha_i|\mathbf{x})$$

$$g_i(\mathbf{x}) = P(\omega_i|\mathbf{x})$$

Equivalent Discriminants for Zero-One Loss (Minimum-Error-Rate)

Trade-off

Simplicity of understanding vs. computation

$$g_i(\mathbf{x}) = P(\omega_i|\mathbf{x}) = \frac{p(\mathbf{x}|\omega_i)P(\omega_i)}{\sum_{j=1}^c p(\mathbf{x}|\omega_j)P(\omega_j)}$$

$$g_i(\mathbf{x}) = p(\mathbf{x}|\omega_i)P(\omega_i)$$

$$g_i(\mathbf{x}) = \ln p(\mathbf{x}|\omega_i) + \ln P(\omega_i)$$

Discriminants for Two Categories

For Two Categories

We can use a single discriminant function, with decision rule: choose class one if the discriminant returns a value > 0 .

Example: Zero-One Loss

$$g(\mathbf{x}) = P(\omega_1|\mathbf{x}) - P(\omega_2|\mathbf{x})$$

$$g(\mathbf{x}) = \ln \frac{p(\mathbf{x}|\omega_1)}{p(\mathbf{x}|\omega_2)} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

Example: Decision Regions for Binary Classifier

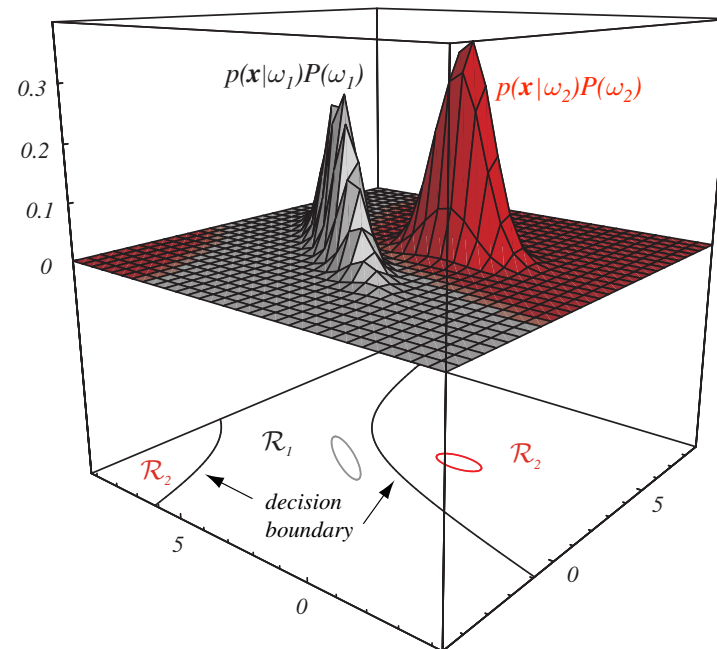


FIGURE 2.6. In this two-dimensional two-category classifier, the probability densities are Gaussian, the decision boundary consists of two hyperbolas, and thus the decision region \mathcal{R}_2 is not simply connected. The ellipses mark where the density is $1/e$ times that at the peak of the distribution. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

The (Univariate) Normal Distribution

Why are Gaussians so Useful?

They represent many probability distributions in nature quite accurately. In our case, when patterns can be represented as random variations of an ideal prototype (represented by the mean feature vector)

- Everyday examples: height, weight of a population

Univariate Normal Distribution

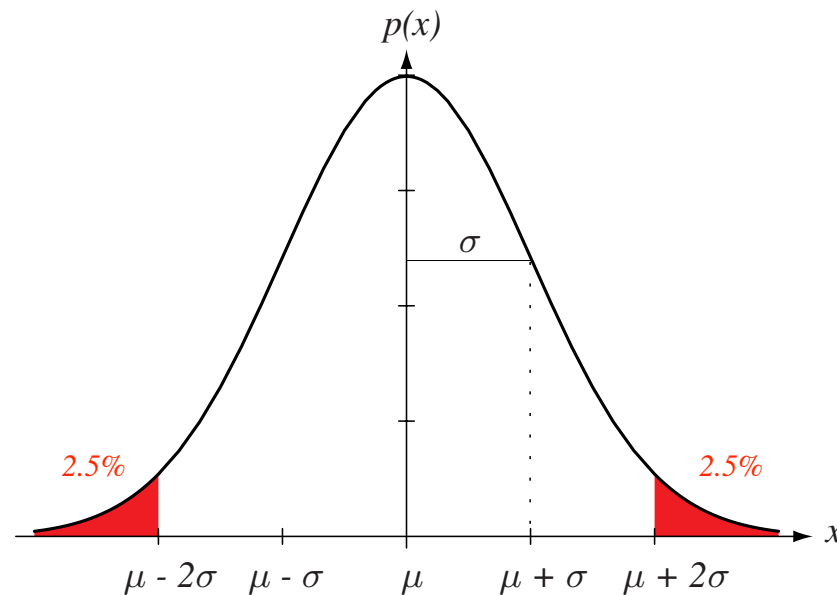


FIGURE 2.7. A univariate normal distribution has roughly 95% of its area in the range $|x - \mu| \leq 2\sigma$, as shown. The peak of the distribution has value $p(\mu) = 1/\sqrt{2\pi}\sigma$. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Formal Definition

Peak of the Distribution (the mean)

Has value: $\frac{1}{\sqrt{2\pi}\sigma}$

Definition for Univariate Normal

$$P(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right]$$

Def. for mean, variance

$$\mu = \int_{-\infty}^{\infty} x p(x) dx$$

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx$$

Multivariate Normal Density

Informal Definition

A normal distribution over two or more variables (d variables/dimensions)

Formal Definition

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu)^t \Sigma^{-1} (\mathbf{x} - \mu) \right]$$

$$\mu = \int_{-\infty}^{\infty} \mathbf{x} p(\mathbf{x}) d\mathbf{x}$$

$$\Sigma = \int (\mathbf{x} - \mu)(\mathbf{x} - \mu)^t p(\mathbf{x}) d\mathbf{x}$$

The Covariance Matrix

(Σ)

For our purposes...

Assume matrix is positive definite, so the determinant of the matrix is always positive

Matrix Elements

- Main diagonal: variances for each individual variable
- Off-diagonal: covariances of each variable pairing i & j (note: values are repeated, as matrix is symmetric)

Independence and Correlation

For multivariate normal covariance matrix

- Off-diagonal entries with a value of 0 indicate uncorrelated variables, that are *statistically* independent (variables likely do not influence one another)
- Roughly speaking, covariance positive if two variables increase together (positive correlation), negative if one variable decreases when the other increases (negative correlation)

A Two-Dimensional Gaussian Distribution, with Samples Shown

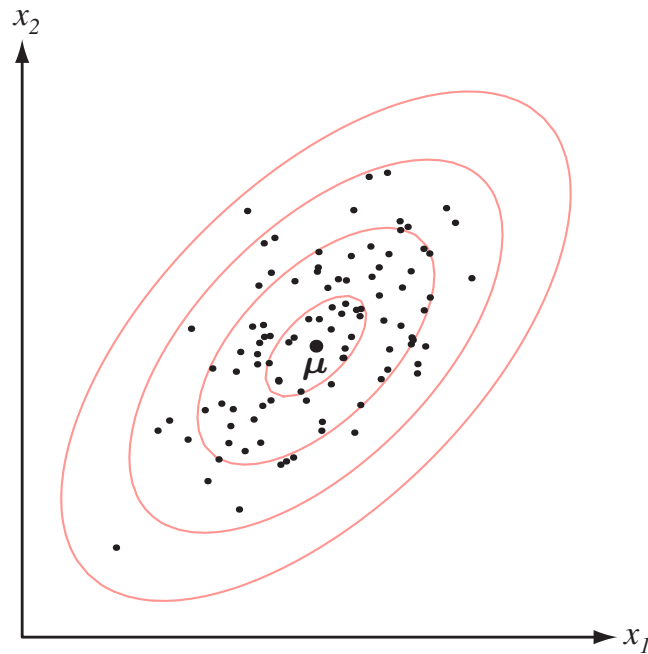
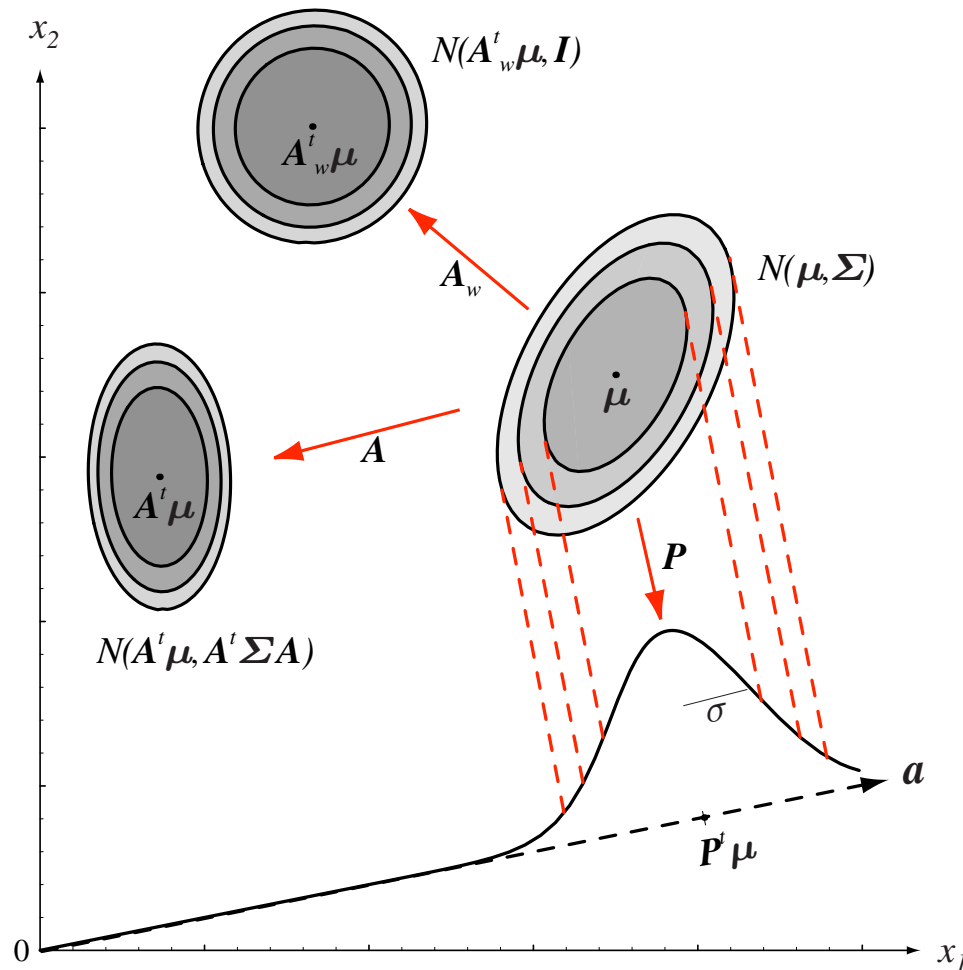


FIGURE 2.9. Samples drawn from a two-dimensional Gaussian lie in a cloud centered on the mean μ . The ellipses show lines of equal probability density of the Gaussian. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Linear Transformations in a 2D Feature Space



Discriminant Functions ($g_i(\mathbf{x})$) for the Normal Density

Discriminant Functions

We will consider three special cases for:

- normally distributed features, and
- minimum-error-rate classification (0-1 loss)

Recall: $g_i(\mathbf{x}) = \ln p(\mathbf{x}|\omega_i) + \ln P(\omega_i)$

if $p(x|\omega_i) \sim N(\mu_i, \Sigma_i)$ then approx. $p(\mathbf{x}|\omega_i)$

using: $p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mu)^t \Sigma^{-1} (\mathbf{x} - \mu) \right]$

Minimum Error-Rate Discriminant Function for Multivariate Gaussian Feature Distributions

In (natural log) of

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

gives a general form for our discriminant functions:

$$g_i(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^t \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i)$$

Special Cases for Binary Classification

Purpose

Overview of commonly assumed cases for feature likelihood densities, $p(\mathbf{x}|\omega_i)$

- **Goal:** eliminate common additive constants in discriminant functions. These do not affect the classification decision (i.e. define $g_i(\mathbf{x})$ providing “just the differences”)
- Also, look at resulting **decision surfaces** (defined by $g_i(\mathbf{x}) = g_j(\mathbf{x})$)

Three Special Cases

1. Statistically independent features, identically distributed Gaussians for each class
2. Identical covariances for each class
3. Arbitrary covariances

Case I: $\Sigma_i = \sigma^2 I$

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

Remove:

Items in red: same across classes (“unimportant additive constants”)

Inverse of Covariance Matrix: $\Sigma_i^{-1} = (1/\sigma^2)I$

Only effect is to scale vector product by $1/\sigma^2$

Discriminant function:

$$g_i(x) = -\frac{(\mathbf{x} - \mu_i)^t (\mathbf{x} - \mu_i)}{2\sigma^2} + \ln P(\omega_i)$$

$$g_i(x) = -\frac{1}{2\sigma^2} [\mathbf{x}^t \mathbf{x} - 2\mu_i^t \mathbf{x} + \mu_i^t \mu_i] + \ln P(\omega_i)$$

Case I: $\Sigma_i = \sigma^2 I$

Linear Discriminant Function

Produced by factoring the previous form

$$g_i(x) = \mathbf{w}_i^t \mathbf{x} + \omega_{i0}$$

$$g_i(x) = \frac{1}{\sigma^2} \mu_i^t \mathbf{x} - \frac{1}{2\sigma^2} \mu_i^t \mu_i + \ln P(\omega_i)$$

Threshold or Bias for Class i : ω_{i0}

Change in prior translates decision boundary

Case I: $\Sigma_i = \sigma^2 I$

Decision Boundary: $g_i(x) = g_j(x)$

$$\mathbf{w}^t (\mathbf{x} - \mathbf{x}_0) = 0$$

$$(\mu_i - \mu_j)^t \left(\mathbf{x} - \left(\frac{1}{2}(\mu_i + \mu_j) - \frac{\sigma^2}{(\mu_i - \mu_j)^t (\mu_i - \mu_j)} \ln \frac{P(\omega_i)}{P(\omega_j)} (\mu_i - \mu_j) \right) \right)$$

- Decision boundary goes through \mathbf{x}_0 along line between means, orthogonal to this line
- If priors equal, \mathbf{x}_0 between means (*minimum distance classifier*), otherwise \mathbf{x}_0 shifted
- If variance small relative to distance between means, priors have limited effect on boundary location

Case I: Statistically Independent Features with Identical Variances

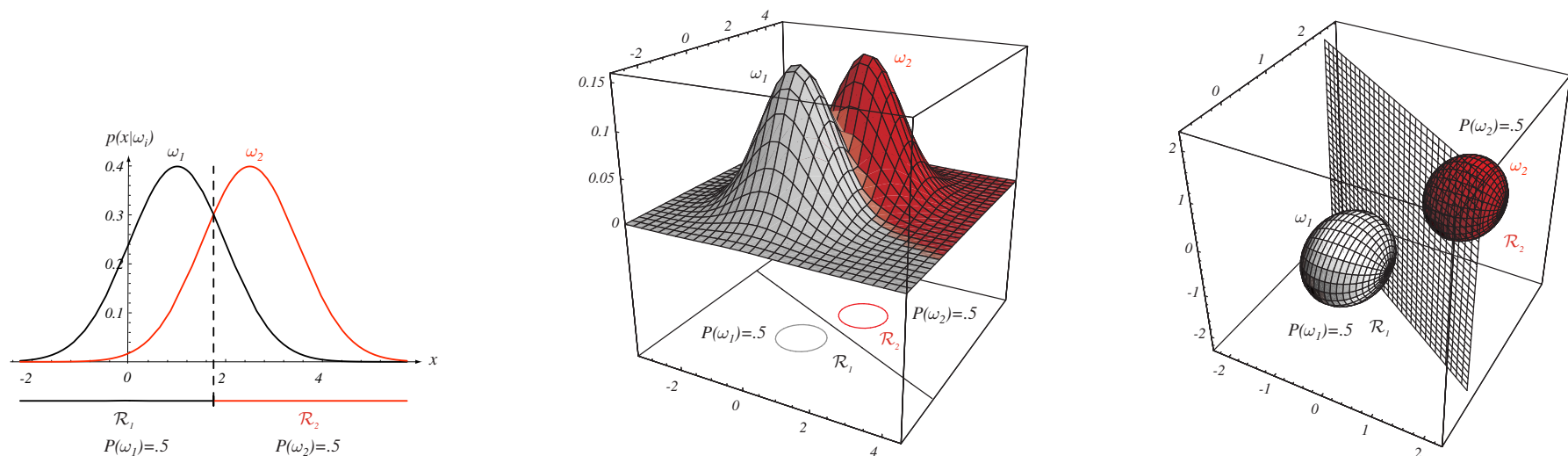
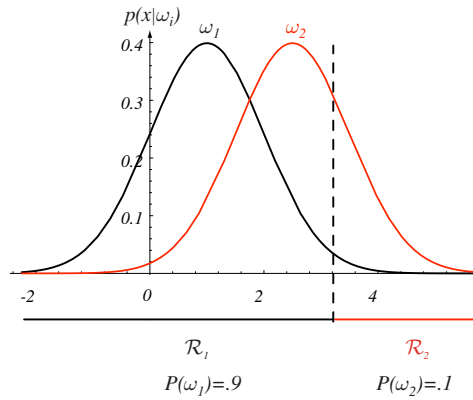
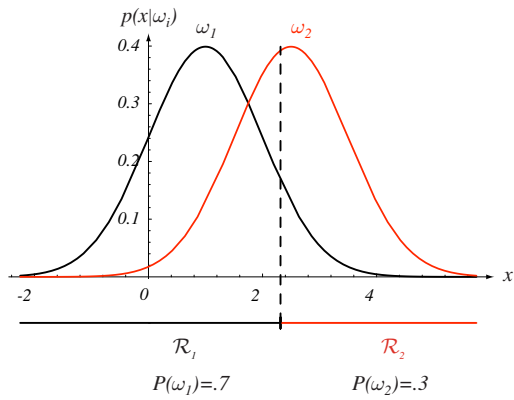
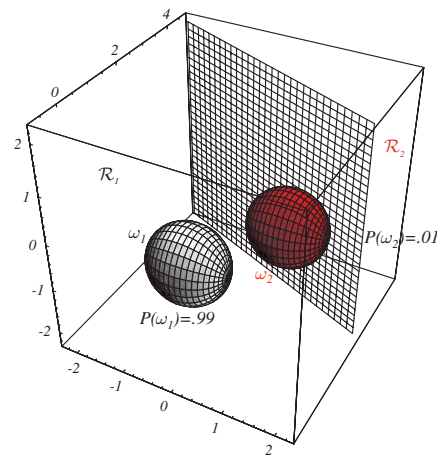
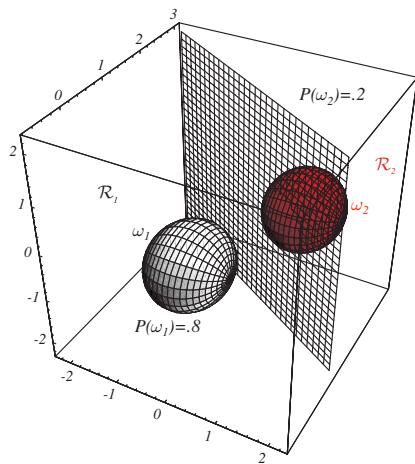
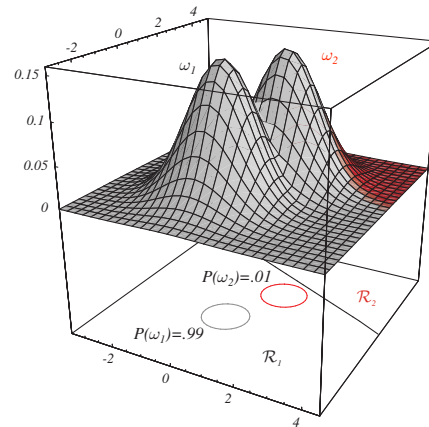
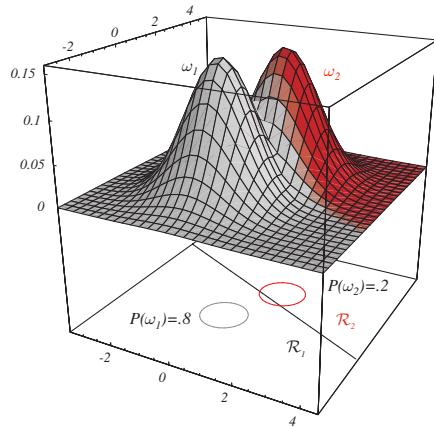


FIGURE 2.10. If the covariance matrices for two distributions are equal and proportional to the identity matrix, then the distributions are spherical in d dimensions, and the boundary is a generalized hyperplane of $d - 1$ dimensions, perpendicular to the line separating the means. In these one-, two-, 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.



Example: Translation of Decision Boundaries Through Changing Priors



Case II: Identical Covariances, $\Sigma_i = \Sigma$

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

Remove

Terms in red; as in Case I these can be ignored (same across classes)

Squared Mahalanobis Distance (yellow)

Distance from \mathbf{x} to mean for class i , taking covariance into account; defines contours of fixed density

Case II: Identical Covariances, $\Sigma_i = \Sigma$

Expansion of squared Mahalanobis distance

$$\begin{aligned} & (\mathbf{x} - \mu_i)^t \Sigma^{-1} (\mathbf{x} - \mu_i) \\ &= \mathbf{x}^t \Sigma^{-1} \mathbf{x} - \mathbf{x}^t \Sigma^{-1} \mu_i - \mu_i^t \Sigma^{-1} \mathbf{x} + \mu_i^t \Sigma^{-1} \mu_i \\ &= \mathbf{x}^t \Sigma^{-1} \mathbf{x} - 2(\Sigma^{-1} \mu_i)^t \mathbf{x} + \mu_i^t \Sigma^{-1} \mu_i \end{aligned}$$

the last step comes from symmetry of the covariance matrix and thus its inverse:

$$\Sigma^t = \Sigma, (\Sigma^{-1})^t = \Sigma^{-1}$$

Once again, term above in red is an additive constant independent of class, and can be removed

Case II: Identical Covariances, $\Sigma_i = \Sigma$

Linear Discriminant Function

$$g_i(x) = \mathbf{w}_i^t \mathbf{x} + \omega_{i0}$$

$$g_i(\mathbf{x}) = (\Sigma^{-1} \mu_i)^t \mathbf{x} - \frac{1}{2} \mu_i^t \Sigma^{-1} \mu_i + \ln P(\omega_i)$$

Decision Boundary: $g_i(x) = g_j(x)$

$$\mathbf{w}^t (\mathbf{x} - \mathbf{x}_0) = 0$$

$$(\Sigma^{-1} (\mu_i - \mu_j))^t \left(\mathbf{x} - \left(\frac{1}{2} (\mu_i + \mu_j) - \frac{\ln [P(\omega_i) / P(\omega_j)]}{(\mu_i - \mu_j)^t \Sigma^{-1} (\mu_i - \mu_j)} (\mu_i - \mu_j) \right) \right) = 0$$

Case II: Identical Covariances, $\Sigma_i = \Sigma$

Notes on Decision Boundary

- As for Case I, passes through point x_0 lying on the line between the two class means. Again, x_0 in the middle if priors identical
- Hyperplane defined by boundary generally not orthogonal to the line between the two means

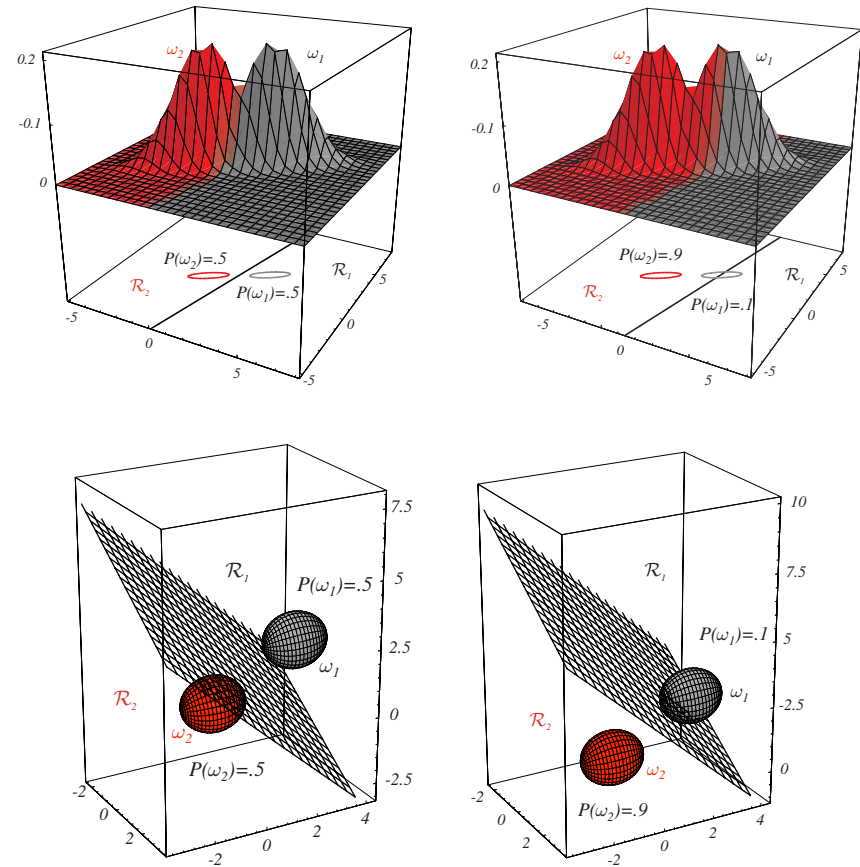


FIGURE 2.12. Probability densities (indicated by the surfaces in two dimensions as ellipsoidal surfaces in three dimensions) and decision regions for equal but asymmetric Gaussian distributions. The decision hyperplanes need not be perpendicular to the line connecting the means. From: Richard O. Duda,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

Case III: arbitrary Σ_i

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

Remove

Can only remove the one term in red above

Discriminant Function (quadratic)

$$g_i(x) = x^t W_i x + w_i^t x + w_{i0}$$

$$g_i(x) = x^t \left(-\frac{1}{2} \Sigma_i^{-1}\right) x + \left(\Sigma_i^{-1} \mu_i\right)^t x - \frac{1}{2} \mu_i^t \Sigma_i^{-1} \mu_i - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

Case III: arbitrary Σ_i

Decision Boundaries

Are hyperquadrics: can be hyperplanes, hyperplane pairs, hyperspheres, hyperellipsoids, hyperparabaloids, hyperhyperparabaloids

Decision Regions

Need not be simply connected, even in one dimension (next slide)

Case 3: Arbitrary Covariances

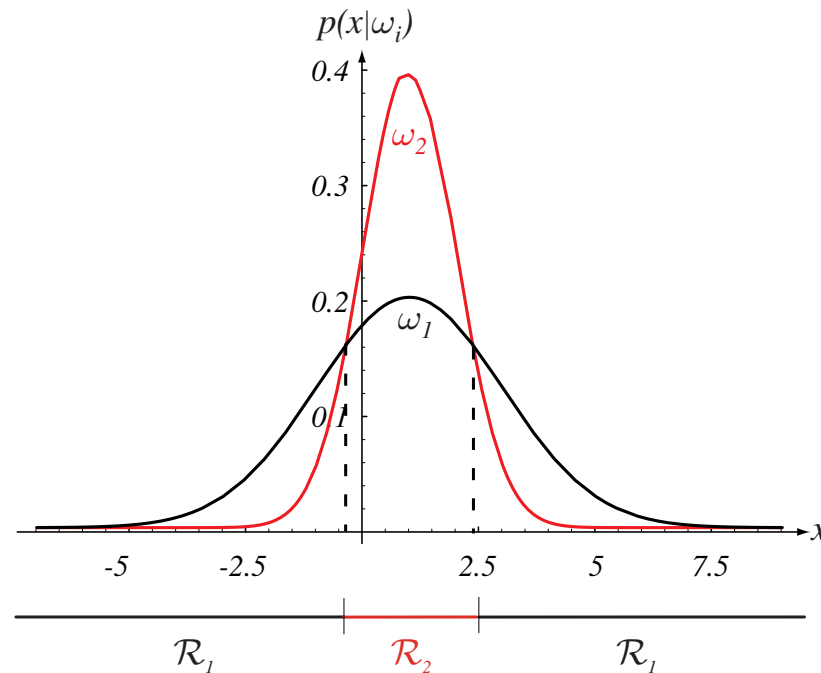
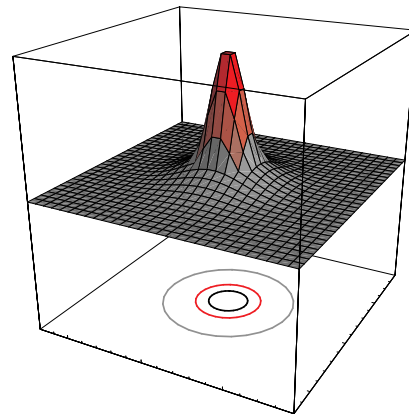
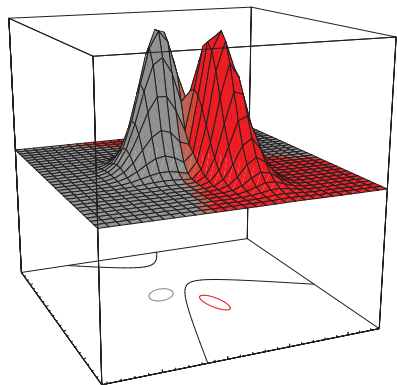
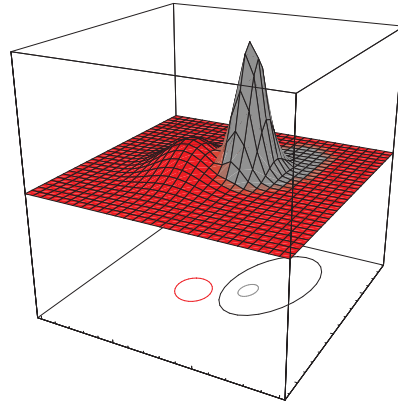
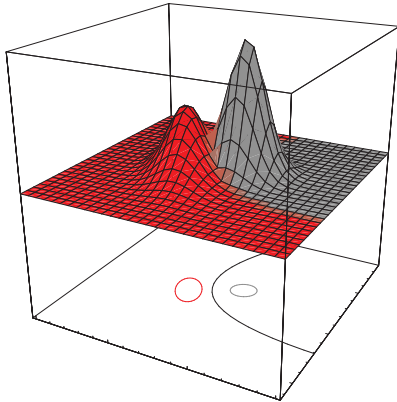
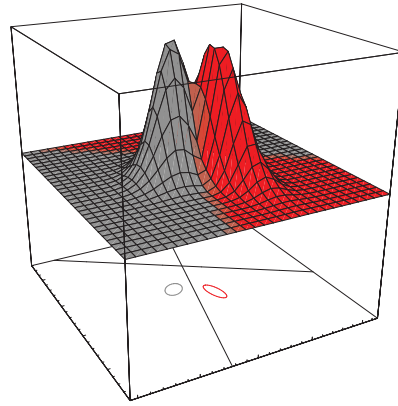
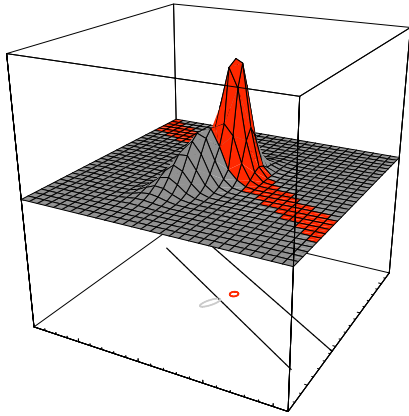


FIGURE 2.13. Non-simply connected decision regions can arise in one dimensions for Gaussians having unequal variance. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



See Fig. 2.15 for
3D cases

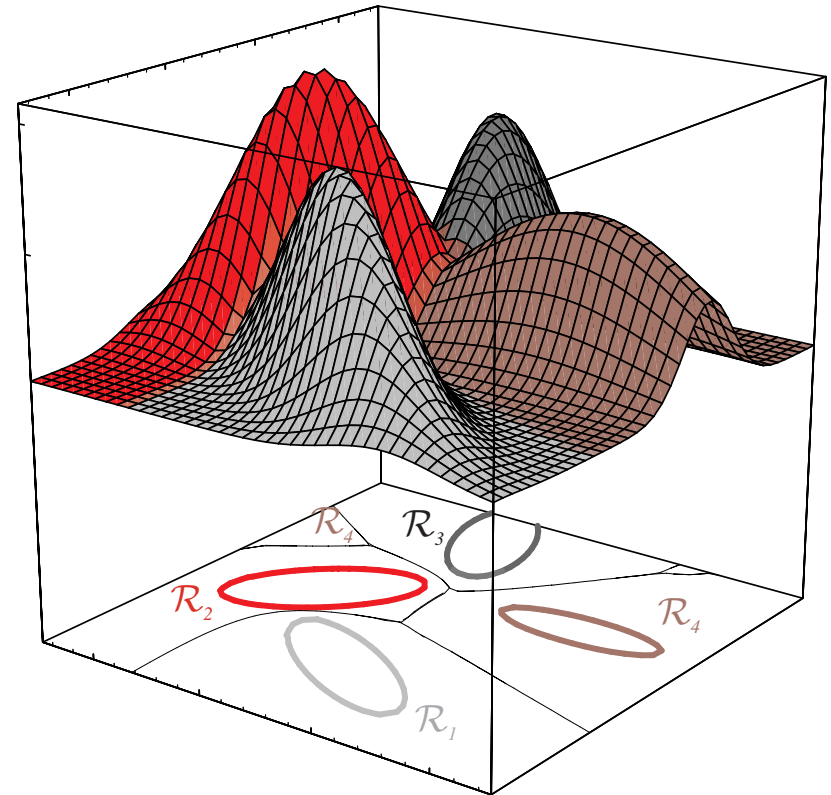
More than Two Categories

Decision Boundary

Defined by two most likely classes for each segment

Other Distributions

Possible; underlying Bayesian Decision Theory is unmodified, however



Discrete Features

Roughly speaking...

Replace probability densities by probability mass functions. Expressions using integrals are changed to use summations, e.g.

$$\int p(\mathbf{x}|\omega_j) d\mathbf{x} \quad \sum_x P(\mathbf{x}|\omega_j)$$

Bayes Formula $P(\omega_j|\mathbf{x}) = \frac{P(\mathbf{x}|\omega_j)P(\omega_j)}{P(\mathbf{x})}$

$$P(\mathbf{x}) = \sum_{j=1}^c P(\mathbf{x}|\omega_j)P(\omega_j)$$

Example: Independent Binary Features

Binary Feature Vector

$\mathbf{x} = \{x_1, \dots, x_d\}$ of 0/1 -valued features, where each x_i is 0/1 with probability: $p_i = Pr[x_i = 1 | \omega_1]$

Conditional Independence

Assume that *given a class*, the features are independent

Likelihood Function

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