

Bayes Classification

Richa Singh

Probability

- Conditional probability of A given B:

$$P(A/B) = \frac{P(A,B)}{P(B)}$$

- Deriving chain rule from above:

$$P(A,B) = P(A/B)P(B) = P(B/A)P(A)$$

Bayes Theorem

$$P(A / B) = \frac{P(B / A)P(A)}{P(B)}$$

$$P(B) = P(B, A) + P(B, \bar{A}) = P(B / A)P(A) + P(B / \bar{A})P(\bar{A})$$

Bayes Theorem

$$P(A/B) = \frac{P(B/A)P(A)}{P(B)}$$

$$P(B) = P(B, A) + P(B, \bar{A}) = P(B/A)P(A) + P(B/\bar{A})P(\bar{A})$$

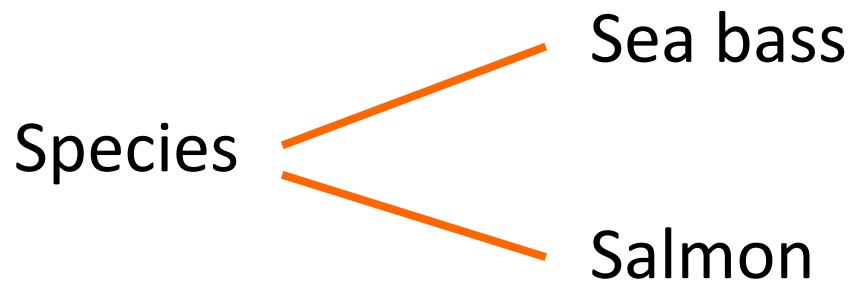
$$P(Disease/Symptom) = \frac{P(Symptom/Disease)P(Disease)}{P(Symptom)}$$

$$P(Symptom) = P(Symptom/Disease)P(Disease) + \\ P(Symptom/NoDisease)P(NoDisease)$$

Bayes Classification

An Example

- “Sorting incoming fish on a conveyor according to species using optical sensing”



Let us build a machine learning system that classifies between Sea Bass and Salmon

Fish Classification: Salmon vs. Sea Bass

- Set up a camera and take some sample images
- Preprocessing involves image enhancement and segmentation;
 - separate touching or occluding fishes and
 - extract fish contour

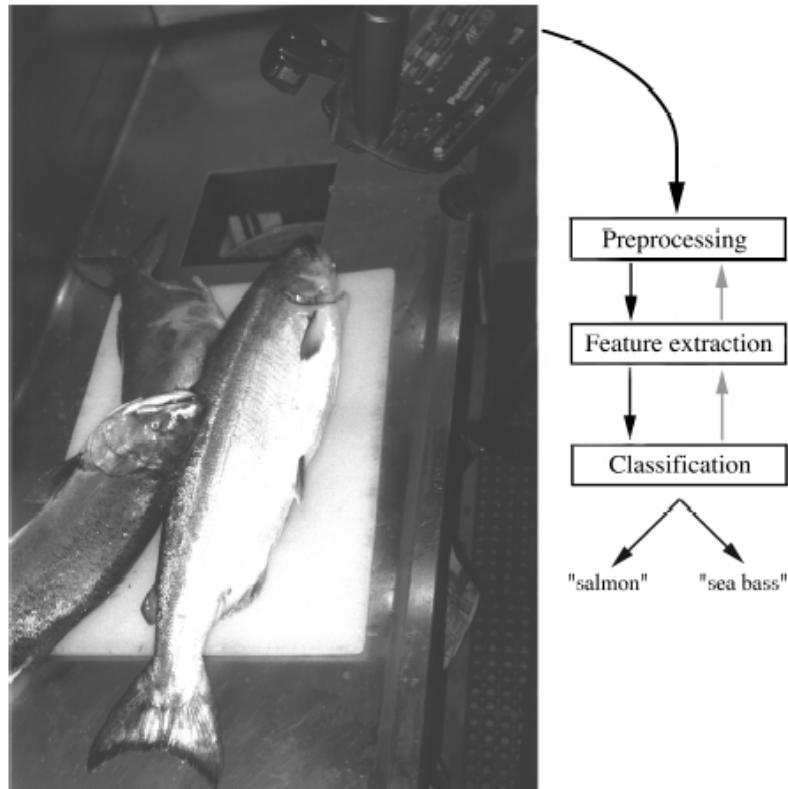


FIGURE 1.1. The objects to be classified are first sensed by a transducer (camera), whose signals are preprocessed. Next the features are extracted and finally the classification is emitted, here either "salmon" or "sea bass." Although the information flow is often chosen to be from the source to the classifier, some systems employ information flow in which earlier levels of processing can be altered based on the tentative or preliminary response in later levels (gray arrows). Yet others combine two or more stages into a unified step, such as simultaneous segmentation and feature extraction. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

State of Nature/Prior

- Prior probabilities reflect domain expert's knowledge of *how likely it is that each type of fish will appear*, before we actually see it.
 - State of nature is a random variable: $P(\omega_1)$, $P(\omega_2)$
 - Uniform priors: The catch of salmon and sea bass is equiprobable ($P(\omega_1) = P(\omega_2)$)
 - $P(\omega_1) + P(\omega_2) = 1$ (exclusivity and exhaustivity)

Problem Analysis

- Extract features from the images
 - Length
 - Lightness
 - Width
 - Number and shape of fins
 - Position of the mouth, etc...
- This is the set of all suggested features to explore for use in our classifier

Representation: Fish Length as Feature

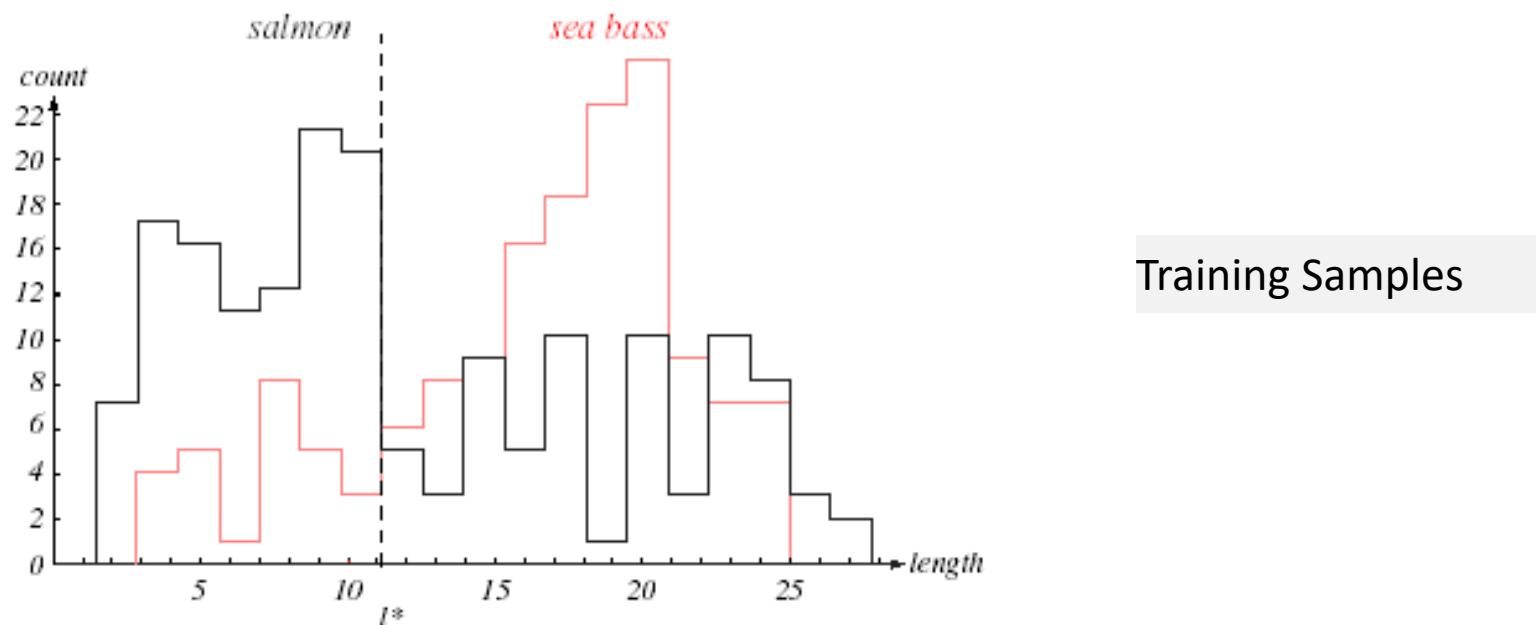


FIGURE 1.2. Histograms for the length feature for the two categories. No single threshold value of the length will serve to unambiguously discriminate between the two categories; using length alone, we will have some errors. The value marked l^* will lead to the smallest number of errors, on average. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Class-conditional Probabilities

- Use of the class-conditional information
- $P(x|\omega_1)$ and $P(x|\omega_2)$ describe the difference in feature (length or lightness) between the populations of sea-bass and salmon

Class-conditional PDF

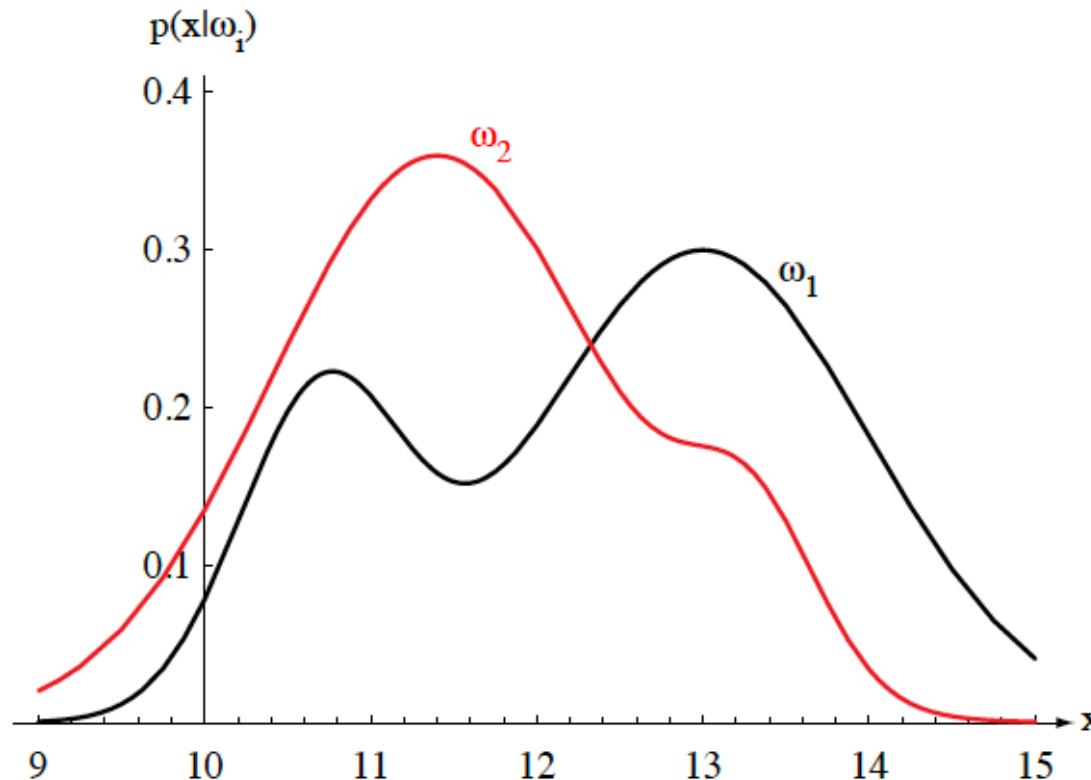


Figure 2.1: Hypothetical class-conditional probability density functions show the probability density of measuring a particular feature value x given the pattern is in category ω_i . If x represents the length of a fish, the two curves might describe the difference in length of populations of two types of fish. Density functions are normalized, and thus the area under each curve is 1.0.

The Bayes Classifier

- Use Bayes Rule!

$$P(Y|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y)P(Y)}{P(X_1, \dots, X_n)}$$

↑
Posterior Likelihood Prior
Normalization Constant

Bayes' Classification

- Posterior, likelihood, prior, evidence

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

- Evidence: In case of two categories

$$p(x) = \sum_{j=1}^2 p(x|\omega_j)P(\omega_j)$$

$$posterior = \frac{likelihood \times prior}{evidence}$$

Posterior Probabilities

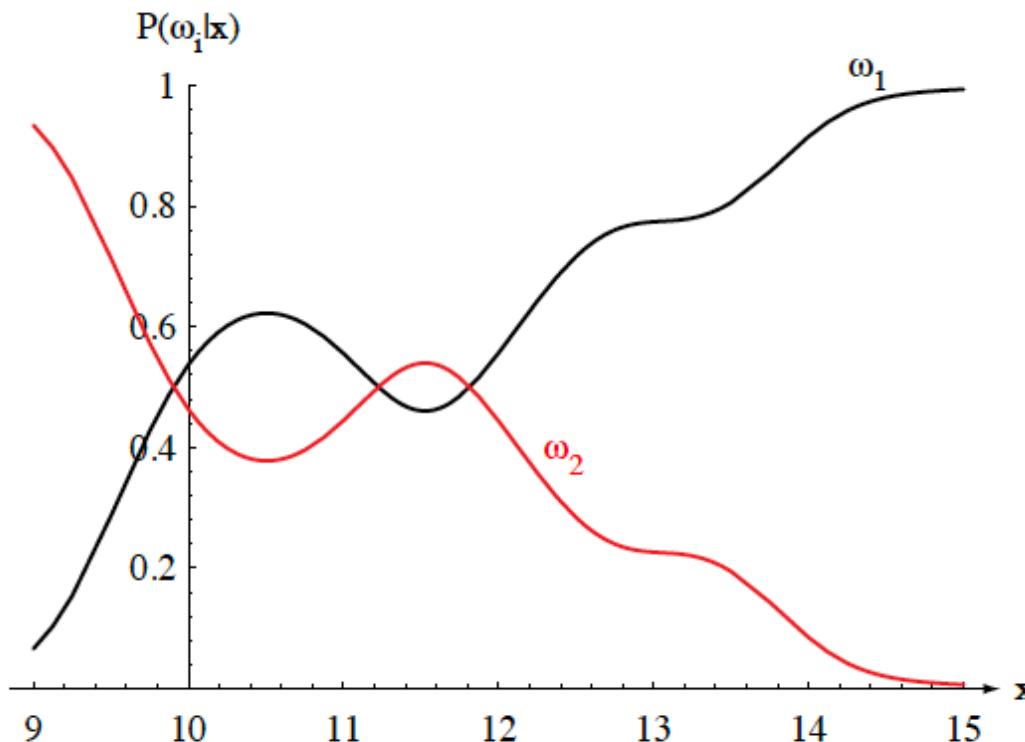


Figure 2.2: Posterior probabilities for the particular priors $P(\omega_1) = 2/3$ and $P(\omega_2) = 1/3$ for the class-conditional probability densities shown in Fig. 2.1. Thus in this case, given that a pattern is measured to have feature value $x = 14$, the probability it is in category ω_2 is roughly 0.08, and that it is in ω_1 is 0.92. At every x , the posteriors sum to 1.0.

Bayes' Decision

- Decision given the posterior probabilities

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

- Therefore, whenever we observe a particular x , the probability of error is :

$$P(error|x) = \min [P(\omega_1|x), P(\omega_2|x)]$$

Questions

Review

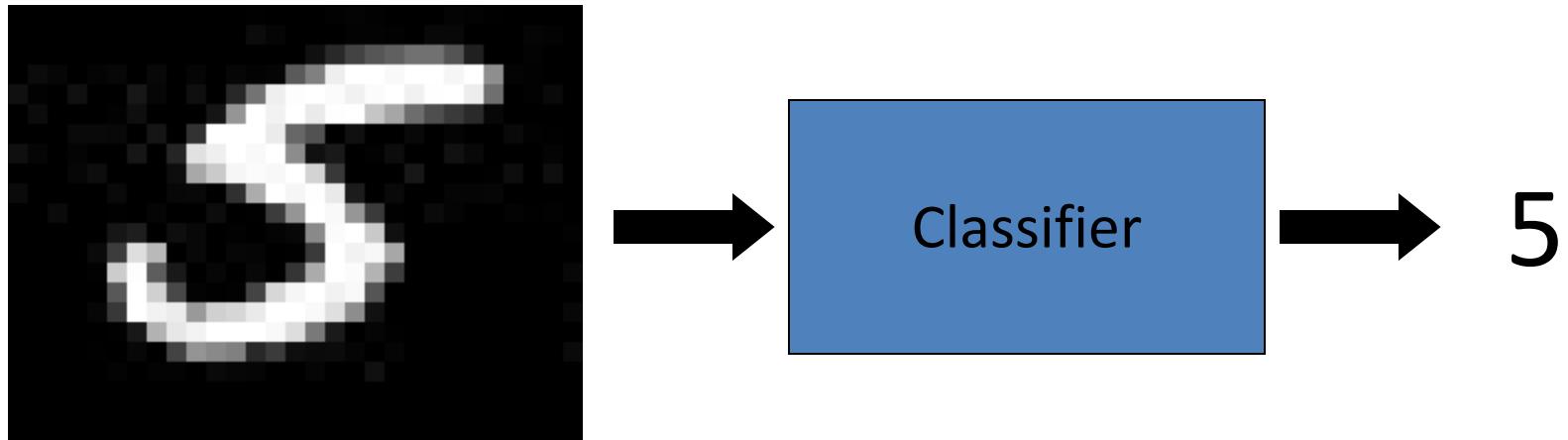
- Classification based on a single feature
- Two class classification
- Sample is assigned to one of the two classes
- The cost of making a false accept or a false reject is same

Bayesian Decision Theory

- Generalization of the preceding ideas
 - Use of more than one feature
 - Use more than two states of nature
 - Allowing actions other than decide on the state of nature
 - Allowing actions other than classification primarily allows the possibility of rejection
 - Refusing to make a decision in close or bad cases!
 - Introduce a loss function which is more general than the probability of error
 - The loss function states how costly each action taken is

Another Application

- **Digit Recognition**



- $X_1, \dots, X_n \in \{0,1\}$ (Black vs. White pixels)
- $Y \in \{5,6\}$ (predict whether a digit is a 5 or a 6)

The Bayes Classifier

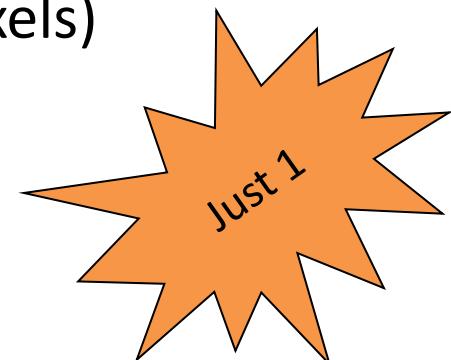
- Let's expand this for our digit recognition task:

$$P(Y = 5|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y = 5)P(Y = 5)}{P(X_1, \dots, X_n|Y = 5)P(Y = 5) + P(X_1, \dots, X_n|Y = 6)P(Y = 6)}$$
$$P(Y = 6|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y = 6)P(Y = 6)}{P(X_1, \dots, X_n|Y = 5)P(Y = 5) + P(X_1, \dots, X_n|Y = 6)P(Y = 6)}$$

- To classify, we will simply compute these two probabilities and predict based on which one is greater

Model Parameters

- For the Bayes classifier, we need to “learn” two functions, the likelihood and the prior
- How many parameters are required to specify the prior for our digit recognition example?
- How many parameters are required to specify the likelihood?
 - (Supposing that each image is 30x30 pixels)
- # of parameters for modeling $P(X_1, \dots, X_n | Y)$:
 - $2(2^n - 1)$



Model Parameters

- The problem with explicitly modeling $P(X_1, \dots, X_n | Y)$ is that there are usually way too many parameters:
 - We'll run out of space
 - We'll run out of time
 - And we'll need tons of training data (which is usually not available)

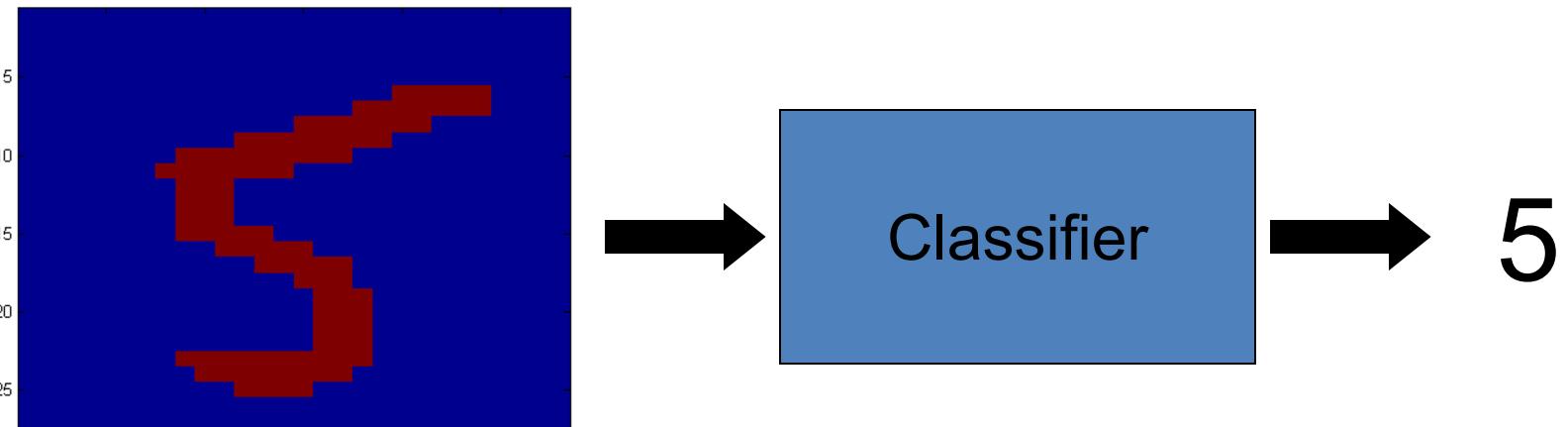
The Naïve Bayes Model

- The *Naïve Bayes Assumption*: Assume that all features are independent **given the class label Y**
- Equationally speaking:

$$P(X_1, \dots, X_n | Y) = \prod_{i=1}^n P(X_i | Y)$$

Another Application

- **Digit Recognition**



- $X_1, \dots, X_n \in \{0,1\}$ (Black vs. White pixels)
- $Y \in \{5,6\}$ (predict whether a digit is a 5 or a 6)

The Bayes Classifier

- A good strategy is to predict:

$$\arg \max_Y P(Y|X_1, \dots, X_n)$$

- (for example: what is the probability that the image represents a 5 given its pixels?)

- So ... How do we compute that?

The Bayes Classifier

- Use Bayes Rule!

$$P(Y|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y)P(Y)}{P(X_1, \dots, X_n)}$$

Likelihood

Prior

Normalization Constant

- Why did this help? Well, we think that we might be able to specify how features are “generated” by the class label

The Bayes Classifier

- Let's expand this for our digit recognition task:

$$P(Y = 5|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y = 5)P(Y = 5)}{P(X_1, \dots, X_n|Y = 5)P(Y = 5) + P(X_1, \dots, X_n|Y = 6)P(Y = 6)}$$
$$P(Y = 6|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y = 6)P(Y = 6)}{P(X_1, \dots, X_n|Y = 5)P(Y = 5) + P(X_1, \dots, X_n|Y = 6)P(Y = 6)}$$

- To classify, we'll simply compute these two probabilities and predict based on which one is greater

Model Parameters

- For the Bayes classifier, we need to “learn” two functions, the likelihood and the prior
- How many parameters are required to specify the prior for our digit recognition example?

Model Parameters

- How many parameters are required to specify the likelihood?
 - (Supposing that each image is 30x30 pixels)

?

Model Parameters

- The problem with explicitly modeling $P(X_1, \dots, X_n | Y)$ is that there are usually way too many parameters:
 - We'll run out of space
 - We'll run out of time
 - And we'll need tons of training data (which is usually not available)

The Naïve Bayes Model

- The *Naïve Bayes Assumption*: Assume that all features are independent **given the class label Y**
- Equationally speaking:

$$P(X_1, \dots, X_n | Y) = \prod_{i=1}^n P(X_i | Y)$$

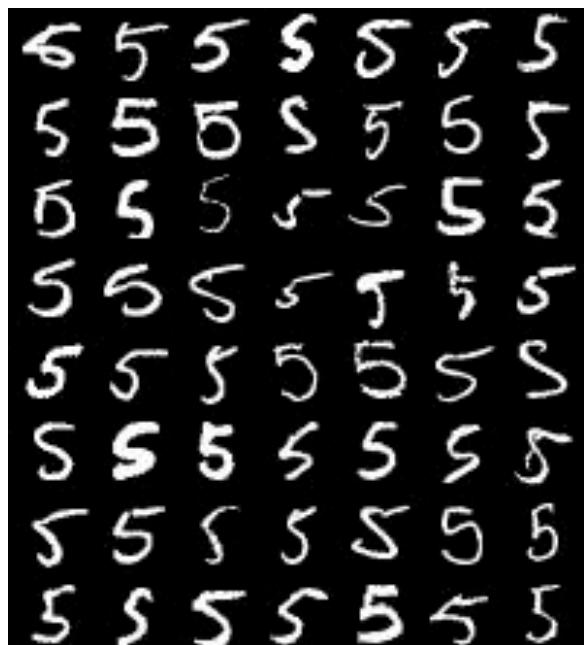
- (We will discuss the validity of this assumption later)

Why is this useful?

- # of parameters for modeling $P(X_1, \dots, X_n | Y)$:
 - Given each x_i is a binary attribute and y is boolean
 - $2(2^n - 1)$
- # of parameters for modeling $P(X_1 | Y), \dots, P(X_n | Y)$
 - $2n$

Naïve Bayes Training

- Now that we've decided to use a Naïve Bayes classifier, we need to train it with some data:



MNIST Training Data

Naïve Bayes Training

- Training in Naïve Bayes is **easy**:
 - Estimate $P(Y=v)$ as the fraction of records with $Y=v$

$$P(Y = v) = \frac{\text{Count}(Y = v)}{\# \text{ records}}$$

- Estimate $P(X_i=u | Y=v)$ as the fraction of records with $Y=v$ for which $X_i=u$
- $$P(X_i = u | Y = v) = \frac{\text{Count}(X_i = u \wedge Y = v)}{\text{Count}(Y = v)}$$
- (This corresponds to Maximum Likelihood estimation of model parameters)

Naïve Bayes Training

- In practice, some of these counts can be zero
- Fix this by adding “virtual” counts:

$$P(X_i = u|Y = v) = \frac{\text{Count}(X_i = u \wedge Y = v) + 1}{\text{Count}(Y = v) + 2}$$

- (This is like putting a prior on parameters and doing MAP estimation instead of MLE)
- This is called *Smoothing*

Bayesian Decision Theory

- Generalization of the preceding ideas
 - Use of more than one feature
 - Use more than two states of nature
 - Allowing actions other than decide on the state of nature
 - Allowing actions other than classification primarily allows the possibility of rejection
 - Refusing to make a decision in close or bad cases!
 - Introduce a loss function which is more general than the probability of error
 - The loss function states how costly each action taken is

Bayesian Decision Theory – Continuous Features...

- Let $\{\omega_1, \omega_2, \dots, \omega_c\}$ be the set of c states of nature (or “categories”)
- Let $\{\alpha_1, \alpha_2, \dots, \alpha_a\}$ be the set of possible actions
- Let $\lambda(\alpha_i | \omega_j)$ be the loss incurred for taking action α_i when the true state of nature is ω_j

Two-category Classification

- α_1 : deciding ω_1
- α_2 : deciding ω_2
- $\lambda_{ij} = \lambda(\alpha_i | \omega_j)$
- Loss incurred for deciding α_i when the true state of nature is ω_j

Two-category Classification

- α_1 : deciding ω_1
- α_2 : deciding ω_2
- $\lambda_{ij} = \lambda(\alpha_i | \omega_j)$
- Loss incurred for deciding α_i when the true state of nature is ω_j
- Conditional risk:

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

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

Two-category Classification

- Our rule is the following:

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

- Action α_1 : “decide ω_1 ” is taken
- This results in the equivalent rule :
- Decide ω_1 if:

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

- and decide ω_2 otherwise

Bayesian Decision Theory – Continuous Features...

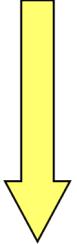
- Overall risk

$$R = \text{Sum of all } \underbrace{R(\alpha_i | x)}_{\text{Conditional risk}} \text{ for } i = 1, \dots, a$$

- Minimizing R \longleftrightarrow Minimizing $R(\alpha_i | x)$ for $i = 1, \dots, a$

$$\bullet \quad R(\alpha_i | x) = \sum_{j=1}^{j=c} \lambda(\alpha_i | \omega_j) P(\omega_j | x) \quad \text{for } i = 1, \dots, a$$

Bayesian Decision Theory – Continuous Features...

- Select the action α_i for which $R(\alpha_i | x)$ is minimum
- 
- R is minimum and is called the Bayes risk = best performance that can be achieved!

Likelihood Ratio

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

- The preceding rule is equivalent to the following rule:
- If $\frac{p(\mathbf{x}|\omega_1)}{p(\mathbf{x}|\omega_2)} > \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \frac{P(\omega_2)}{P(\omega_1)}$
- Then take action α_1 (decide ω_1)
- Otherwise take action α_2 (decide ω_2)

Likelihood Ratio...

- Optimal decision property

“If the likelihood ratio exceeds a threshold value independent of the input pattern x , we can take optimal actions”

Example: Checking on a course

- A student needs to make a decision which courses to take, based only on first lecture's impression
- From student's previous experience:

Quality of the course	good	fair	bad
Probability (prior)	0.2	0.4	0.4

- These are prior probabilities.

Example: Checking on a course

- The student also knows the class-conditionals:

$\Pr(x/w_j)$	good	fair	bad
Interesting lecture	0.8	0.5	0.1
Boring lecture	0.2	0.5	0.9

- The loss function is given by the matrix

$\lambda(\alpha_i \omega_j)$	good course	fair course	bad course
Taking the course	0	5	10
Not taking the course	20	5	0

Example: Checking on a course

- We can take decisions according to two rules:
 - Posterior probability
 - Risk minimization

Example: Checking on a course

- Given that the first class was interesting, the student wants to make an optimal decision.
 - $\Pr(\text{bad} \mid \text{interesting})$
 - $\Pr(\text{fair} \mid \text{interesting})$
 - $\Pr(\text{good} \mid \text{interesting})$
- $\Pr(\text{interesting} \mid \text{bad})$
- $\Pr(\text{interesting} \mid \text{fair})$
- $\Pr(\text{interesting} \mid \text{good})$
- $\Pr(\text{interesting})$

Example: Checking on a course

- The probability to get the “interesting lecture”(x= interesting):
- $\Pr(\text{interesting}) = \Pr(\text{interesting} | \text{good course}) * \Pr(\text{good course})$
+ $\Pr(\text{interesting} | \text{fair course}) * \Pr(\text{fair course})$
+ $\Pr(\text{interesting} | \text{bad course}) * \Pr(\text{bad course})$
 $= 0.8 * 0.2 + 0.5 * 0.4 + 0.1 * 0.4 = 0.4$
- Consequently, $\Pr(\text{boring}) = 1 - 0.4 = 0.6$
- Suppose the lecture was interesting. Then we want to compute the **posterior** probabilities of each one of the 3 possible “states of nature”.

Example: Checking on a course

$$\Pr(\text{good course}|\text{interesting lecture})$$

$$= \frac{\Pr(\text{interesting}|\text{good})\Pr(\text{good})}{\Pr(\text{interesting})} = \frac{0.8 * 0.2}{0.4} = 0.4$$

$$\Pr(\text{fair}|\text{interesting})$$

$$= \frac{\Pr(\text{interesting}|\text{fair})\Pr(\text{fair})}{\Pr(\text{interesting})} = \frac{0.5 * 0.4}{0.4} = 0.5$$

- We can get $\Pr(\text{bad}|\text{interesting}) = 0.1$ either by the same method, or by noting that it complements the above two to 1.
- Now, we have all we need to make an intelligent decision about an optimal action

Example: Checking on a course

- The student needs to minimize the conditional risk

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

$R(\text{taking} | \text{interesting}) =$

$\Pr(\text{good} | \text{interesting}) \lambda(\text{taking good course}) +$

$\Pr(\text{fair} | \text{interesting}) \lambda(\text{taking fair course}) +$

$\Pr(\text{bad} | \text{interesting}) \lambda(\text{taking bad course})$

$$= 0.4 * 0 + 0.5 * 5 + 0.1 * 10 = 3.5$$

Example: Checking on a course

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

R(not taking | interesting) =

$$\begin{aligned} & \Pr(\text{good} | \text{interesting}) \lambda(\text{not taking good course}) + \\ & \Pr(\text{fair} | \text{interesting}) \lambda(\text{not taking fair course}) + \\ & \Pr(\text{bad} | \text{interesting}) \lambda(\text{not taking bad course}) \end{aligned}$$

$$= 0.4 * 20 + 0.5 * 5 + 0.1 * 0 = 10.5$$

Constructing an optimal decision function

- So, if the first lecture was interesting, the student will minimize the conditional risk by taking the course.
- In order to construct the full decision function, we need to define the risk minimization action for the case of boring lecture, as well.

Exercise

- Select the optimal decision where: $W = \{\omega_1, \omega_2\}$
 - $P(x | \omega_1) \rightarrow N(2, 0.5)$ (Normal distribution)
 - $P(x | \omega_2) \rightarrow N(1.5, 0.2)$
 - $P(\omega_1) = 2/3$
 - $P(\omega_2) = 1/3$
- $$\lambda = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

Bayesian Decision Theory

- Generalization of the preceding ideas
 - Use of more than one feature
 - Use more than two states of nature
 - Allowing actions other than decide on the state of nature
 - Allowing actions other than classification primarily allows the possibility of rejection
 - Refusing to make a decision in close or bad cases!
 - Introduce a loss function which is more general than the probability of error
 - The loss function states how costly each action taken is

The Normal Density

- Univariate density
 - Continuous density
 - A lot of processes are asymptotically Gaussian
 - Handwritten characters, speech sounds are ideal or prototype corrupted by random process (central limit theorem)

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

Univariate density

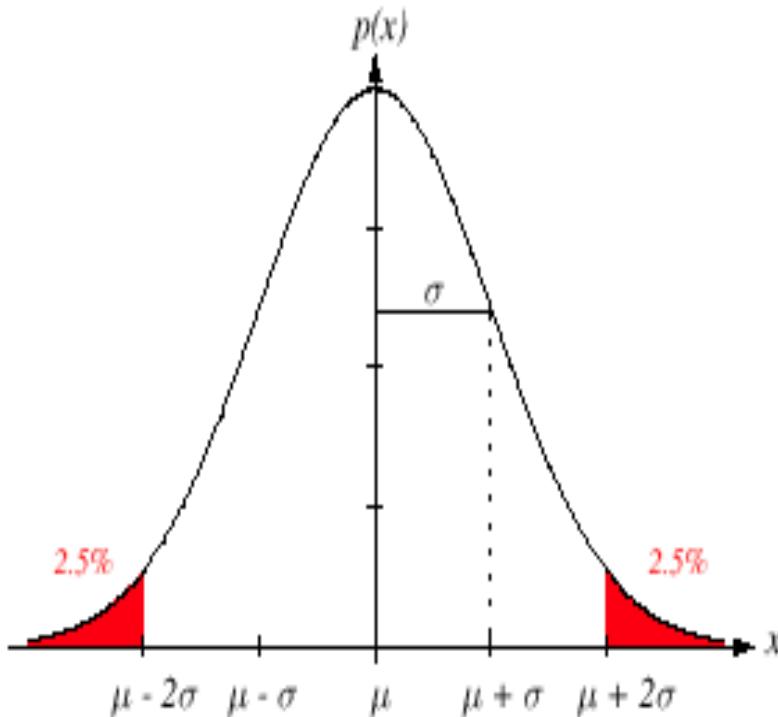


FIGURE 2.7. A univariate normal distribution has roughly 95% of its area in the range $|x - \mu| \leq 2\sigma$, as shown. The peak of the distribution has value $p(\mu) = 1/\sqrt{2\pi}\sigma$. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John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

The Normal Density

- Multivariate density: Multivariate normal density in d dimensions is:

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

$$x = (x_1, x_2, \dots, x_d)^t$$

$$cov_{x,y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N - 1}$$

$\mu = (\mu_1, \mu_2, \dots, \mu_d)^t$ mean vector

$\Sigma = d^*d$ covariance matrix

$|\Sigma|$ and Σ^{-1} are determinant and inverse respectively

Multivariate density

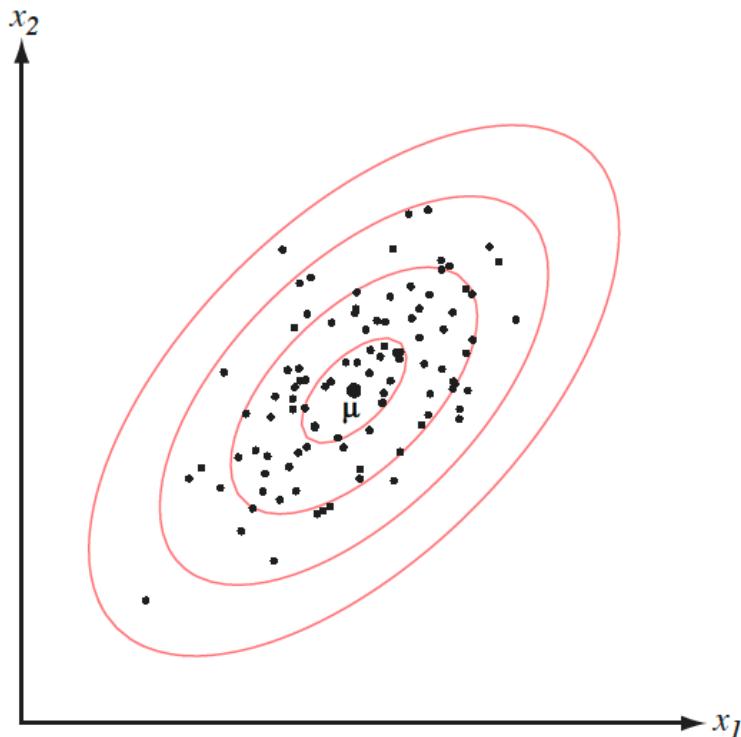


Figure 2.9: Samples drawn from a two-dimensional Gaussian lie in a cloud centered on the mean μ . The red ellipses show lines of equal probability density of the Gaussian.

Discriminant Functions for the Normal Density

$$P(\omega_j|x) = \frac{p(x|\omega_j)P(\omega_j)}{p(x)} \quad P(error|x) = \min [P(\omega_1|x), P(\omega_2|x)]$$

- Minimum error-rate classification can be achieved by the discriminant function
- $g1(x) = \ln p(x | \omega_1) + \ln P(\omega_1) - \ln p(x)$
- $g2(x) = \ln p(x | \omega_2) + \ln P(\omega_2) - \ln p(x)$

Discriminant Functions for the Normal Density

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

- Minimum error-rate classification can be achieved by the discriminant function
- $g_i(x) = \ln P(x | \omega_i) + \ln P(\omega_i)$
- Case of multivariate normal

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

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

Analyzing Covariance Matrix

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

- Two cases:
 - Features are independent
 - Features are dependent

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

Case $\Sigma_i = \Sigma$ (covariance of all classes are identical but arbitrary!)

Case Σ_i = actual covariance

Discriminant Functions for the Normal Density

- Case $\Sigma_i = \sigma^2 \cdot I$ (I stands for the identity matrix)
 - $\sigma_{ij} = 0$: Features are statistically independent
 - σ_{ii} is same for all the features

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

1/ σ^2 *Constant for all
the classes* *Constant for all
the classes*

Discriminant Functions for the Normal Density

- Case $\Sigma_i = \sigma^2 \cdot I$ (I stands for the identity matrix)
 - $\sigma_{ij} = 0$: Features are statistically independent)
 - σ_{ii} is same for all the features

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

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

Discriminant Functions for the Normal Density...

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

- Disregarding $\mathbf{x}^t \mathbf{x}$, we get a linear discriminant function

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

where :

$$w_i = \frac{\mu_i}{\sigma^2}; \quad w_{i0} = -\frac{1}{2\sigma^2} \mu_i^t \mu_i + \ln P(\omega_i)$$

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

Discriminant Functions for the Normal Density...

- A classifier that uses linear discriminant functions is called “a linear machine”
- The decision surfaces for a linear machine are hyperplanes defined by $g_i(x) = g_j(x)$

$$\mathbf{w}^t(\mathbf{x} - \mathbf{x}_0) = 0 \quad \mathbf{w} = \boldsymbol{\mu}_i - \boldsymbol{\mu}_j$$

$$\mathbf{x}_0 = \frac{1}{2}(\boldsymbol{\mu}_i + \boldsymbol{\mu}_j) - \frac{\sigma^2}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2} \ln \frac{P(\omega_i)}{P(\omega_j)} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)$$

Discriminant Functions for the Normal Density...

- The hyperplane separating \mathcal{R}_i and \mathcal{R}_j

$$x_0 = \frac{1}{2}(\mu_i + \mu_j) - \frac{\sigma^2}{\|\mu_i - \mu_j\|^2} \ln \frac{P(\omega_i)}{P(\omega_j)} (\mu_i - \mu_j)$$

is always orthogonal to the line linking the means!

if $P(\omega_i) = P(\omega_j)$ then $x_0 = \frac{1}{2}(\mu_i + \mu_j)$

$$g_i(x) = -\|x - \mu_i\|^2$$

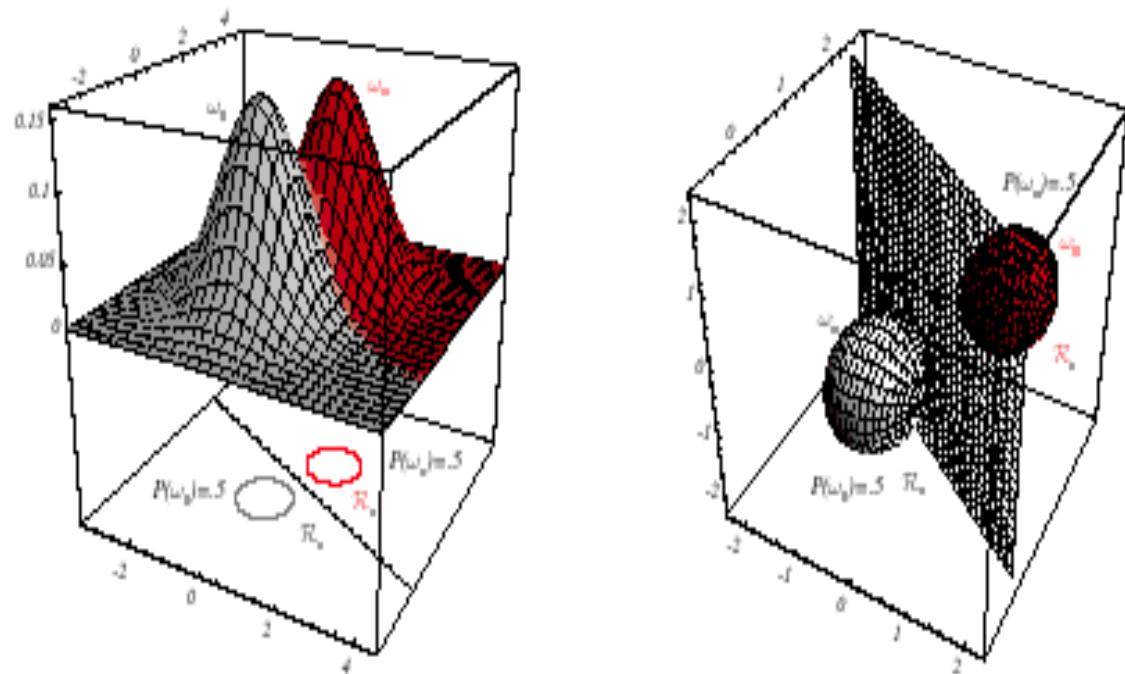
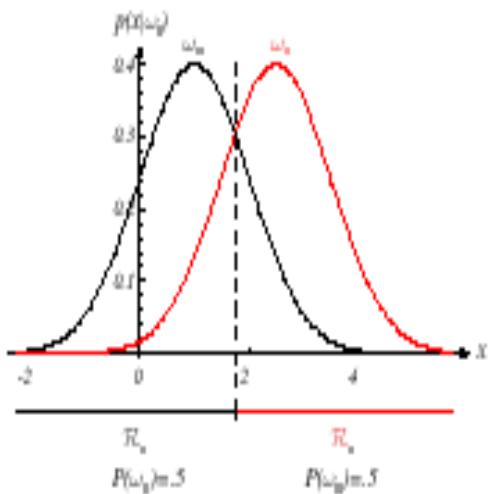


FIGURE 2.10. If the covariance matrices for two distributions are equal and proportional to the identity matrix, then the distributions are spherical in d dimensions, and the boundary is a generalized hyperplane of $d - 1$ dimensions, perpendicular to the line separating the means. In these one-, two-,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

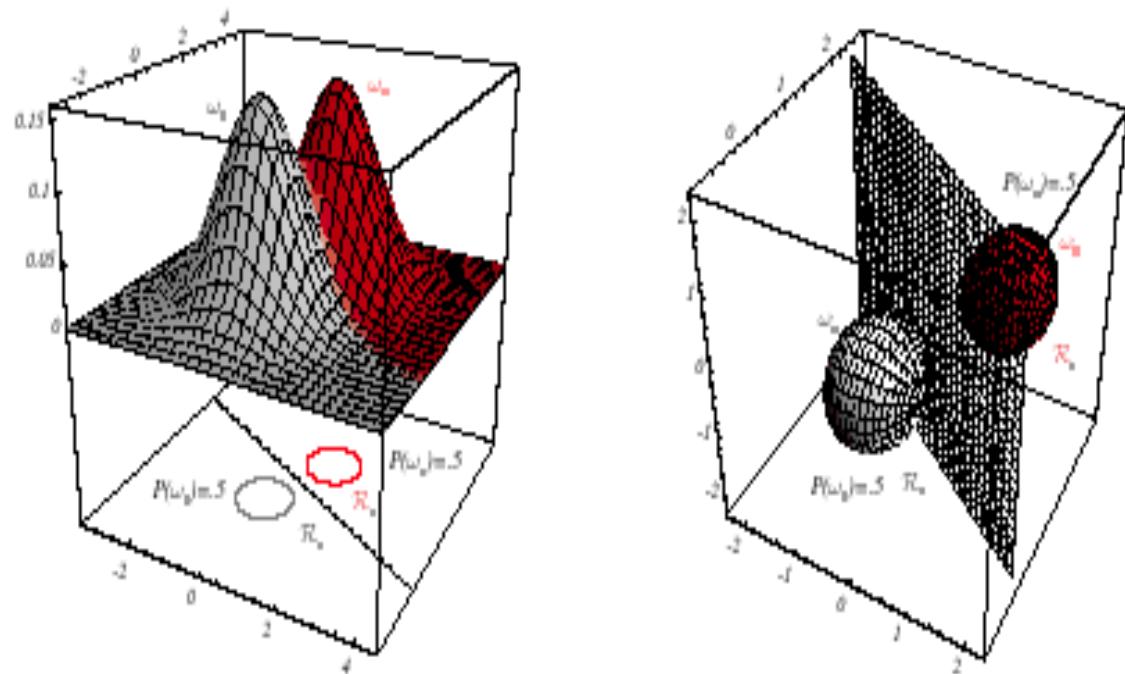
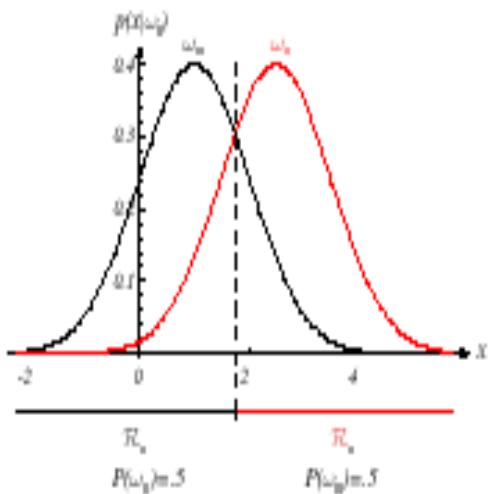
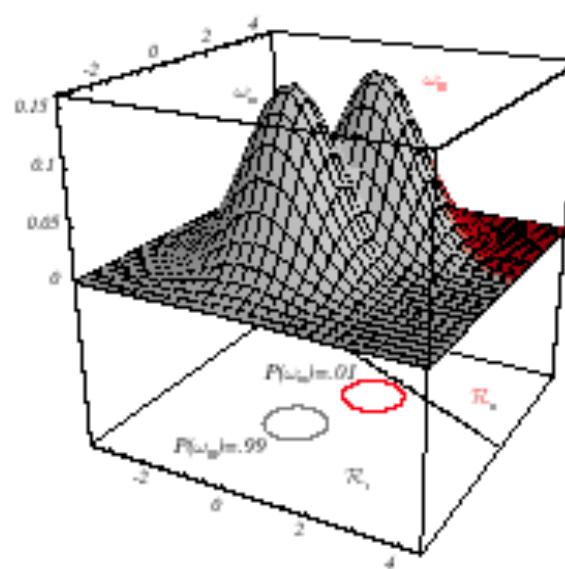
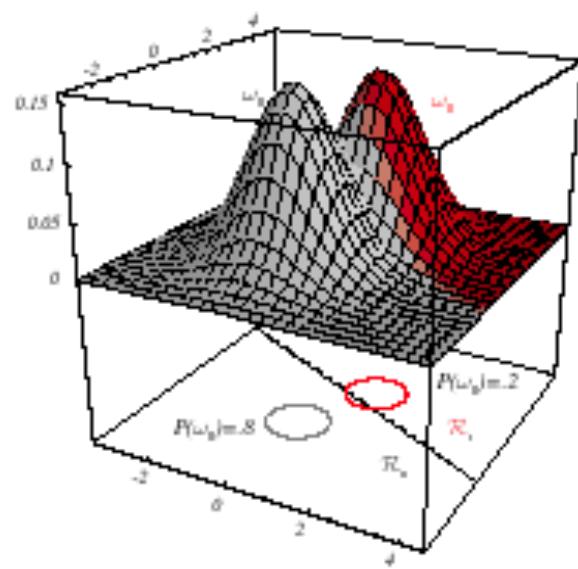
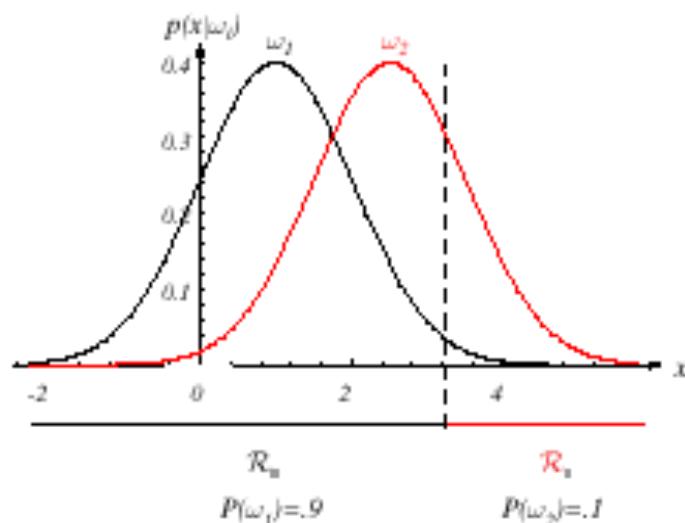
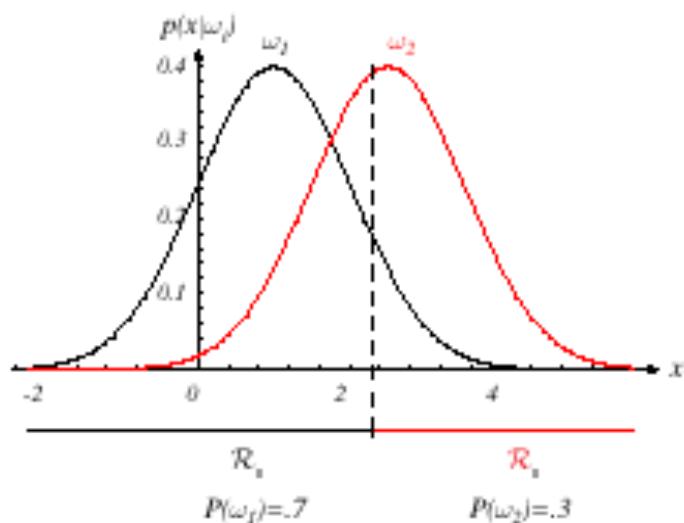


FIGURE 2.10. If the covariance matrices for two distributions are equal and proportional to the identity matrix, then the distributions are spherical in d dimensions, and the boundary is a generalized hyperplane of $d - 1$ dimensions, perpendicular to the line separating the means. In these one-, two-,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



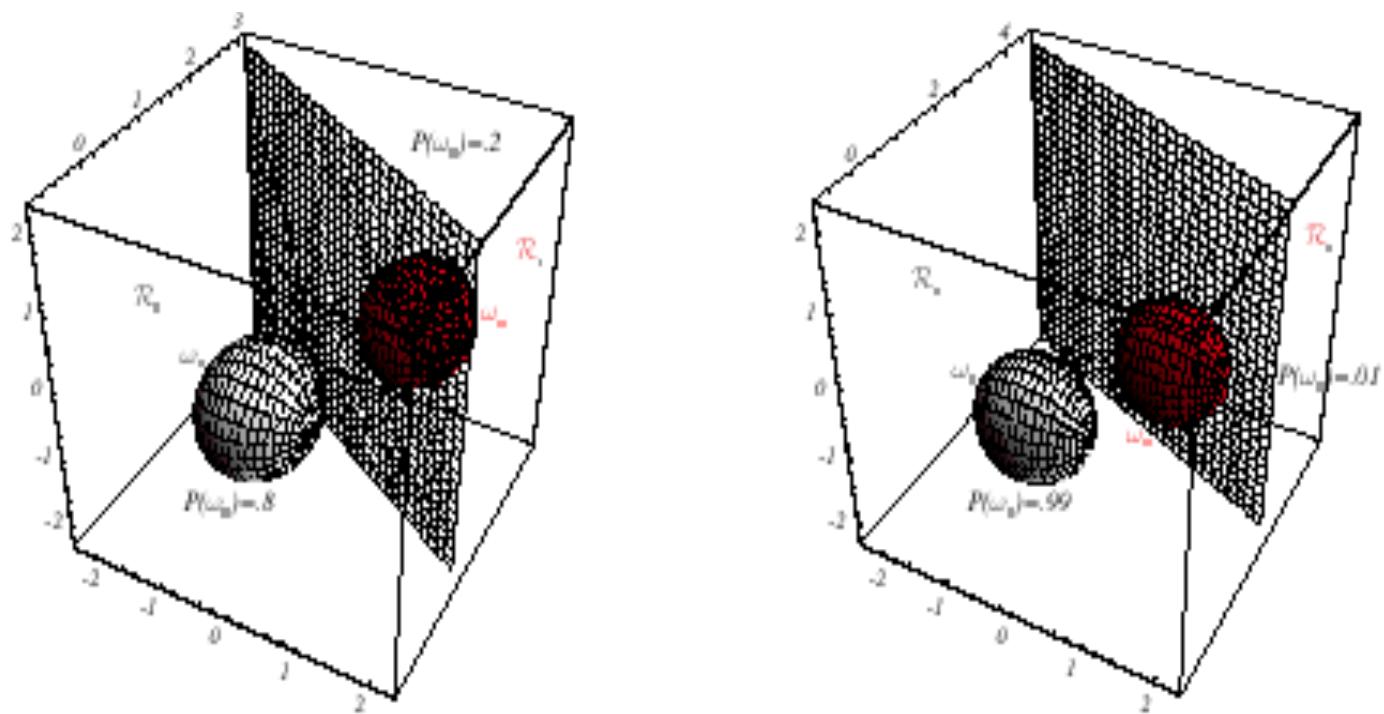


FIGURE 2.11. As the priors are changed, the decision boundary shifts; for sufficiently disparate priors the boundary will not lie between the means of these one-, two- and three-dimensional spherical Gaussian distributions. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*.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

Questions?

Discriminant Functions for the Normal Density...

- Case $\Sigma_i = \Sigma$ (covariance of all classes are identical but arbitrary!)

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

- Expand the term and disregard the quadratic expression

where :

$$g_i(x) = w_i^t x + w_{i0} \quad w_i = \Sigma^{-1} \mu; \quad w_{i0} = -\frac{1}{2} \mu_i^t \Sigma^{-1} \mu_i + \ln P(\omega_i)$$

Discriminant Functions for the Normal Density...

$$x_0 = \frac{1}{2}(\mu_i + \mu_j) - \frac{\ln[P(\omega_i)/P(\omega_j)]}{(\mu_i - \mu_j)^t \Sigma^{-1} (\mu_i - \mu_j)} \cdot (\mu_i - \mu_j)$$

- Comments about this hyperplane:
 - It passes through \mathbf{x}_0
 - It is NOT orthogonal to the line linking the means.
 - What happens when $P(\omega_i) = P(\omega_j)$?
 - If $P(\omega_i) \neq P(\omega_j)$, then \mathbf{x}_0 shifts away from the more likely mean.

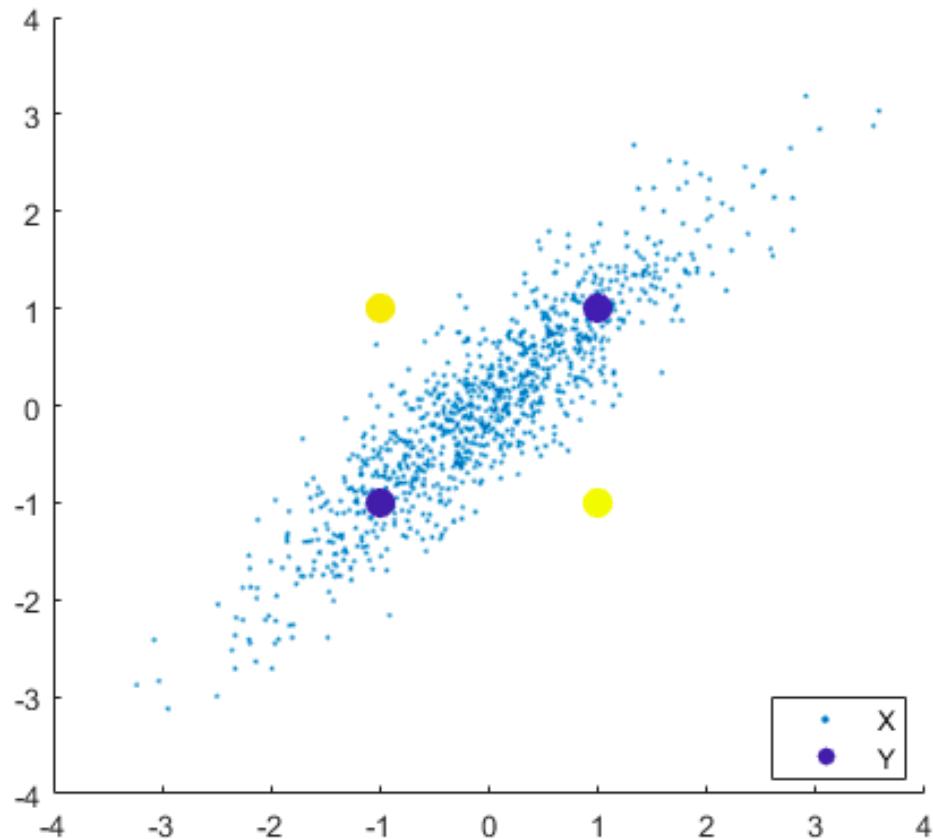
Discriminant Functions for the Normal Density...

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

- When $P(\omega_i)$ is the same for each of the c classes
- Mahalanobis distance classifier

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^t \Sigma_i^{-1} (x - \mu_i)$$

Mahalanobis Distance



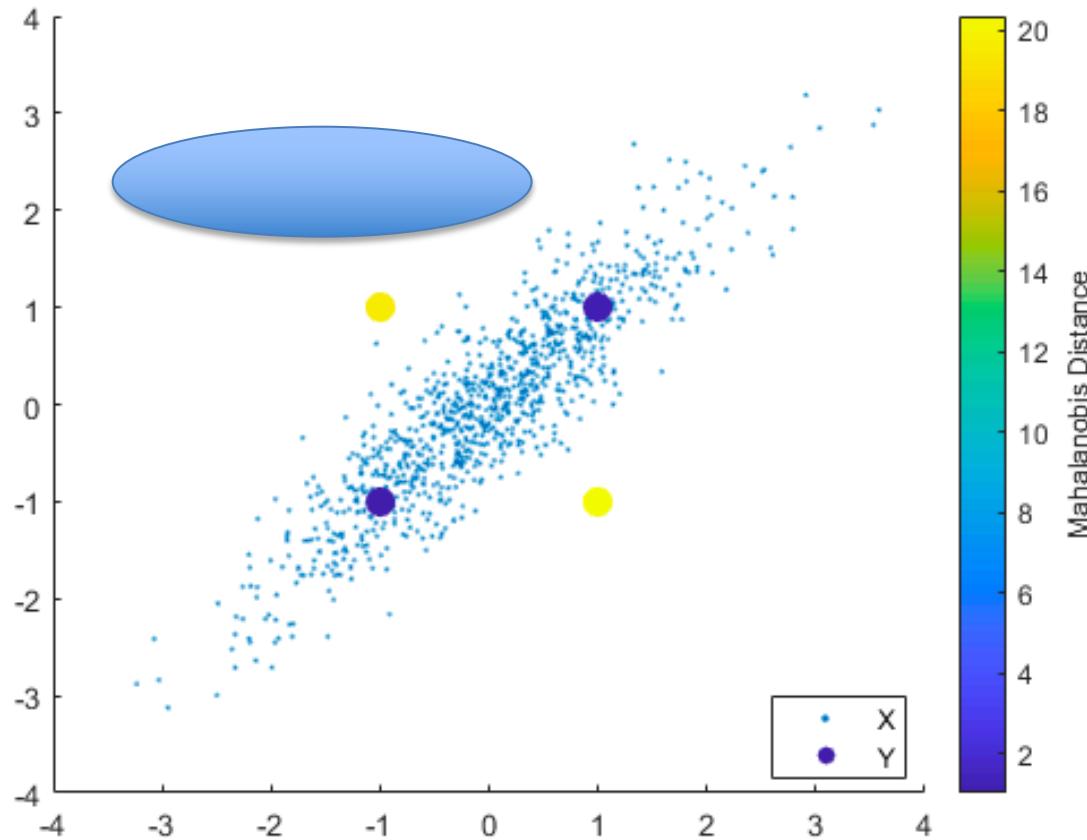
Euclidean Distance:

$$d = \sqrt{\sum_{i=1}^n (X_i - Y_i)^2}$$

Mahalanobis Distance:

$$d(\vec{x}, \vec{y}) = \sqrt{(\vec{x} - \vec{y})^\top S^{-1} (\vec{x} - \vec{y})}$$

Mahalanobis Distance



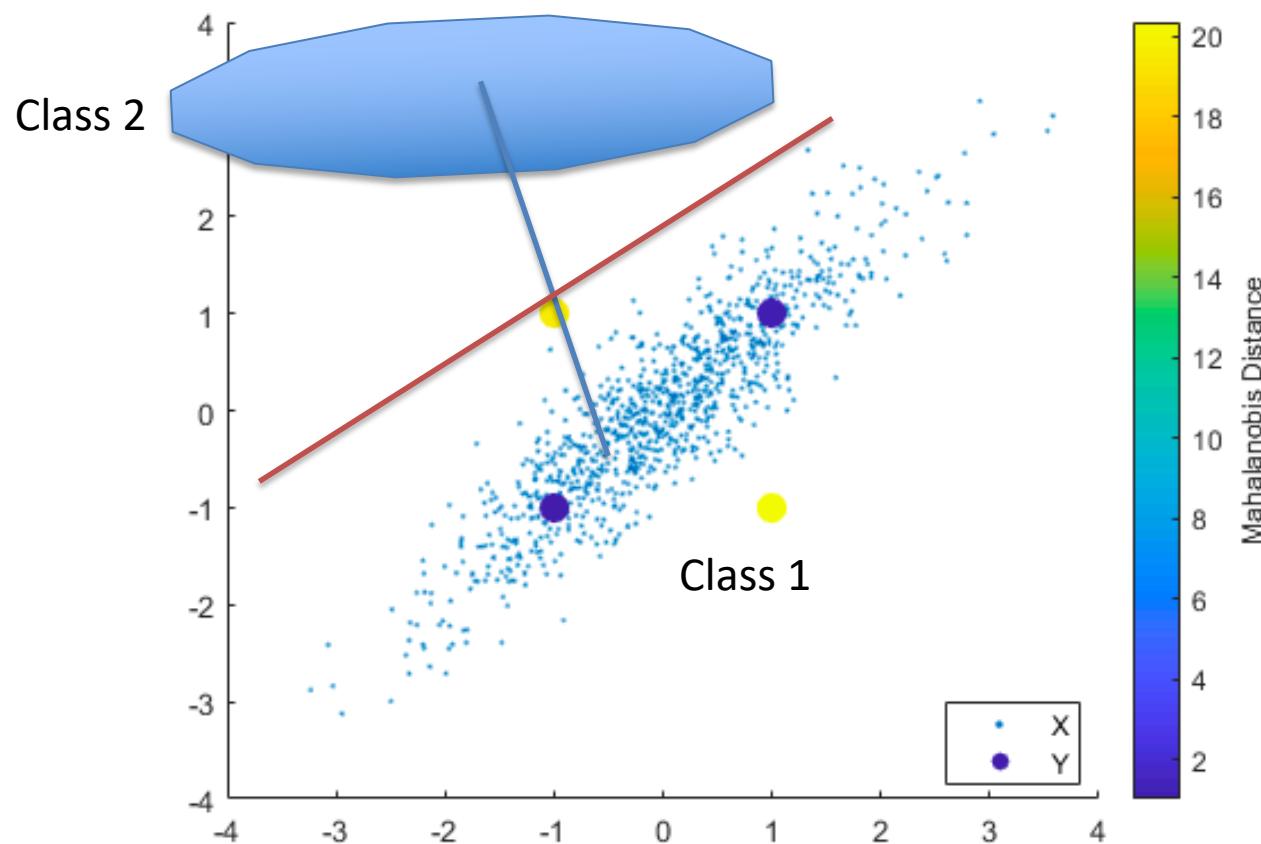
$$d = \sqrt{\sum_{i=1}^n (X_i - Y_i)^2}$$

Euclidean Distance:

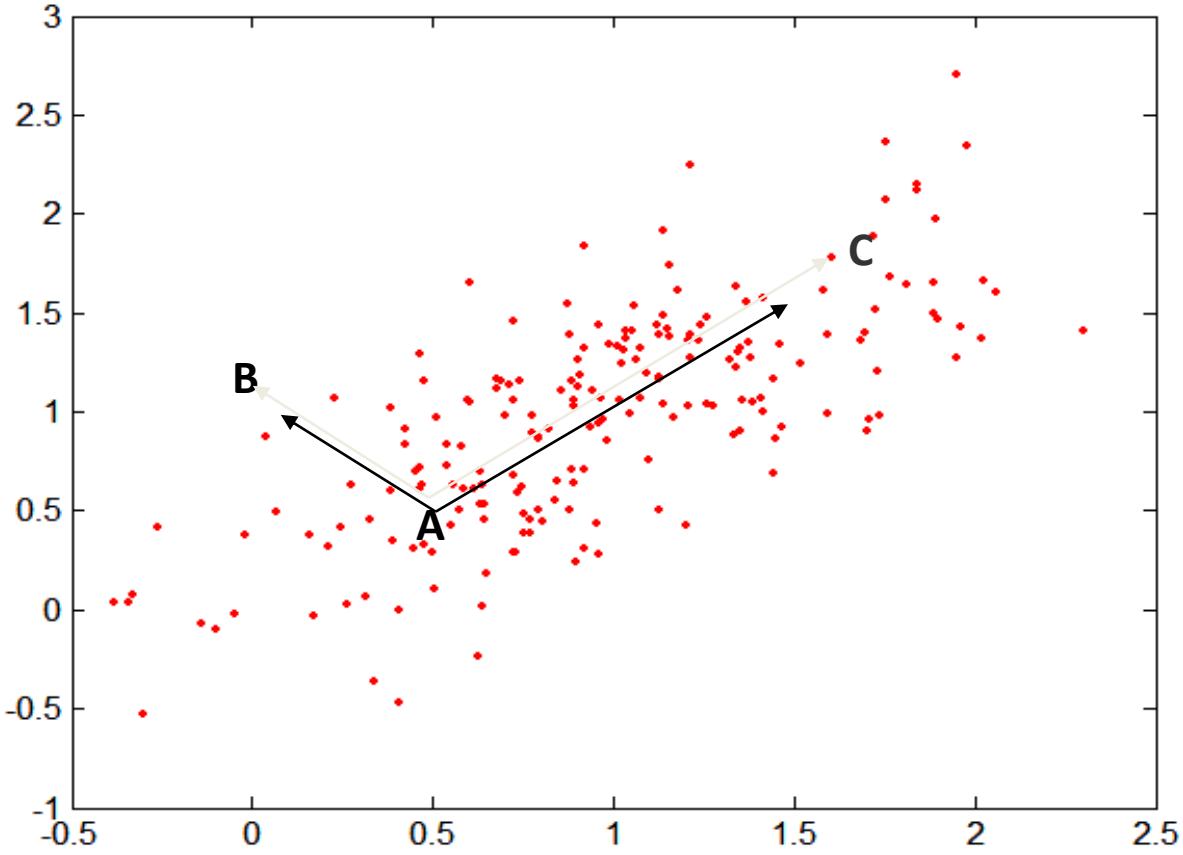
All points are equidistant

$$d(\vec{x}, \vec{y}) = \sqrt{(\vec{x} - \vec{y})^\top S^{-1} (\vec{x} - \vec{y})}$$

Mahalanobis Distance



Mahalanobis Distance



Covariance Matrix:

$$\Sigma = \begin{bmatrix} 0.3 & 0.2 \\ 0.2 & 0.3 \end{bmatrix}$$

A: (0.5, 0.5)

B: (0, 1)

C: (1.5, 1.5)

Compute squared versions

$$d = \sqrt{\sum_{i=1}^n (X_i - Y_i)^2}$$

$$d(\vec{x}, \vec{y}) = \sqrt{(\vec{x} - \vec{y})^\top S^{-1} (\vec{x} - \vec{y})}$$

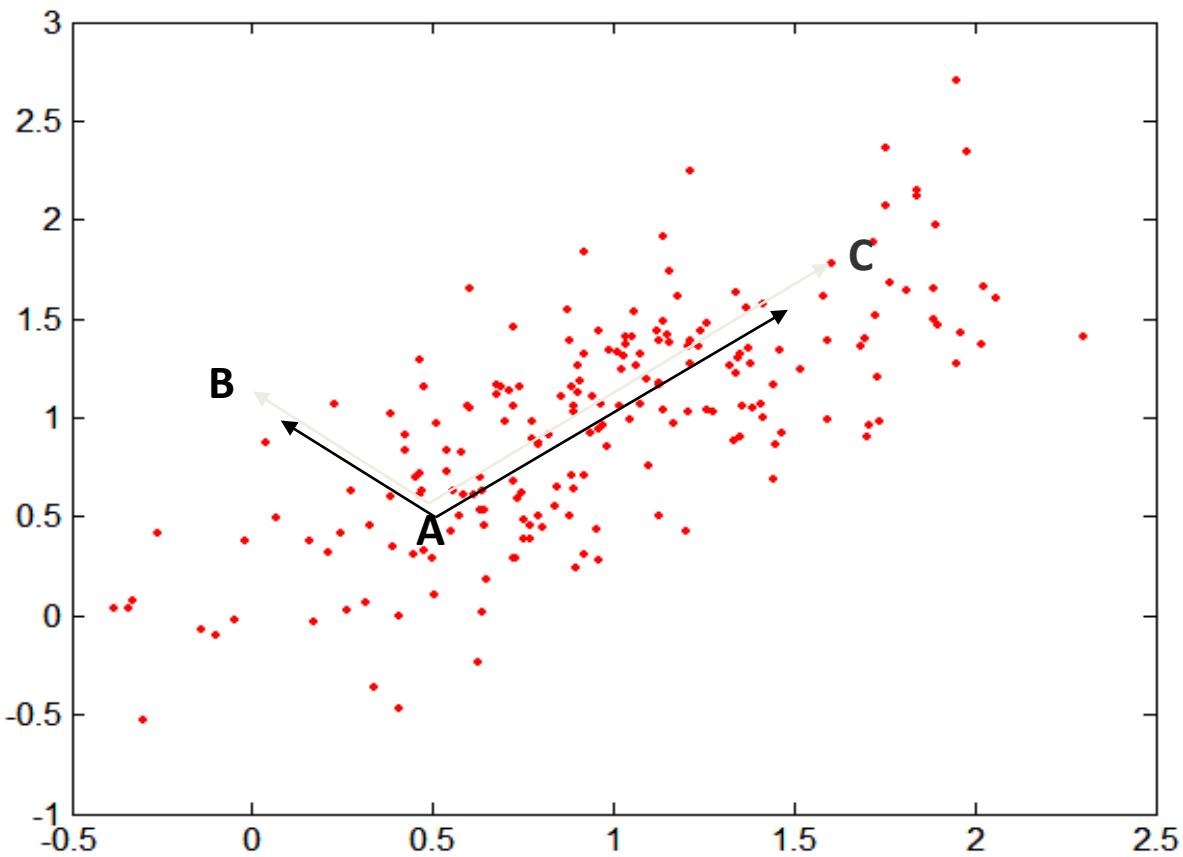
Euclid(A,B)

Euclid(A,C)

Mahal(A,B)

Mahal(A,C)

Mahalanobis Distance



Covariance Matrix:

$$\Sigma = \begin{bmatrix} 0.3 & 0.2 \\ 0.2 & 0.3 \end{bmatrix}$$

A: (0.5, 0.5)

B: (0, 1)

C: (1.5, 1.5)

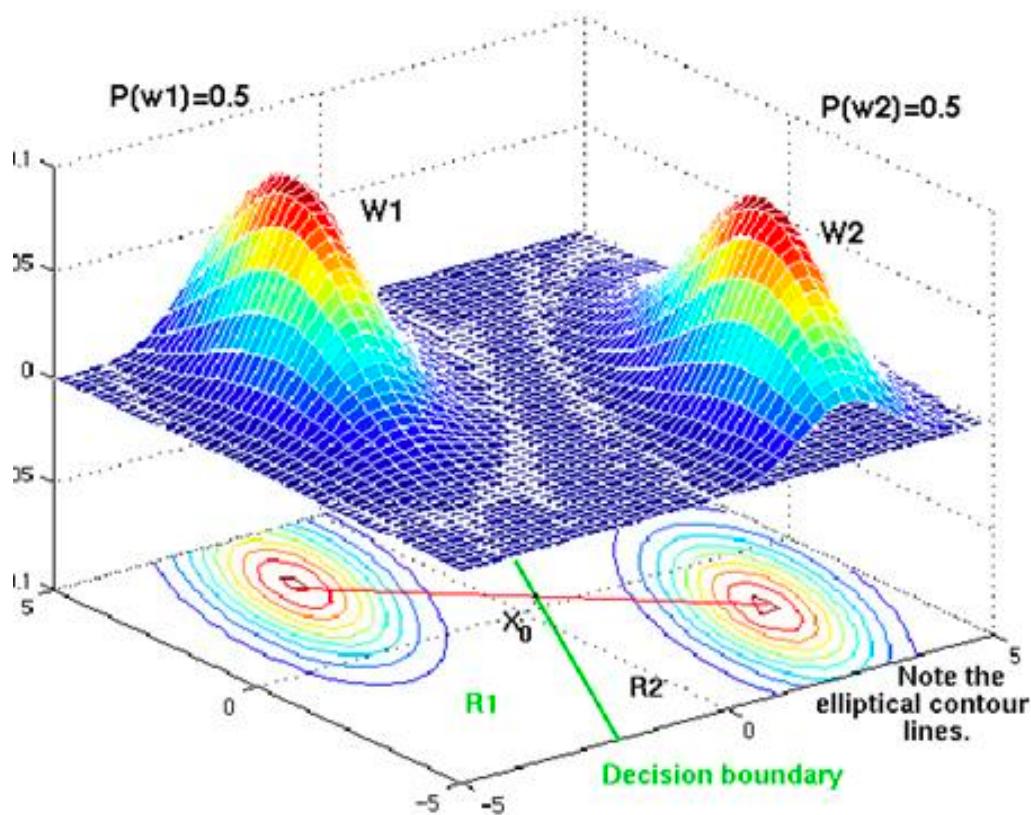
$\text{Euclid}(A,B) = 0.5$

$\text{Euclid}(A,C) = 2$

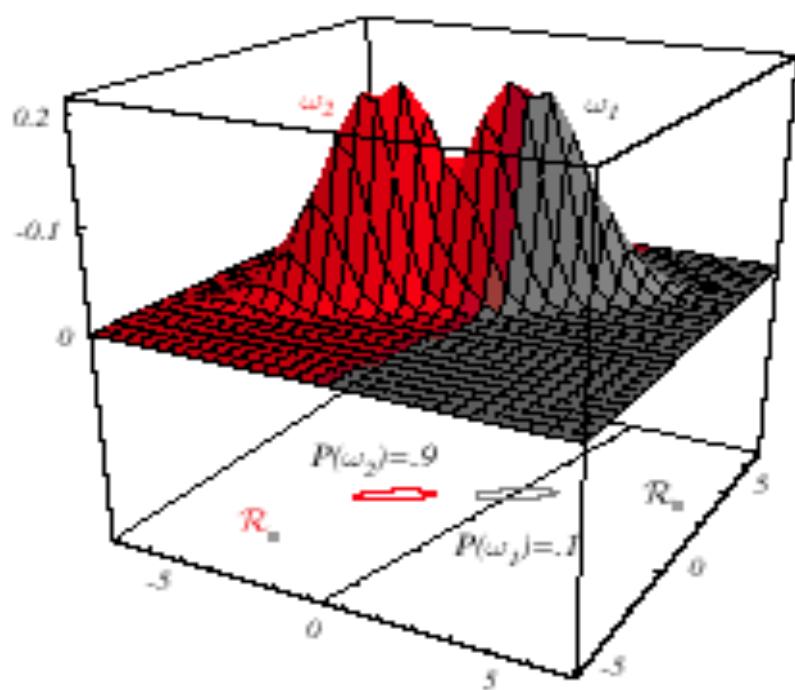
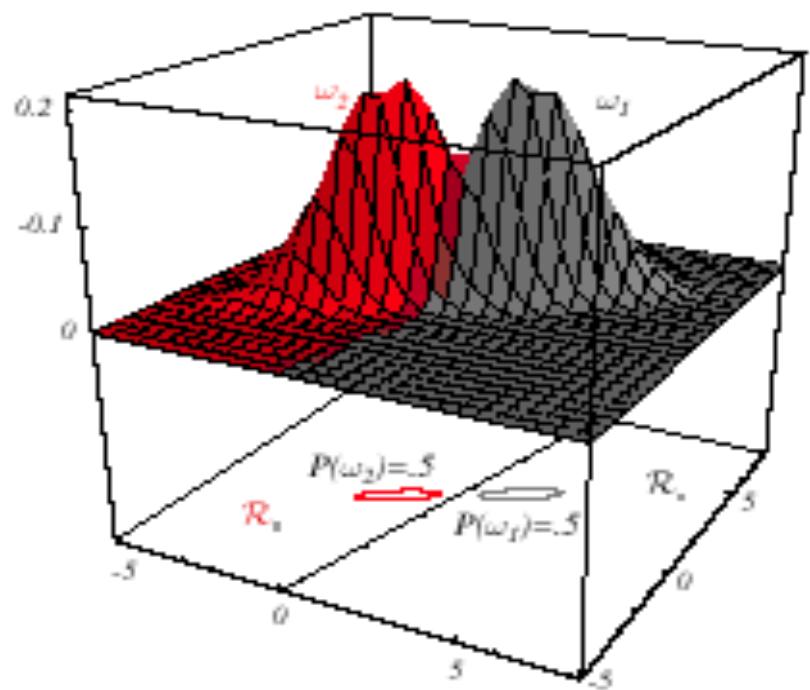
$\text{Mahal}(A,B) = 5$

$\text{Mahal}(A,C) = 4$

Discriminant Functions for the Normal Density...



The contour lines are elliptical in shape because the covariance matrix is not diagonal. However, both densities show the same elliptical shape. The prior probabilities are the same, and so the point x_0 lies halfway between the 2 means. The decision boundary is not orthogonal to the red line. Instead, it is tilted so that its points are of equal distance to the contour lines in w_1 and those in w_2 .



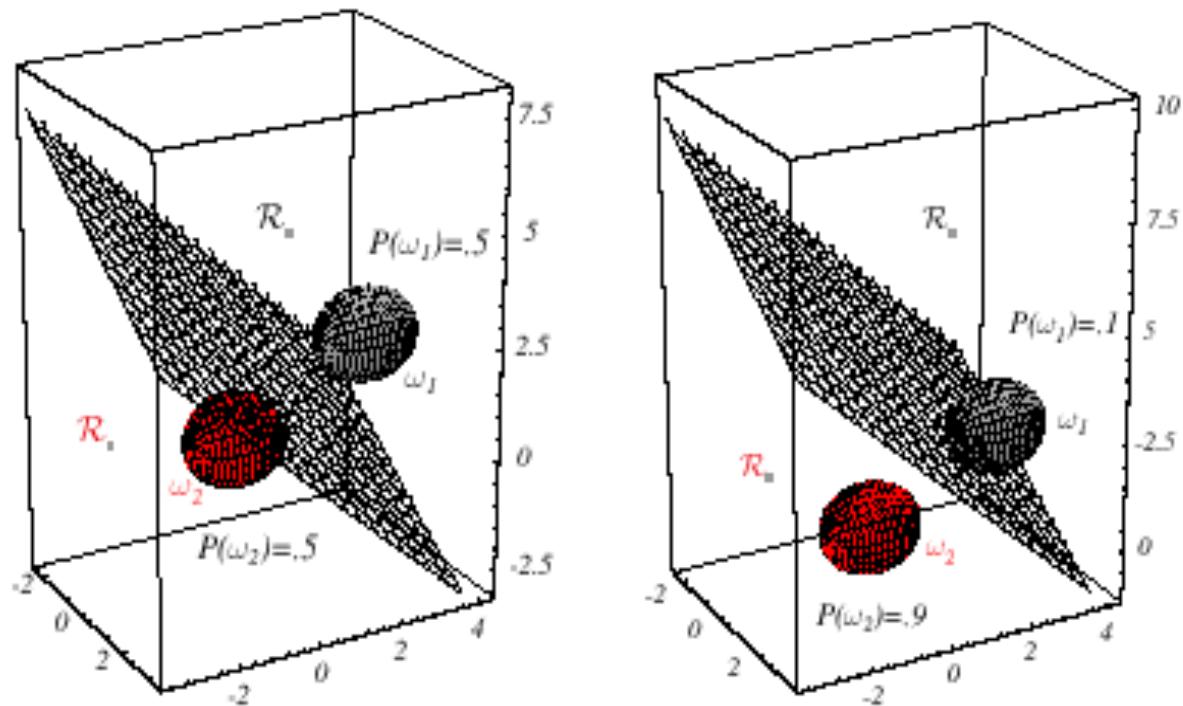


FIGURE 2.12. Probability densities (indicated by the surfaces in two dimensions and ellipsoidal surfaces in three dimensions) and decision regions for equal but asymmetric Gaussian distributions. The decision hyperplanes need not be perpendicular to the line connecting the means. From: Richard O. Duda,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

Discriminant Functions for the Normal Density...

- 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 :

$$W_i = -\frac{1}{2} \Sigma_i^{-1}$$

$$w_i = \Sigma_i^{-1} \mu_i$$

$$w_{i0} = -\frac{1}{2} \mu_i^t \Sigma_i^{-1} \mu_i - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

Discriminant Functions for the Normal Density...

Disconnected
decision
regions

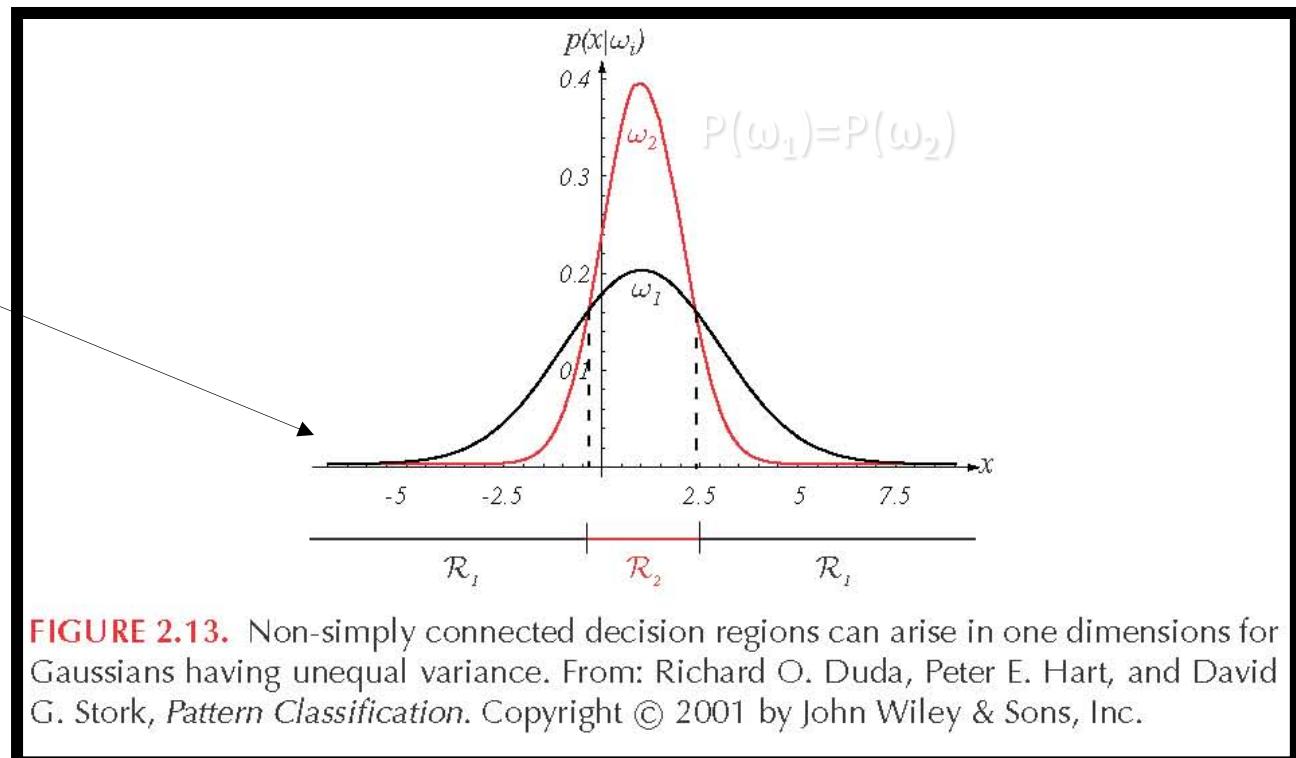
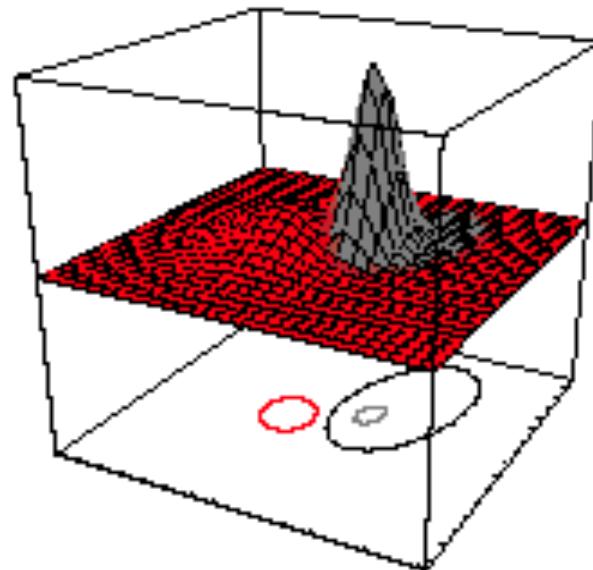
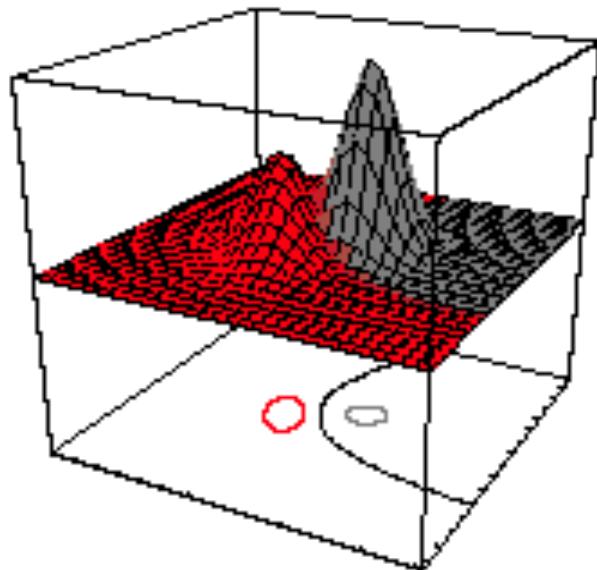
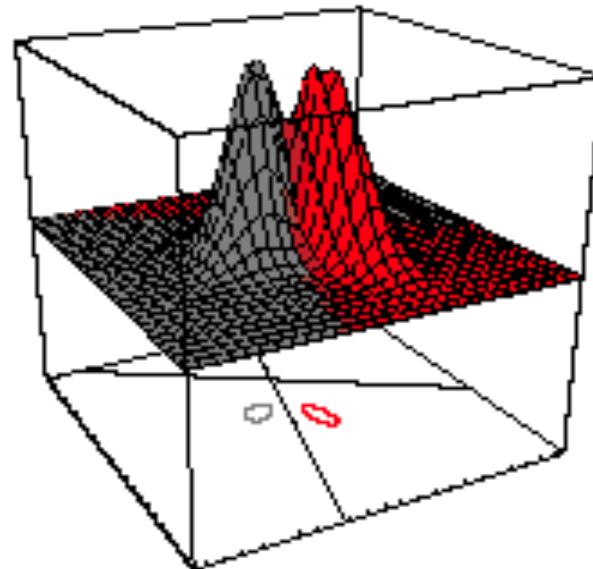
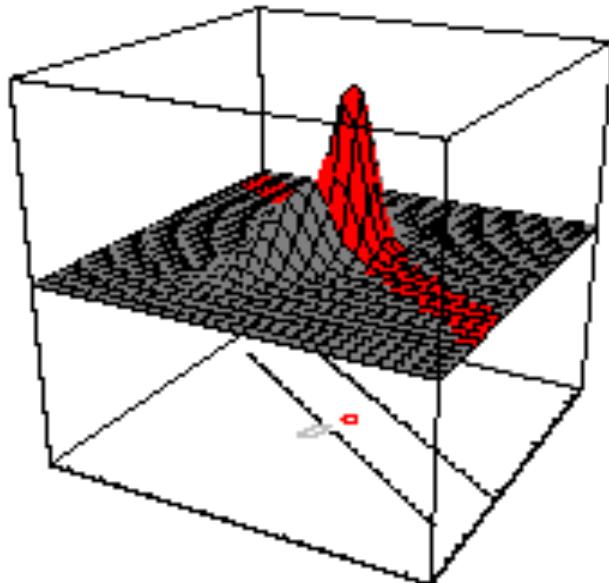


FIGURE 2.13. Non-simply connected decision regions can arise in one dimensions for Gaussians having unequal variance. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



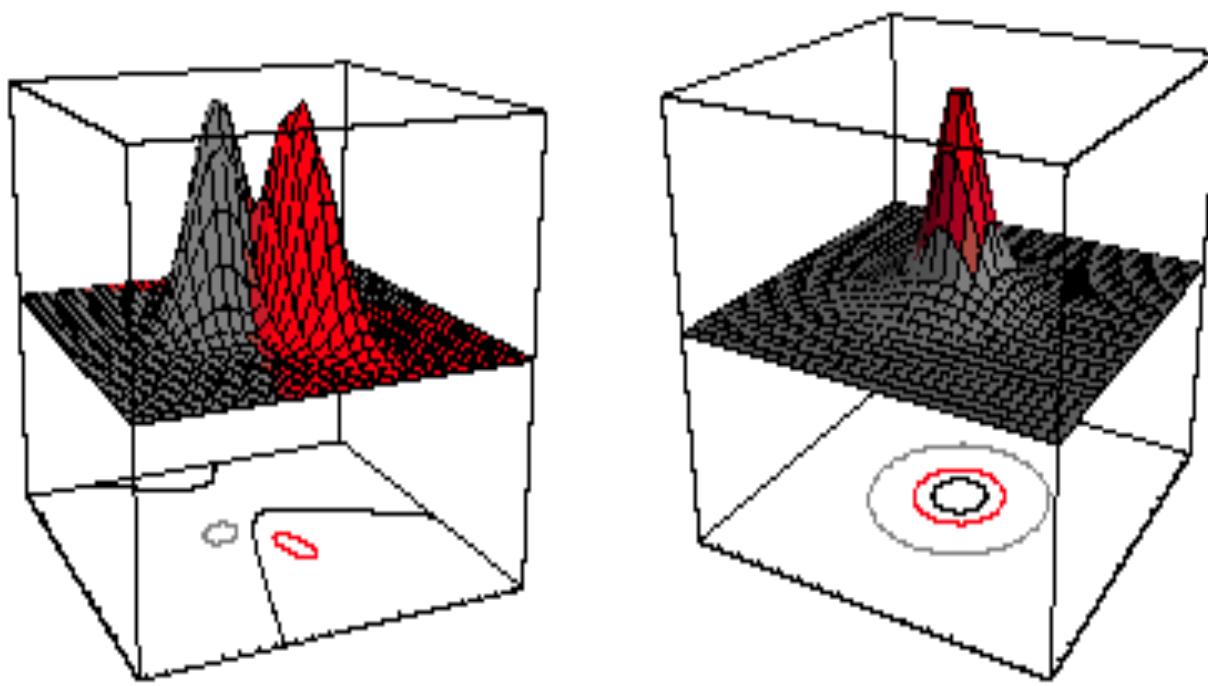


FIGURE 2.14. Arbitrary Gaussian distributions lead to Bayes decision boundaries that are general hyperquadrics. Conversely, given any hyperquadric, one can find two Gaussian distributions whose Bayes decision boundary is that hyperquadric. These variances are indicated by the contours of constant probability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

Discriminant Functions for the Normal Density

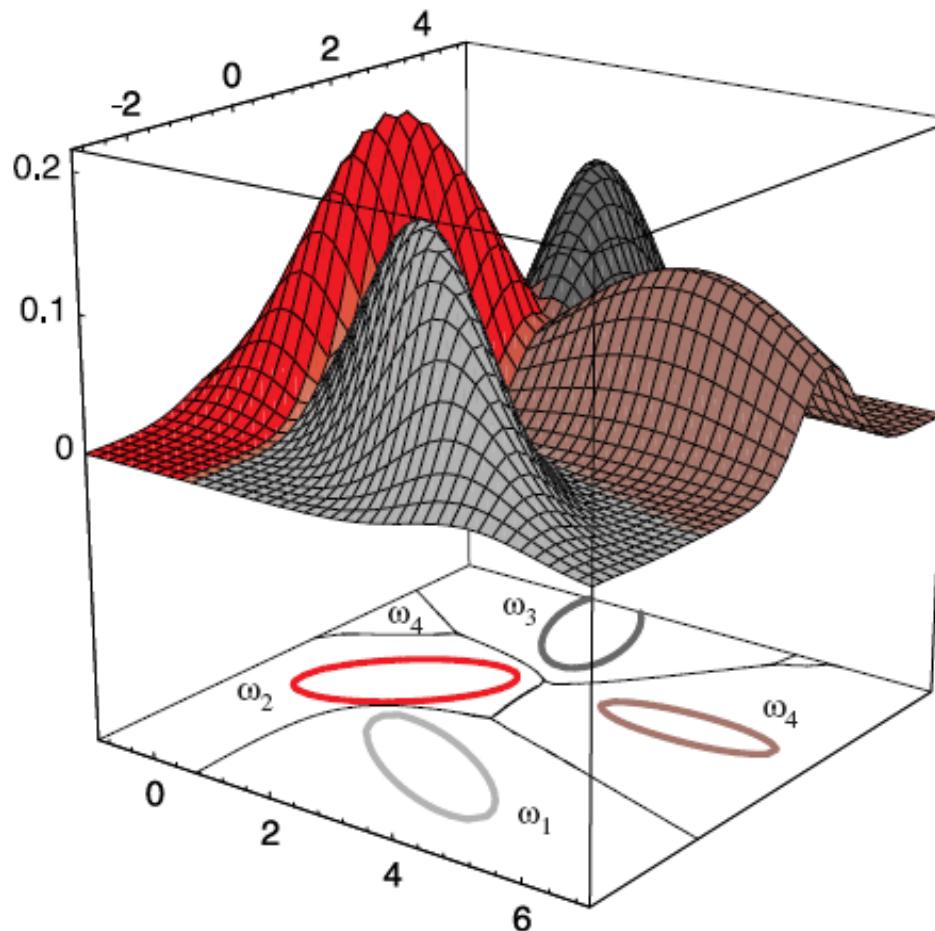
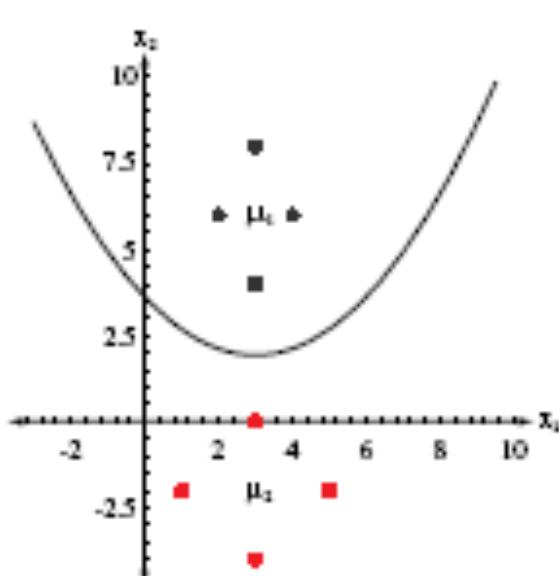


Figure 2.16: The decision regions for four normal distributions. Even with such a low number of categories, the shapes of the boundary regions can be rather complex.

Decision Regions for Two-Dimensional Gaussian Data

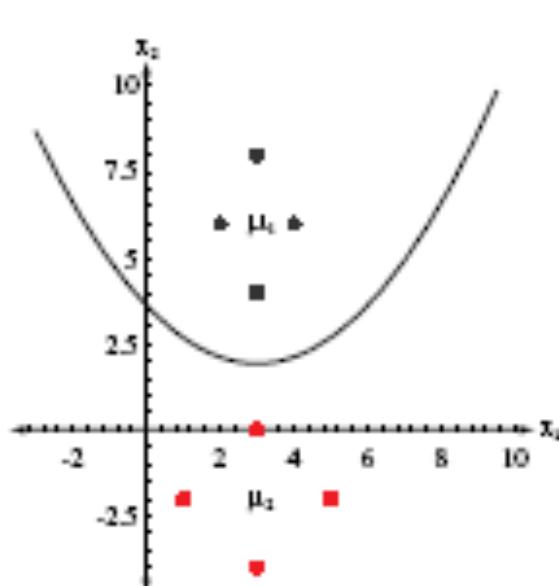


$$\mu_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix}; \quad \Sigma_1 = \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix} \text{ and } \mu_2 = \begin{bmatrix} 3 \\ -2 \end{bmatrix}; \quad \Sigma_2 = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}.$$

$$\Sigma_1^{-1} = \begin{pmatrix} 2 & 0 \\ 0 & 1/2 \end{pmatrix} \quad \text{and} \quad \Sigma_2^{-1} = \begin{pmatrix} 1/2 & 0 \\ 0 & 1/2 \end{pmatrix}.$$

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

Decision Regions for Two-Dimensional Gaussian Data



$$\mu_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix}; \quad \Sigma_1 = \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix} \text{ and } \mu_2 = \begin{bmatrix} 3 \\ -2 \end{bmatrix}; \quad \Sigma_2 = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}.$$

$$\Sigma_1^{-1} = \begin{pmatrix} 2 & 0 \\ 0 & 1/2 \end{pmatrix} \quad \text{and} \quad \Sigma_2^{-1} = \begin{pmatrix} 1/2 & 0 \\ 0 & 1/2 \end{pmatrix}.$$

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

$$x_2 = 3.514 - 1.125x_1 + 0.1875x_1^2$$

Thanks.