

Bayes Decision Theory

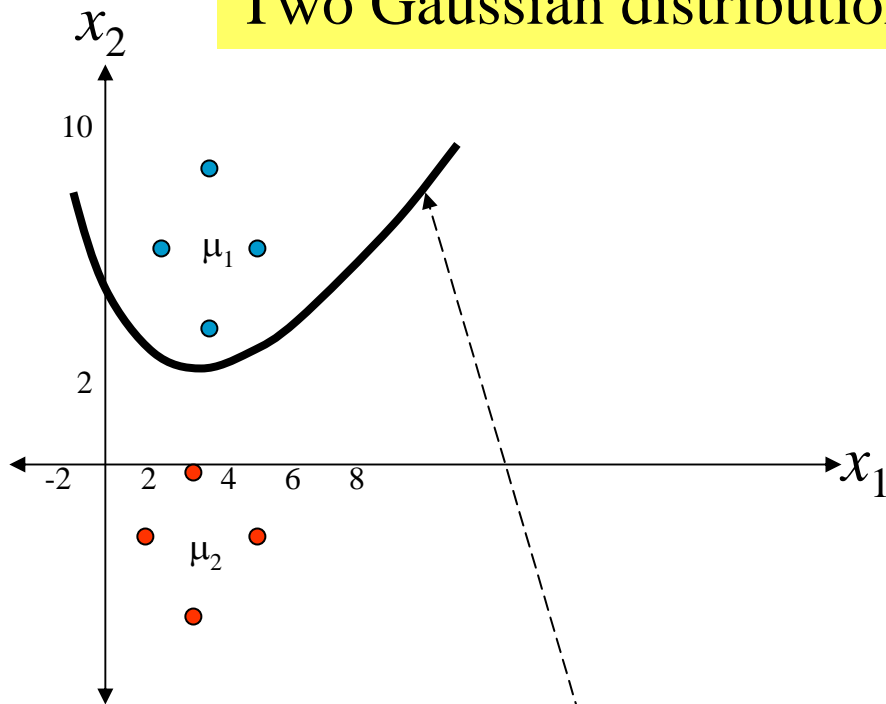
Bayes Error Rate and Error Bounds
Receiver Operating Characteristics

Discrete Features

Missing Features

Example of Bayes Decision Boundary

Two Gaussian distributions each with four data points



$$x_2 = 3.514 - 1.125x_1 + 0.1865x_1^2$$

Decision Boundary assumes
 $P(\omega_1) = P(\omega_2) = 0.5$

$$\mu_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix} \quad \Sigma_1 = \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix}$$

$$\mu_2 = \begin{bmatrix} 3 \\ -2 \end{bmatrix} \quad \Sigma_2 = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$$

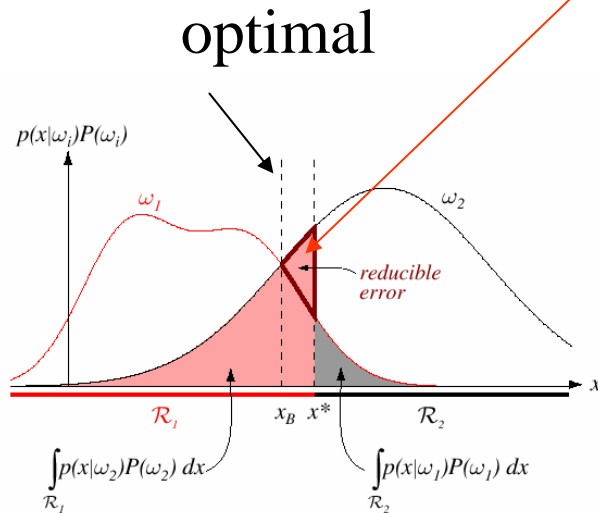
Inverse Matrices

$$\Sigma_1^{-1} = \begin{pmatrix} 2 & 0 \\ 0 & 1/2 \end{pmatrix}$$

$$\Sigma_2^{-1} = \begin{pmatrix} 1/2 & 0 \\ 0 & 1/2 \end{pmatrix}$$

Bayes Error Rate

Two Class Case



$$\begin{aligned}
 P(\text{error}) &= P(\mathbf{x} \in \mathcal{R}_2, \omega_1) + P(\mathbf{x} \in \mathcal{R}_1, \omega_2) \\
 &= P(\mathbf{x} \in \mathcal{R}_2 | \omega_1)P(\omega_1) + P(\mathbf{x} \in \mathcal{R}_1 | \omega_2)P(\omega_2) \\
 &= \int_{\mathcal{R}_2} p(\mathbf{x} | \omega_1)P(\omega_1) d\mathbf{x} + \int_{\mathcal{R}_1} p(\mathbf{x} | \omega_2)P(\omega_2) d\mathbf{x}.
 \end{aligned}$$

Stated for the multidimensional case (regions \mathcal{R}_1 and \mathcal{R}_2 not easy to specify)

$$\begin{aligned}
 P(\text{correct}) &= \sum_{i=1}^c P(\mathbf{x} \in \mathcal{R}_i, \omega_i) \\
 &= \sum_{i=1}^c P(\mathbf{x} \in \mathcal{R}_i | \omega_i)P(\omega_i) \\
 &= \sum_{i=1}^c \int_{\mathcal{R}_i} p(\mathbf{x} | \omega_i)P(\omega_i) d\mathbf{x}.
 \end{aligned}$$

Components of $P(\text{error})$ for equal priors and non-optimal decision point x^*

Error Bounds for Normal Densities

- Evaluation of Error Integrals difficult
 - Discontinuous decision regions
 - Instead, obtain bounds on error rate
- Useful inequality in obtaining a bound
 - Minimum of two integers is smaller than the square root of their product, more generally

$$\min[a, b] \leq a^\beta b^{1-\beta} \quad \text{for } a, b \geq 0 \text{ and } 0 \leq \beta \leq 1$$

- **Proof:** If $a \geq b$ then $(a/b)^\beta \geq 1$
or $(a/b)^\beta b \geq b$
or $a^\beta b^{1-\beta} \geq b$

Chernoff Bound

$$P(\text{error}) \leq P^\beta(\omega_1)P^{1-\beta}(\omega_2) \int p^\beta(x|\omega_1)p^{1-\beta}(x|\omega_2)dx \quad \text{for } 0 \leq \beta \leq 1$$

Note that integral is over all space—no need to impose integration limits. If conditional probabilities are normal, the integral can be evaluated analytically, yielding

$$\int p^\beta(x|\omega_1)p^{1-\beta}(x|\omega_2)dx = e^{-k(\beta)}$$

where

$$k(\beta) = \frac{\beta(1-\beta)}{2} (\mu_2 - \mu_1)' [\beta\Sigma_1 + (1-\beta)\Sigma_2]^{-1} (\mu_2 - \mu_1) + \frac{1}{2} \ln \frac{|\beta\Sigma_1 + (1-\beta)\Sigma_2|}{|\Sigma_1|^\beta |\Sigma_2|^{1-\beta}}. \quad (75)$$

Bhattacharya Bound

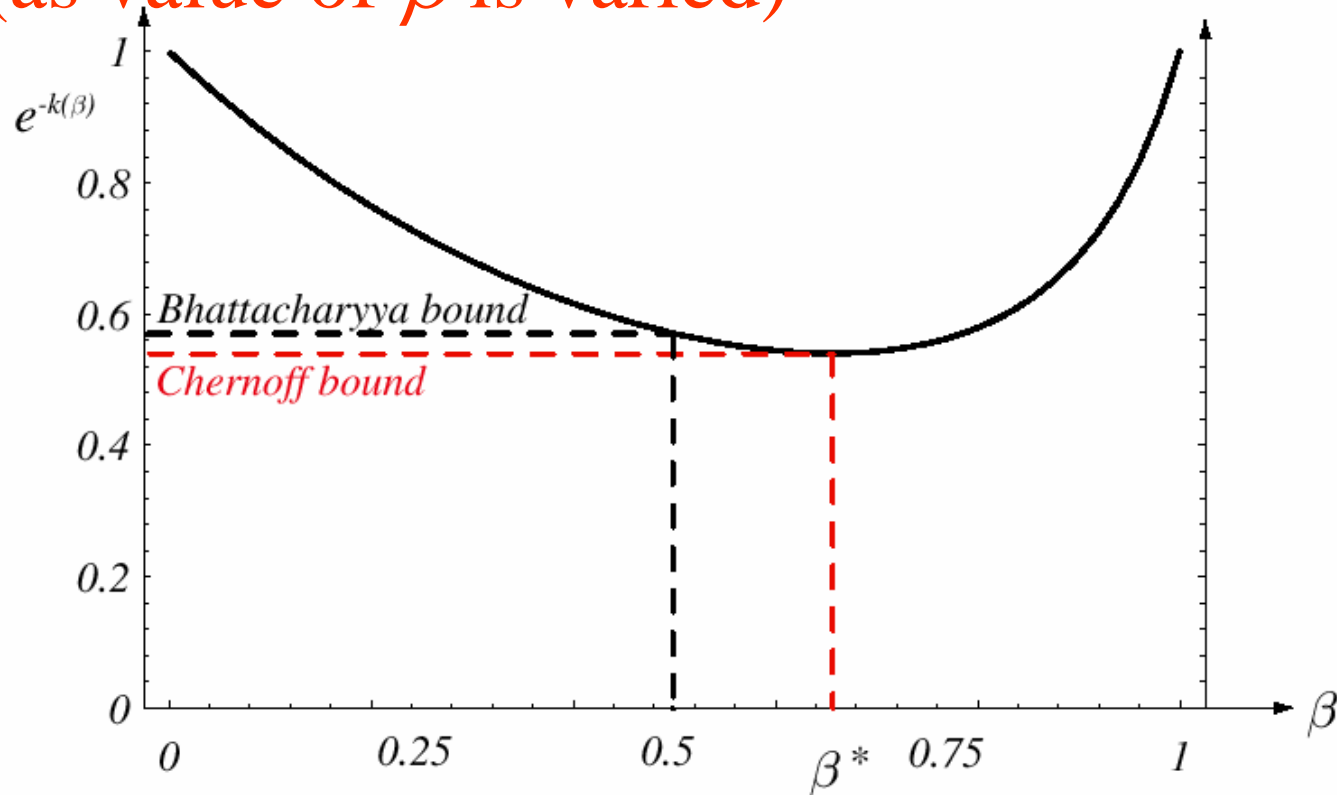
- Special case of Chernoff bound where $\beta = 0.5$

$$\begin{aligned} P(\text{error}) &\leq \sqrt{P(\omega_1)P(\omega_2)} \int \sqrt{p(\mathbf{x}|\omega_1)p(\mathbf{x}|\omega_2)} d\mathbf{x} \\ &= \sqrt{P(\omega_1)P(\omega_2)} e^{-k(1/2)}, \end{aligned}$$

where by Eq. 75 we have for the Gaussian case:

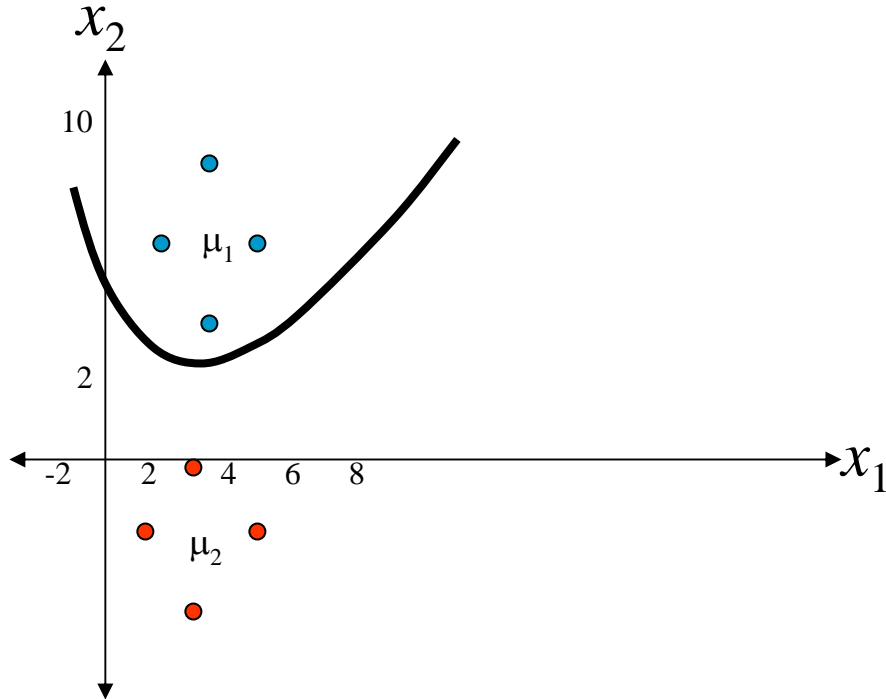
$$k(1/2) = 1/8(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)' \left[\frac{\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2}{2} \right]^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \frac{1}{2} \ln \frac{\left| \frac{\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2}{2} \right|}{\sqrt{|\boldsymbol{\Sigma}_1||\boldsymbol{\Sigma}_2|}}.$$

Bhattacharya versus Chernoff Error Bounds (as value of β is varied)



Chernoff bound is never looser than the Bhattacharyya bound. Here Chernoff bound is at $\beta^* = 0.66$ and is slightly tighter than the Bhattacharyya bound ($\beta = 0.5$)

Example of Bhattacharya bound with normal densities



$$\mu_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix} \quad \Sigma_1 = \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix}$$

$$\mu_2 = \begin{bmatrix} 3 \\ -2 \end{bmatrix} \quad \Sigma_2 = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$$

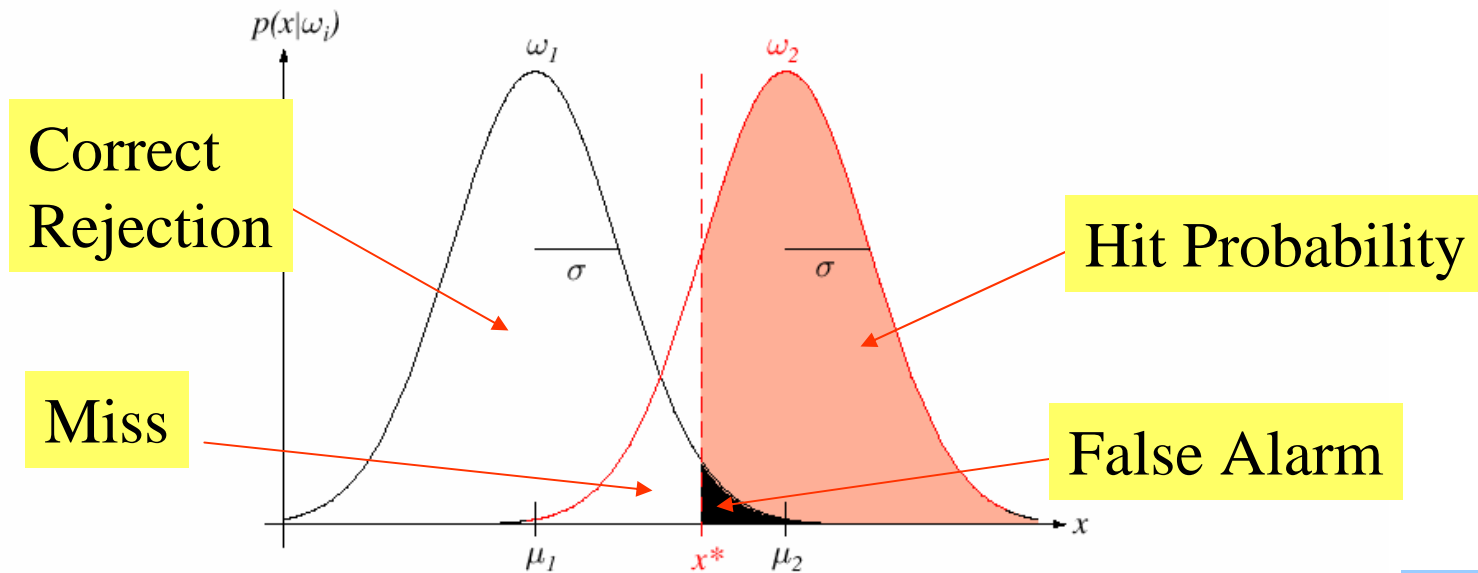
$$k(1/2) = 4.06$$
$$P(\text{error}) \leq 0.0087$$

Receiver Operating Characteristics

- Distance between Gaussian distributions useful in experimental psychology, radar detection, medical diagnosis
- Interested in detecting a weak pulse, or dim flash of light
- Detector detects a signal whose mean value is μ_1 when signal is absent and μ_2 when signal is present

$$p(x / \omega_i) \sim N(\mu_i, \sigma^2)$$

Four Types of Probability in Two-Class Discrimination



When no signal present: $p(x/\omega_1) \sim N(\mu_1, \sigma^2)$

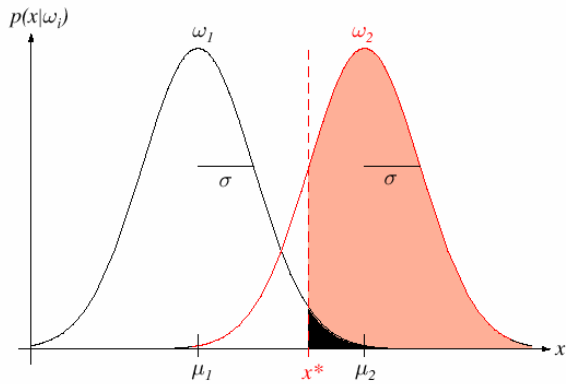
When signal present: $p(x/\omega_1) \sim N(\mu_1, \sigma^2)$

Decision threshold determines probabilities of hit and false alarm

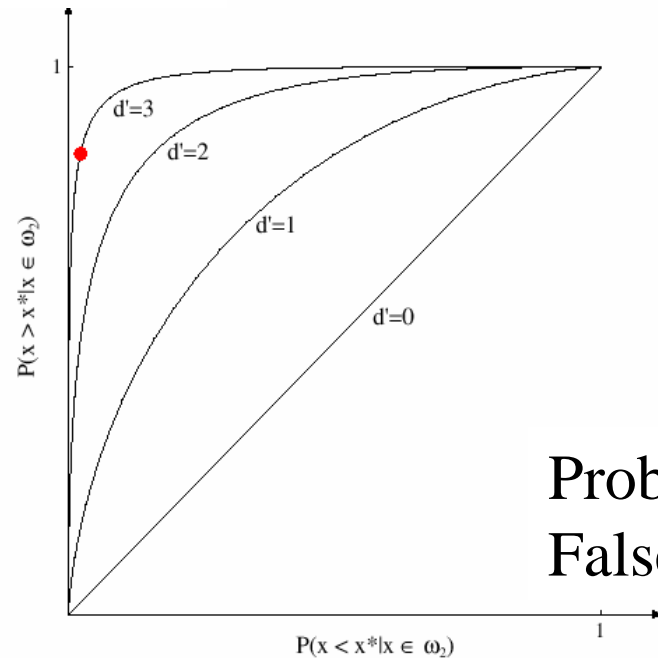
Discriminability

$$d' = \frac{|\mu_2 - \mu_1|}{\sigma}$$

ROC Curve



Probability
of Hit

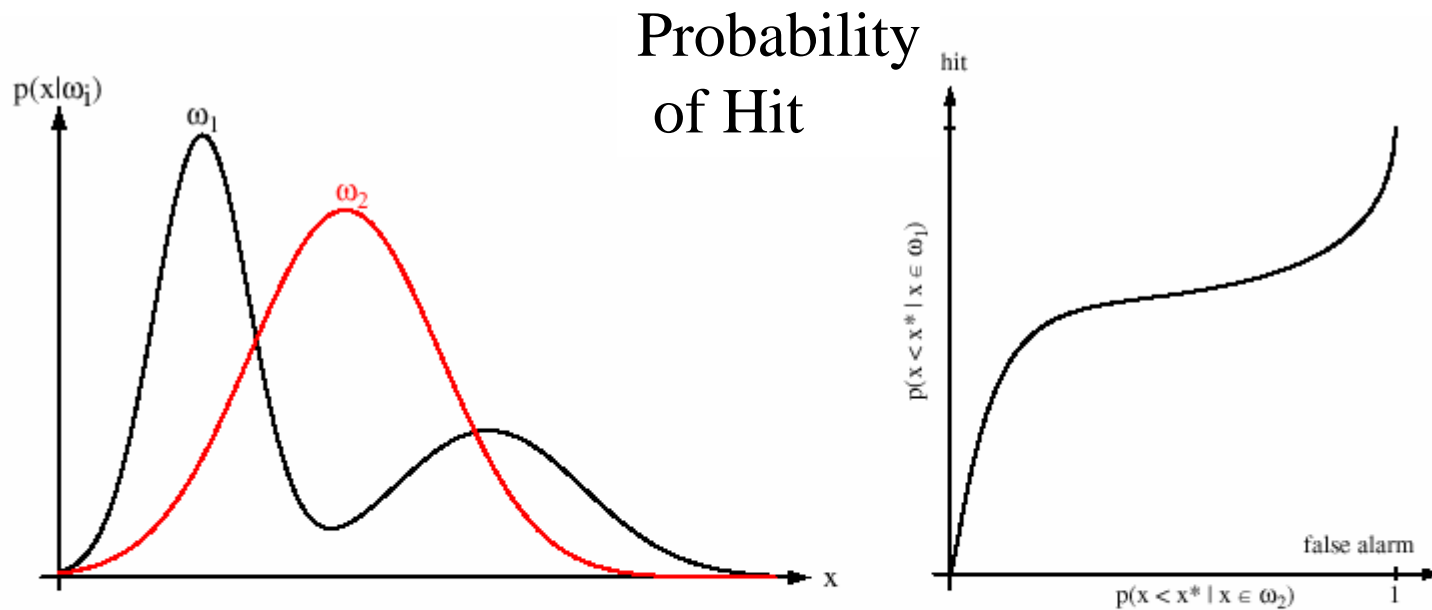


Probability of
False Alarm

When no signal present: $p(x/\omega_1) \sim N(\mu_1, \sigma^2)$
When signal present: $p(x/\omega_2) \sim N(\mu_2, \sigma^2)$
Decision threshold determines probabilities of hit
and false alarm

$$d' = \frac{|\mu_2 - \mu_1|}{\sigma}$$

ROCs need not be symmetric when distributions are not Gaussian



Probability of
False Alarm

Bayes Decision Theory: Discrete Features

Bayes Decision Theory – Discrete Features

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

$$v_1, v_2, \dots, v_m$$

- Probability Density Functions replaced by Probabilities

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

where

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

Independent Binary Features

2 category problem

Let $x = [x_1, x_2, \dots, x_d]^t$

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

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

and

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

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

Bayes discriminant function for Independent Binary Features:

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

where :

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

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 $g(\mathbf{x}) > 0$ and ω_2 if $g(\mathbf{x}) \leq 0$

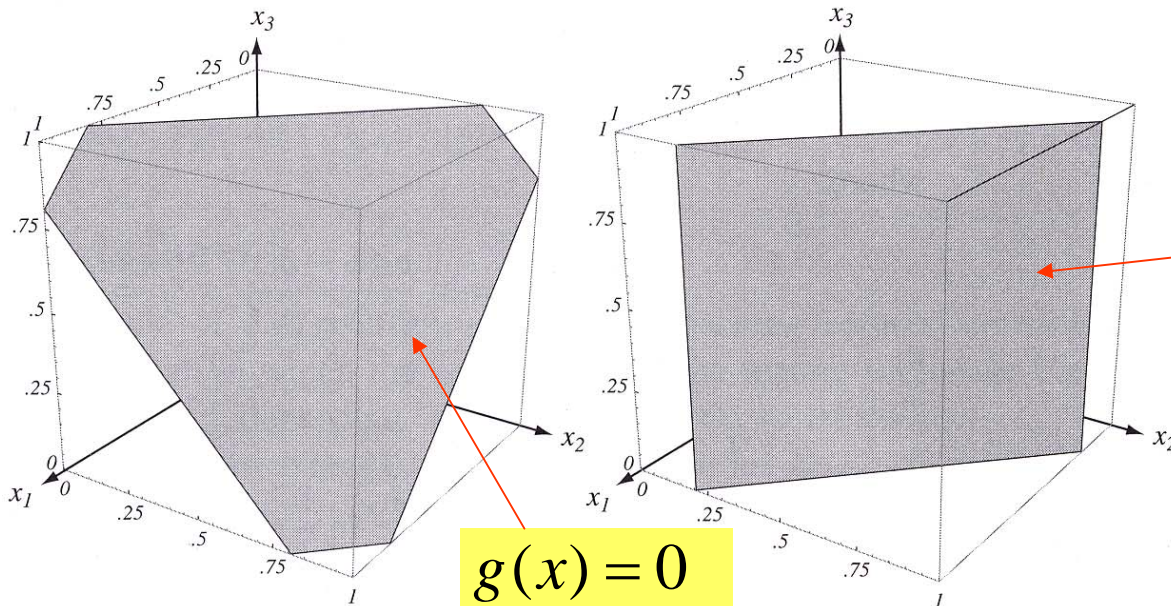
Bayesian Decisions for 3-D Binary Data

$$P(\omega_1) = P(\omega_2) = 0.5$$

$$p_i = 0.8 \quad q_i = 0.5 \quad \text{for } i = 1, 2, 3$$

$$w_i = \ln \frac{.8(1 - .5)}{.5(1 - .8)} = 1.3863$$

$$w_0 = \sum_{i=1}^3 \ln \frac{1 - .8}{1 - .5} + \ln \frac{.5}{.5} = 1.2.$$



$$g(x) = 0$$

$$p_3 = q_3$$

$$w_3 = 0$$

Missing and Noisy Features

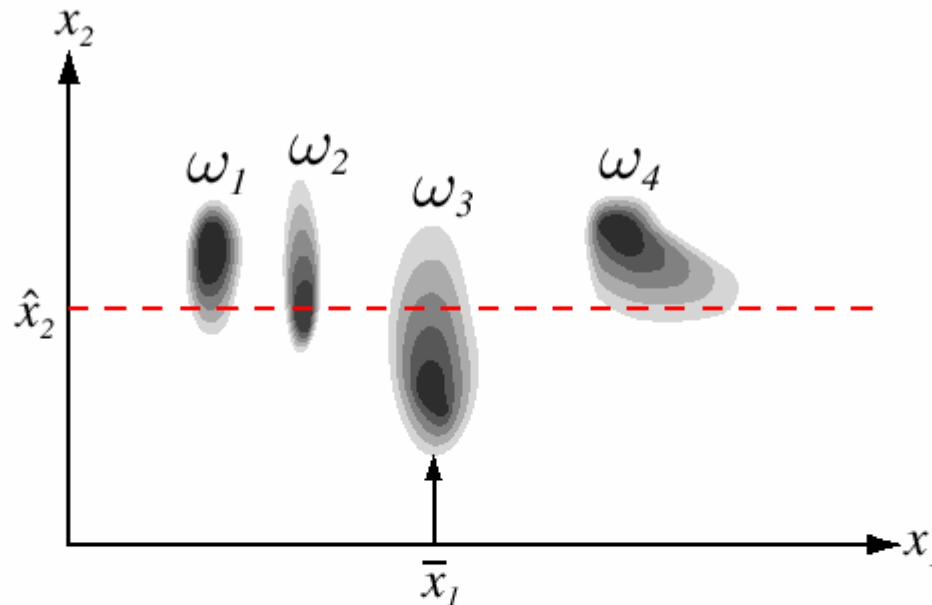
Features are corrupted by a known noise source

Ex: variability of light source may degrade measurement of lightness

Features are missing

Ex: occlusion prevents measurement of length

Example of Missing Feature



Four categories with equal priors and class-conditional distributions. Here x_1 is missing and the other has value x_2 . We want to classify as ω_2 since it has the largest likelihood.

Missing Feature Analysis

good features x_g

bad features (unknown or missing) x_b

$$\begin{aligned} P(\omega_i | \mathbf{x}_g) &= \frac{p(\omega_i, \mathbf{x}_g)}{p(\mathbf{x}_g)} = \frac{\int p(\omega_i, \mathbf{x}_g, \mathbf{x}_b) d\mathbf{x}_b}{p(\mathbf{x}_g)} \\ &= \frac{\int P(\omega_i | \mathbf{x}_g, \mathbf{x}_b) p(\mathbf{x}_g, \mathbf{x}_b) d\mathbf{x}_b}{p(\mathbf{x}_g)} \\ &= \frac{\int g_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}_b}{\int p(\mathbf{x}) d\mathbf{x}_b}, \end{aligned}$$

Marginalize over
all values of
missing feature

This is the Bayes Discriminant Function

Noisy Features

- Uncorrupted good features \mathbf{x}_g
- Noise model $p(\mathbf{x}_b|\mathbf{x}_t)$
 - \mathbf{x}_t = True value of the Observed value \mathbf{x}_b
 - Assume if \mathbf{x}_t were known \mathbf{x}_b would be independent of w_i and \mathbf{x}_g

$$P(\omega_i|\mathbf{x}_g, \mathbf{x}_b) = \frac{\int p(\omega_i, \mathbf{x}_g, \mathbf{x}_b, \mathbf{x}_t) d\mathbf{x}_t}{p(\mathbf{x}_g, \mathbf{x}_b)}.$$

Integral is weighted by the noise model

$$\begin{aligned} P(\omega_i|\mathbf{x}_g, \mathbf{x}_b) &= \frac{\int P(\omega_i|\mathbf{x}_g, \mathbf{x}_t)p(\mathbf{x}_g, \mathbf{x}_t)p(\mathbf{x}_b|\mathbf{x}_t) d\mathbf{x}_t}{\int p(\mathbf{x}_g, \mathbf{x}_t)p(\mathbf{x}_b|\mathbf{x}_t) d\mathbf{x}_t} \\ &= \frac{\int g_i(\mathbf{x})p(\mathbf{x})p(\mathbf{x}_b|\mathbf{x}_t) d\mathbf{x}_t}{\int p(\mathbf{x})p(\mathbf{x}_b|\mathbf{x}_t) d\mathbf{x}_t}, \end{aligned}$$