

# Introduction to Machine Learning and Pattern Recognition

David Brady<sup>1</sup>

<sup>1</sup>ECE Department  
Northeastern University

Fall 2018

Next:

## Intuition

## Bayes Decision Rule Intuition

## Minimum Risk Decisions

## Minimax Bayes Decision Rule

## Neyman-Pearson Detection Rule

## The Normal Distribution

## a Priori Probability

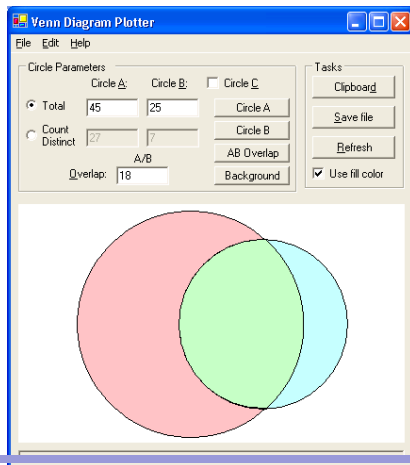
How do we use *a priori* information in classifying sea bass or salmon?

What if we know that 90% of fish are bass and 10% are salmon?

How does the best classifier change if 50% are bass and 50% are salmon?

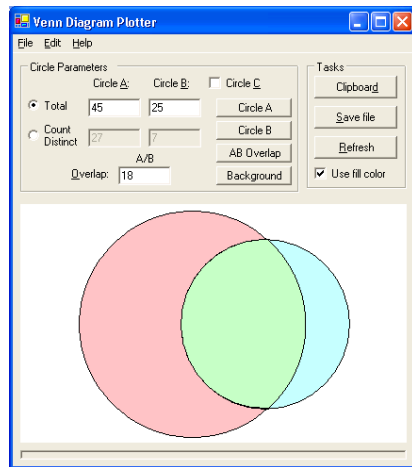
Some notation:

- ▶  $x$  is a column vector of new features
- ▶  $x$  is a point in plot for each "experiment"
- ▶ *a priori* = "before the observation"
- ▶  $\omega$  is an experimental outcome (new feature)
  - ▶  $\omega = \omega_1 \mapsto x$  in right circle
  - ▶  $\omega = \omega_2 \mapsto x$  in left circle
- ▶ *a priori* probabilities
  - ▶  $P[\omega = \omega_1] = \text{Probability } x \text{ falls in right circle}$
  - ▶  $P[\omega = \omega_2] = \text{Probability } x \text{ falls in left circle}$
- ▶  $P[\omega = \omega_1] < P[\omega = \omega_2]$
- ▶  $P[\omega_1]$  is not related to area of right circle



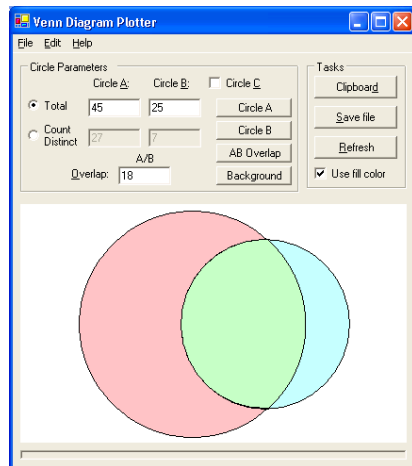
# A Posteriori Probability

- ▶ *a posteriori* = “after the observation”
- ▶ *a posteriori* probabilities
  - ▶  $P[\omega_1|x]$  = chance  $x$  in right circle after observing  $x$
  - ▶  $P[\omega_2|x]$  = chance  $x$  in left circle after observing  $x$
- ▶  $P[\omega_2] > P[\omega_1]$  (*a priori*)
- ▶  $P[\omega_2|x \text{ in green}] ? P[\omega_1|x \text{ in green}]$  (*a posteriori*)
- ▶  $P[\omega_2|x \text{ in red}] ? P[\omega_1|x \text{ in red}]$  (*a posteriori*)
- ▶  $P[\omega_2|x \text{ in blue}] ? P[\omega_1|x \text{ in blue}]$  (*a posteriori*)



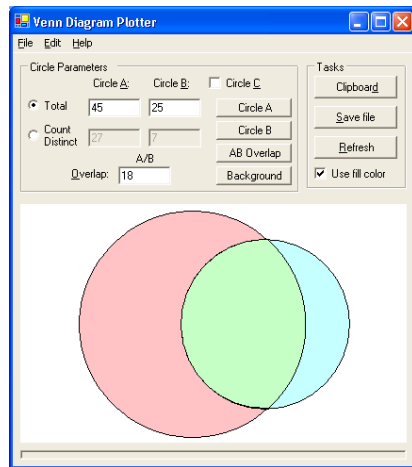
# Likelihood

- ▶  $p[\mathbf{x}|\omega_1]$  is called the *likelihood of  $\omega_1$  with respect to  $\mathbf{x}$*
- ▶ for a continuous random vector  $\mathbf{x}$ ,  $p[\mathbf{x}|\omega_1]$  is the feature vector probability density function if  $\omega = \omega_1$ 
  - ▶  $\int_{\mathbf{x} \in \mathcal{X}} p[\mathbf{x}|\omega_j] d\mathbf{x} = 1$
- ▶ for a discrete random vector  $\mathbf{x}$ ,  $p[\mathbf{x}|\omega_j]$  is the feature vector probability mass function if  $\omega = \omega_j$ 
  - ▶  $\sum_{\mathbf{x} \in \mathcal{X}} p[\mathbf{x}|\omega_j] = 1$



## Likelihood (2)

- ▶  $p[x|\omega_1] = 0$  if  $x$  outside of the right circle
- ▶ special case: uniform distributions
  - ▶  $p[x|\omega_1] = 1 / \text{area of right circle}$
  - ▶  $p[x|\omega_2] = 1 / \text{area of left circle}$
- ▶ for  $x$  in the green area,  $p[x|\omega_1] > p[x|\omega_2]$  (we say that  $\omega_1$  is more likely there)
- ▶ so, if  $x$  falls into the green area, which  $\omega_i$  would you choose?



# Bayes Formula

- ▶  $p[\mathbf{x}]$  is the *evidence* of  $\mathbf{x}$

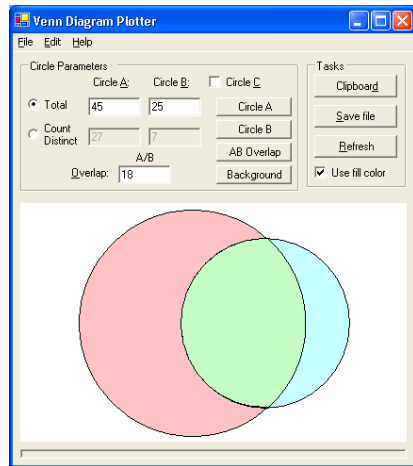


$$p[\mathbf{x}] = \sum_j P[\omega_j] p[\mathbf{x}|\omega_j]$$

- ▶  $p[\mathbf{x}]$  is the pdf (or pmf) of  $\mathbf{x}$
- ▶ Bayes Formula: *posterior = likelihood · prior / evidence*



$$P[\omega_j|\mathbf{x}] = p[\mathbf{x}|\omega_j] \cdot P[\omega_j] / p[\mathbf{x}]$$



Next:

## Intuition

## Bayes Decision Rule Intuition

## Minimum Risk Decisions

## Minimax Bayes Decision Rule

## Neyman-Pearson Detection Rule

## The Normal Distribution

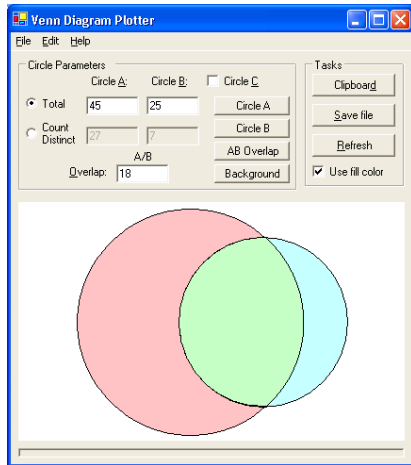


# Bayes Decision Rule

- ▶ so, if  $x$  falls into the green area, which  $\omega_j$  would you choose?
- ▶ Bayes' Decision Rule for Uniform Costs: choose  $j$  to maximize  $P[\omega_j|x]$

special case: uniform distribution, uniform costs

- ▶ choose  $\omega_1$  if  $P[\omega_1|x] = p[x|\omega_1] \cdot P[\omega_1]/p[x] > p[x|\omega_2] \cdot P[\omega_1]/p[x] = P[\omega_2|x]$
- ▶  $p[x|\omega_1] \cdot P[\omega_1] > p[x|\omega_2] \cdot P[\omega_1]$
- ▶ three cases:
- ▶ area of left circle / area of right circle  $> P[\omega_1]/P[\omega_2]$

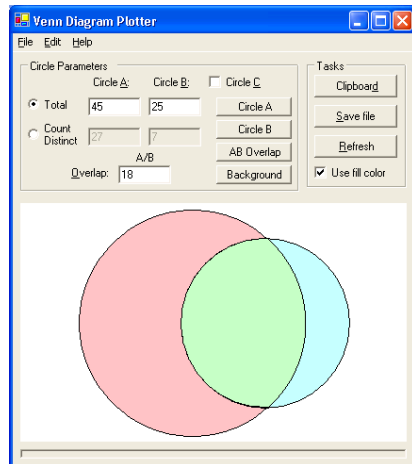


## Bayes Decision Rule (2)

- ▶ Bayes' Decision Rule for Uniform Costs: choose  $j$  to maximize  $P[\omega_j|\mathbf{x}]$

special case: uniform distribution, uniform costs

- ▶ choose  $\omega_1$  if 
$$P[\omega_1|\mathbf{x}] = p[\mathbf{x}|\omega_1] \cdot P[\omega_1]/p[\mathbf{x}] > p[\mathbf{x}|\omega_2] \cdot P[\omega_2]/p[\mathbf{x}] = P[\omega_2|\mathbf{x}]$$
- ▶ if  $p[\mathbf{x}|\omega_1] \cdot P[\omega_1] > p[\mathbf{x}|\omega_2] \cdot P[\omega_2]$ , decide  $\omega = \omega_1$
- ▶ what to do with "=" ? (later)
  - ▶ (think of each side as functions of  $\mathbf{x}$ , which is observed)



## Bayes Decision Rule (3)

- Bayes' Decision Rule says choose  $j$  to maximize  $P[\omega_j | \mathbf{x}]$

- 4 cases:

### 1. $\mathbf{x}$ falls in red area

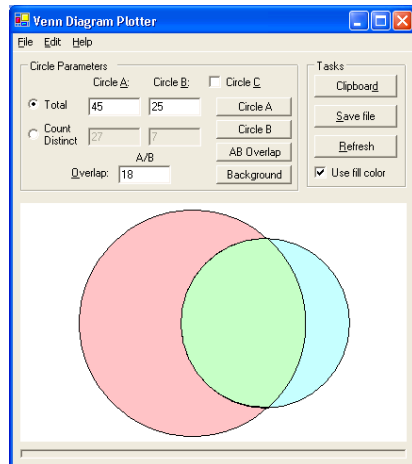
1.1  $0 \cdot P[\omega_1] > ? P[\omega_2] / \text{area of left circle} \leftarrow$   
never!

1.2 decide  $\omega = \omega_2$

### 2. $\mathbf{x}$ falls in blue area

2.1  $P[\omega_1] / \text{area of right circle} > ? 0 \cdot P[\omega_2] \leftarrow$ always!

2.2 decide  $\omega = \omega_1$



## Bayes Decision Rule (4)

- ▶ Bayes' Decision Rule : choose  $j$  to maximize  $P[\omega_j | \mathbf{x}]$

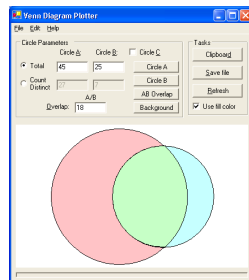
- ▶ 2 more cases:

### 3. $\mathbf{x}$ falls in white area

1.  $0 \cdot P[\omega_1] = P[\omega_1] \cdot 0$  tie!
2. decide either  $\omega = \omega_1$  or  $\omega_2$

### 4. $\mathbf{x}$ falls in green area

1.  $P[\omega_1] / \text{area of left circle} > P[\omega_2] / \text{area of right circle}$  ← depends on priors!  
right!  
left!
2. decide  $\omega = \omega_1$  if above inequality is satisfied



Next:

## Intuition

## Bayes Decision Rule Intuition

## Minimum Risk Decisions

## Minimax Bayes Decision Rule

## Neyman-Pearson Detection Rule

## The Normal Distribution

# Action Function, Conditional Loss, Conditional Risk, Average Risk

- ▶ we design *action function*
- ▶ say  $\omega_1 \iff$  action function  
 $\alpha(\mathbf{x}) = \alpha_1$
- ▶  $\alpha(\mathbf{x})$  depends only on  $\mathbf{x}$

$\lambda(\alpha_i, \omega_j)$  is the *conditional loss* of action  $\alpha(\mathbf{x}) = \alpha_i$  when  $\omega_j$  occurs

$\lambda(\alpha_i   \omega_j)$	$\omega_1$	$\omega_2$
$\alpha_1$	0USD	43USD
$\alpha_2$	9700USD	0USD

- ▶ some mistakes are more

- ▶ Conditional Risk is  
 $R(\alpha(\mathbf{x}) | \mathbf{x})$
- ▶ depends on observation and action rule
- ▶  $R(\alpha(\mathbf{x}) = \alpha_i | \mathbf{x}) = \sum_j \lambda(\alpha_i | \omega_j) P[\omega_j | \mathbf{x}]$
- ▶ average risk  
 $R = \int R(\alpha(\mathbf{x}) | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} = E[R(\alpha(\mathbf{x}) | \mathbf{x})]$
- ▶ Bayes Decision Rule  
 $\alpha_{\text{Bayes}}(\mathbf{x})$  minimizes *average risk*

## Minimum Risk Solutions

- ▶ special case: two classes  $\omega_1$  and  $\omega_2$
- ▶ two actions:  $\alpha_1$  and  $\alpha_2$
- ▶ Bayes classifiers minimize conditional risk
- ▶  $R(\alpha_i|\mathbf{x}) = \lambda(\alpha_i|\omega_1)P[\omega_1|\mathbf{x}] + \lambda(\alpha_i|\omega_2)P[\omega_2|\mathbf{x}]$
- ▶ for each  $\mathbf{x}$ ,  $\alpha_{\text{Bayes}}(\mathbf{x}) = \text{argmin} R(\alpha_i|\mathbf{x})$
- ▶ say  $\omega_1$  if:  $R(\alpha_1|\mathbf{x}) < R(\alpha_2|\mathbf{x})$ , or
- ▶  $\lambda(\alpha_1|\omega_1)P[\omega_1|\mathbf{x}] + \lambda(\alpha_1|\omega_2)P[\omega_2|\mathbf{x}] <$   
 $\lambda(\alpha_2|\omega_1)P[\omega_1|\mathbf{x}] + \lambda(\alpha_2|\omega_2)P[\omega_2|\mathbf{x}]$
- ▶  $P[\omega_1|\mathbf{x}](\lambda(\alpha_2|\omega_1) - \lambda(\alpha_1|\omega_1)) > P[\omega_2|\mathbf{x}](\lambda(\alpha_1|\omega_2) - \lambda(\alpha_2|\omega_2))$
- ▶  $\frac{P[\omega_1|\mathbf{x}]}{P[\omega_2|\mathbf{x}]} > \frac{\lambda(\alpha_1|\omega_2) - \lambda(\alpha_2|\omega_2)}{\lambda(\alpha_2|\omega_1) - \lambda(\alpha_1|\omega_1)}$
- ▶  $\frac{P[\mathbf{x}|\omega_1]}{P[\mathbf{x}|\omega_2]} > \frac{P[\omega_2]}{P[\omega_1]} \frac{\lambda(\alpha_1|\omega_2) - \lambda(\alpha_2|\omega_2)}{\lambda(\alpha_2|\omega_1) - \lambda(\alpha_1|\omega_1)} \implies \text{say } \omega_1$
- ▶  $\frac{P[\mathbf{x}|\omega_1]}{P[\mathbf{x}|\omega_2]} \stackrel{\text{red}}{=} \frac{P[\omega_2]}{P[\omega_1]} \frac{\lambda(\alpha_1|\omega_2) - \lambda(\alpha_2|\omega_2)}{\lambda(\alpha_2|\omega_1) - \lambda(\alpha_1|\omega_1)} \implies \text{say } \omega_1 \text{ with prob. } p$

## Randomization

- ▶ randomization:
- ▶  $\frac{P[\mathbf{x}|\omega_1]}{P[\mathbf{x}|\omega_2]} = \frac{P[\omega_2]}{P[\omega_1]} \frac{\lambda(\alpha_1|\omega_2) - \lambda(\alpha_2|\omega_2)}{\lambda(\alpha_2|\omega_1) - \lambda(\alpha_1|\omega_1)} \implies$  say  $\omega_1$  with prob.  $p$
- ▶ all values of  $p$  produce the same average risk, so  $p = 0, 1$  is fine (no randomization)
- ▶  $p$  is sometime used to ease calculations (later)



## Interpretation of Bayes Decision Rule

- ▶  $\frac{P[\mathbf{x}|\omega_1]}{P[\mathbf{x}|\omega_2]} > \frac{P[\omega_2]}{P[\omega_1]} \frac{\lambda(\alpha_1|\omega_2) - \lambda(\alpha_2|\omega_2)}{\lambda(\alpha_2|\omega_1) - \lambda(\alpha_1|\omega_1)}$
- ▶  $\frac{P[\mathbf{x}|\omega_1]}{P[\mathbf{x}|\omega_2]}$  is the likelihood ratio
- ▶  $\frac{P[\omega_2]}{P[\omega_1]}$  is the ratio of a priori probabilities
- ▶  $\frac{\lambda(\alpha_1|\omega_2) - \lambda(\alpha_2|\omega_2)}{\lambda(\alpha_2|\omega_1) - \lambda(\alpha_1|\omega_1)}$  is the ratio of relative costs
- ▶ if the likelihood ratio is sufficiently large, accept  $\omega_1$
- ▶ if the ratio of prior probabilities is sufficiently large, reject  $\omega_1$
- ▶ if the relative cost of accepting  $\omega_1$  is sufficiently large, reject  $\omega_1$

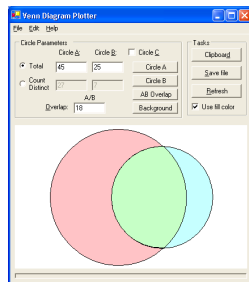
## Example of Bayes Decision Rule

- ▶ right circle ( $\omega_1$ ) has area = 25
- ▶ left circle ( $\omega_2$ ) has area = 45
- ▶ overlap has area=18
- ▶ uniform distribution on circle for  $\omega_j$

$\lambda(\alpha_i \omega_j)$	$\omega_1$	$\omega_2$
$\alpha_1$	0	1
$\alpha_2$	1	0

uniform costs

- ▶ ratio of relative costs = 1
- ▶ likelihood ratio( $x \in \text{red}$ ) = 0
- ▶ likelihood ratio( $x \in \text{blue}$ ) =  $\infty$
- ▶ likelihood ratio( $x \in \text{white}$ ) undefined

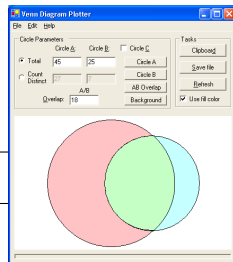


## Example of Bayes Decision Rule (cont'd)

- ▶  $\frac{P[\omega_2]}{P[\omega_1]}$  is variable
- ▶ Bayes Decision Rule:
- ▶  $\frac{P[x|\omega_1]}{P[x|\omega_2]} > \frac{P[\omega_2]}{P[\omega_1]} \frac{\lambda(\alpha_1|\omega_2) - \lambda(\alpha_2|\omega_2)}{\lambda(\alpha_2|\omega_1) - \lambda(\alpha_1|\omega_1)} \Rightarrow \text{say } \omega_1$

$$\frac{P[x|\omega_1]}{P[x|\omega_2]} > \frac{P[\omega_2]}{P[\omega_1]} \cdot 1 \Rightarrow \text{say } \omega_1$$

- ▶  $x \in \text{red}$ , accept  $\omega_2$
- ▶  $x \in \text{blue}$ , accept  $\omega_1$
- ▶  $x \in \text{white}$ , undefined
- ▶  $x \in \text{green}$ ,  $45/25 \stackrel{?}{>} \frac{P[\omega_2]}{P[\omega_1]}$ . If so, accept  $\omega_1$ .
- ▶  $x \in \text{green}$ ,  $45/25 \stackrel{?}{>} \frac{P[\omega_2]}{P[\omega_1]}$ . If so, accept  $\omega_1$  with probability  $p$



- ▶ let  $P[\omega_1] = P[\omega_2] = 0.5$ , so  $\frac{P[\omega_2]}{P[\omega_1]} = 1$ .
- ▶ randomization not needed
- ▶ if  $\omega_2$  is sufficiently rare, accept  $\omega_1$  here

## Conditional Risk For Uniform Costs

conditional risk:

*confusing*

*better*  $r(\alpha_i | \underline{x}) \not\equiv R(\alpha(x) = \alpha_i | \underline{x}) = \sum_j \lambda(\alpha_i | \omega_j) P[\omega_j | \underline{x}]$

uniform costs:  $\lambda(\alpha_i | \omega_j) = 1 - \delta_{ij} = \begin{cases} 1, & \text{if } i \neq j \\ 0, & \text{if } i = j \end{cases}$

conditional risk for uniform costs:

$$r(\alpha_i | \underline{x}) \not\equiv \sum_{j \neq i} P[\omega_j | \underline{x}] = P[\text{error} | \underline{x}]$$

*correct*

$$= 1 - P[\omega_i | \underline{x}]$$

*also good*

For uniform costs, Bayes Decision Rule minimizes the error probability for each  $\underline{x}$ ,  $P[\text{error} | \underline{x}]$ , and the average error probability  $R = E[P[\text{error} | \underline{x}]] = P[\text{error}]$ .

## Computing Average Risk

average risk:

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

*Handwritten notes:*  
 $\alpha(\underline{x}) =$  our decision at  $\underline{x}$   
 $R(\alpha(\underline{x})|\underline{x})$  is the risk of that decision

$p(x) = \sum_j P(\omega_j) p(x|\omega_j)$ , so one way to compute average risk is by conditioning on  $\omega_j$  first, then average over a priori probabilities:

First find

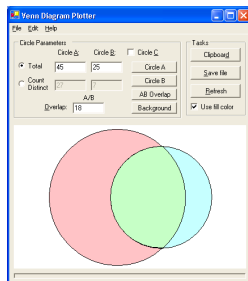
$$R_j = \int R(\alpha(x)|x) p(x|\omega_j) dx$$

$$\text{then } R = \sum_j P(\omega_j) R_j$$

where  $R_j$  is the conditional risk <sup>given</sup>  $\omega_j$  occurs

There are other ways! (homework)

## Example: Computing Average Risk



- ▶ in summary, our Bayes Decision Rule is:
- ▶  $\alpha_1$  if  $\mathbf{x}$  in blue or green
- ▶  $\alpha_2$  if  $\mathbf{x}$  in red
- ▶  $R = P[\text{error}]$  for uniform costs
- ▶ First find risk under each class,  $R_i$
- ▶  $R_1 = P[\text{error}|\omega_1] = 0$
- ▶  $R_2 = P[\text{error}|\omega_2] = P[\omega_2]P[\mathbf{x} \in \text{green}|\omega_2]$
- ▶ average over a priori probabilities
- ▶  $P[\text{error}] = 0.5 \cdot 0 + 0.5 \frac{18}{45} = 0.20$

# Next:

Intuition

Bayes Decision Rule Intuition

Minimum Risk Decisions

**Minimax Bayes Decision Rule**

Neyman-Pearson Detection Rule

Classifiers, Discriminant Functions, Decision Surfaces

The Normal Distribution

Discriminant Functions for the Normal Density

Signal Detection Theory and Operating Characteristics

## If the *a priori* probabilities are not known?

- ▶ Consider two-class case, and suppose  $P[\omega_1]$  is unknown
- ▶ Let  $\lambda_{ij} = \lambda(\alpha_i|\omega_j)$
- ▶ Let  $\mathcal{R}_i$  be the decision region for class  $\omega_i$
- ▶ Rewrite  $R_i$

$$R_2 = \lambda_{22} + (\lambda_{12} - \lambda_{22}) \int_{\mathcal{R}_1} p(\mathbf{x}|\omega_2) d\mathbf{x}$$

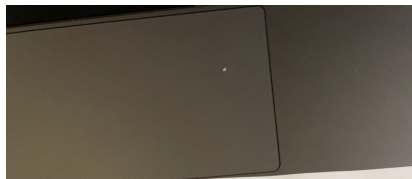
$$R_1 = \lambda_{11} + (\lambda_{21} - \lambda_{11}) \int_{\mathcal{R}_2} p(\mathbf{x}|\omega_1) d\mathbf{x}$$

$$R = P[\omega_1]R_1 + P[\omega_2]R_2 = R_2 + P[\omega_1](R_1 - R_2)$$

- ▶ want risk to be independent of unknown  $P[\omega_1] \implies$  find  $\mathcal{R}_i$  so that  $R_1 = R_2$
- ▶ then  $R_{\min\max} = R_1 = R_2$



# Finding the Minimax Bayes Decision Rule

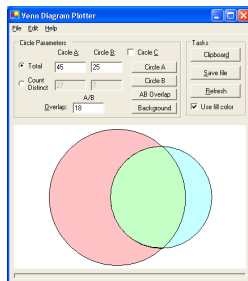


2.4 ■ CLASSIFIERS, DISCRIMINANT FUNCTIONS, AND DECISION SURFACES 29

**FIGURE 2.4.** The curve at the bottom shows the minimum (Bayes) error as a function of prior probability  $P(\omega_1)$  in a two-category classification problem of fixed distributions. For each value of the priors (e.g.,  $P(\omega_1) = 0.25$ ) there is a corresponding optimal decision boundary and associated Bayes error rate. For any (fixed) such boundary, if the priors change, the probability of error will change as a linear function of  $P(\omega_1)$  (shown by the dashed line). The maximum such error will occur at an extreme value of the prior, here at  $P(\omega_1) = 1$ . To minimize the maximum of such error, we should design

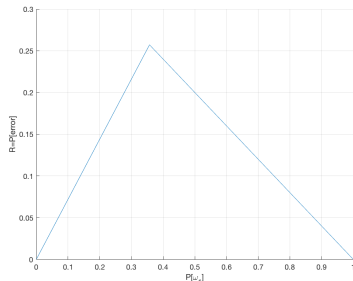
- ▶ Note that when  $R_1 = R_2$ ,  $\frac{\delta R}{\delta P[\omega_1]} = 0$
- ▶  $R_{Bayes}$  is maximized at the minmax Bayes Decision Rule
- ▶ 2 ways to find the minimax decision rule:
  - ▶ fix  $P[\omega_1]$ , and find  $R_i$  for the Bayes Decision Rule ( $\mathcal{R}_i$ ). Then either:
    1. choose the decision rule ( $P[\omega_1], \mathcal{R}_1$ ) so that  $R_1 = R_2$ , or
    2. choose the decision rule so that  $R = P[\omega_1](R_1 - R_2) + R_2$  is maximum

## Example: Minimax Bayes Decision Rule



- ▶  $R = P[\omega_1]R_1 + P[\omega_2]R_2 = P[\text{error}]$
- ▶  $\mathbf{x} \in \text{green}$ , accept  $\omega_1$  if  $\frac{45}{25} > \frac{P[\omega_2]}{P[\omega_1]}$ . (if equality, flip biased coin to accept!)
- ▶  $R_1 = \begin{cases} \frac{18}{25}, & P[\omega_1] < \frac{5}{14} \\ 0, & P[\omega_1] > \frac{5}{14}. \end{cases}$
- ▶  $R_2 = \begin{cases} 0 & P[\omega_1] < \frac{5}{14} \\ \frac{18}{45}, & P[\omega_1] > \frac{5}{14}. \end{cases}$
- ▶  $R = \begin{cases} P[\omega_1]\frac{18}{25} + 0, & P[\omega_1] < \frac{5}{14} \\ 0 + \frac{18}{45}(1 - P[\omega_1]), & P[\omega_1] > \frac{5}{14} \end{cases}$

## Example: Minmax Bayes Decision Rule (2)



- ▶ worst-case risk  $R = 9/35$  at  $P[\omega_1] = 5/14$
- ▶ minmax Bayes detector is Bayes Decision Rule at  $P[\omega_1] = 5/14$ .
- ▶  $\mathbf{x} \in \text{green}$ , accept  $\omega_1$  with probability  $p$  (flip coin)
- ▶ find  $p$ :  $R_1 = (1 - p) \frac{18}{25}$  ;  
 $R_2 = p \frac{18}{45}$
- ▶  $R_1 = R_2$  or  $R = R_2$  implies  
 $p = \frac{9}{14}$

# Next:

Intuition

Bayes Decision Rule Intuition

Minimum Risk Decisions

Minimax Bayes Decision Rule

**Neyman-Pearson Detection Rule**

Classifiers, Discriminant Functions, Decision Surfaces

The Normal Distribution

Discriminant Functions for the Normal Density

Signal Detection Theory and Operating Characteristics

## What to do if costs and priors are unknown?

- ▶ consider the  $j$ -class case
- ▶ let  $\omega_1$  correspond to the *null* class
- ▶ consider all decision rules satisfying a (false-alarm, level) constraint:  $R_1 \leq \text{constant}$
- ▶ create a detector which minimizes the residual risk:  

$$\min \sum_{j \neq 1} R_j$$
- ▶ for  $j=2$ , the NP-optimal detector is a likelihood-ratio test

$$\frac{p(\mathbf{x}|\omega_2)}{p(\mathbf{x}|\omega_1)} > t, \text{ accept } \omega_2$$

$$\frac{p(\mathbf{x}|\omega_2)}{p(\mathbf{x}|\omega_1)} = t, \text{ accept } \omega_2 \text{ with prob } p$$

- ▶ threshold  $t$  and probability  $p$  are set by the above (false-alarm, level) constraint

# Next:

Intuition

Bayes Decision Rule Intuition

Minimum Risk Decisions

Minimax Bayes Decision Rule

Neyman-Pearson Detection Rule

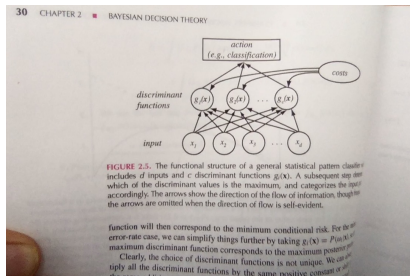
**Classifiers, Discriminant Functions, Decision Surfaces**

The Normal Distribution

Discriminant Functions for the Normal Density

Signal Detection Theory and Operating Characteristics

# Discriminant Functions and Classifiers



- ▶ a *discriminant function*  $g_i(\cdot)$  is a functional mapping a feature  $\mathbf{x}$  to a measure of fit for class  $\omega_i$
- ▶ a *classifier* assigns each feature  $\mathbf{x}$  to one of  $c$  classes using  $c$  discriminator function comparisons
- ▶  $\mathbf{x}$  is assigned to class  $i$  if  $i = \operatorname{argmax}_j g_j(\mathbf{x})$
- ▶ Bayes classifier uses  $g_i(\mathbf{x}) = -R(\alpha_i|\mathbf{x})$

# Faster Code $\Leftrightarrow$ Simplify Discriminant Functions!

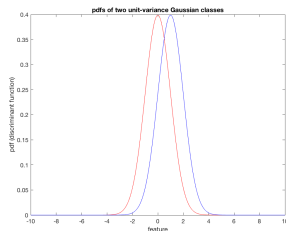
$$\begin{aligned}
 \operatorname{argmax}_j g_j(\mathbf{x}) &= \operatorname{argmax}_j \log(g_j(\mathbf{x})) \\
 &= \operatorname{argmax}_j \exp(g_j(\mathbf{x})) \\
 &= \operatorname{argmax}_j 43 \cdot (g_j(\mathbf{x}))
 \end{aligned}$$

- ▶ coding and analysis are eased by simplifying comparisons
- ▶ the same Bayes classifier can use any of the following sets:

$$\begin{aligned}
 g_i(\mathbf{x}) &= \frac{p(\mathbf{x}|\omega_i)P(\omega_i)}{\sum_{j=1}^c p(\mathbf{x}|\omega_j)P(\omega_j)} \\
 &= p(\mathbf{x}|\omega_i)P(\omega_i) \\
 &= \ln(p(\mathbf{x}|\omega_i)) + \ln(P(\omega_i)) \\
 \text{goal...} &\stackrel{?}{=} \text{simple function of } \mathbf{x}
 \end{aligned}$$



## Example: Bayesian Gaussian Classifier for Equal Prior Probabilities



- ▶  $g_i(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu_i)^2}$ ,  $\mu_1 = 0$ ,  $\mu_2 = 1$
- ▶  $g_1(x) \stackrel{?}{>} g_2(x)$
- ▶  $\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu_1)^2} \stackrel{?}{>} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu_2)^2}$
- ▶  $e^{-\frac{1}{2}(x-\mu_1)^2} \stackrel{?}{>} e^{-\frac{1}{2}(x-\mu_2)^2}$
- ▶  $-\frac{1}{2}(x-\mu_1)^2 \stackrel{?}{>} -\frac{1}{2}(x-\mu_2)^2$
- ▶  $(x-\mu_1)^2 \stackrel{?}{<} (x-\mu_2)^2$
- ▶  $x \stackrel{?}{<} \frac{1}{2}$ , much faster to execute!

## Next:

Intuition

Bayes Decision Rule Intuition

Minimum Risk Decisions

Minimax Bayes Decision Rule

Neyman-Pearson Detection Rule

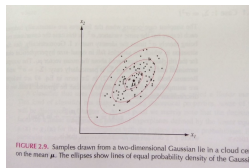
Classifiers, Discriminant Functions, Decision Surfaces

**The Normal Distribution**

Discriminant Functions for the Normal Density

Signal Detection Theory and Operating Characteristics

# Normal Probability Density Function



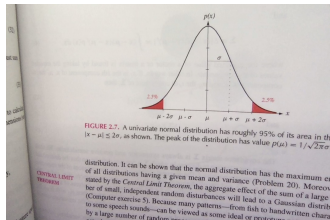
scalar Gaussian random variable is characterized by  
mean:  $\mu$ , variance:  $\sigma^2$

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$$

d-dimensional Gaussian random vector is characterized by its mean vector  $\mu$  and covariance matrix  $\Sigma$

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

# Gaussian Moments



- mean:

$$\mu = E[x] = \int_{-\infty}^{\infty} p(x)x dx = [\mu_1, \mu_2, \dots, \mu_d]^t$$

- covariance matrix:

$$\Sigma = E[(x - \mu)(x - \mu)^t]$$

- $\Sigma = \Sigma^t$  (symmetric, positive semi-definite)

- Average Value of

$$f(x) : E[f(x)] = \int_{-\infty}^{\infty} f(x)p(x)dx$$

- Example:  $f(x) = -\ln(p(x)) \implies$  entropy,  $H(x)$

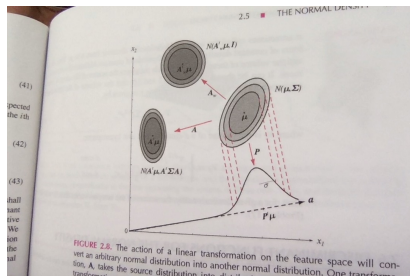
- $H(x) = -E[\ln(p(x))] = \frac{1}{2} + \log_2 \sqrt{2\pi\sigma^2}$

- Gaussian has largest entropy of any continuous r.v. having same mean  $\mu$  and variance  $\sigma^2$

- independence  $\Leftrightarrow$  uncorrelatedness

# Linear Transformations of Gaussian Vectors

- ▶  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- ▶ let  $\mathbf{A}$  be a deterministic matrix, and let  $\mathbf{y} = \mathbf{A}^t \mathbf{x}$
- ▶  $\mathbf{y} \sim \mathcal{N}(\mathbf{A}^t \boldsymbol{\mu}, \mathbf{A}^t \boldsymbol{\Sigma} \mathbf{A})$
- ▶ special case:  $\boldsymbol{\Sigma} = \boldsymbol{\Phi} \boldsymbol{\Lambda} \boldsymbol{\Phi}^t$ , where
  - ▶ columns of  $\boldsymbol{\Phi}$  are eigenvectors of  $\boldsymbol{\Sigma}$  (orthonormal set)
  - ▶ diagonal of  $\boldsymbol{\Lambda}$  contain the eigenvalues of  $\boldsymbol{\Sigma}$
  - ▶ if  $\mathbf{A}_w = \boldsymbol{\Phi} \boldsymbol{\Lambda}^{-1/2}$ , then  $\mathbf{A}_w^t \boldsymbol{\Sigma} \mathbf{A}_w = \mathbf{I}$  (whitening transformation)



## Next:

Intuition

Bayes Decision Rule Intuition

Minimum Risk Decisions

Minimax Bayes Decision Rule

Neyman-Pearson Detection Rule

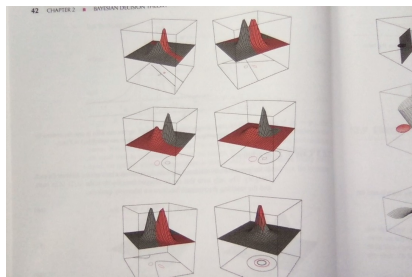
Classifiers, Discriminant Functions, Decision Surfaces

The Normal Distribution

**Discriminant Functions for the Normal Density**

Signal Detection Theory and Operating Characteristics

# Minimum Error Probability Discriminant Functions



- ▶ Bayes Gaussian Discriminant Function:

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

- ▶ for each  $\mathbf{x}$ , choose class  $\omega_k$  if
- ▶  $k = \arg \max_i g_i(\mathbf{x})$

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

- ▶  $(\mathbf{x} - \mu_i)^t \Sigma_i^{-1}(\mathbf{x} - \mu_i)$  is the Mahalanobis <sup>squared</sup> distance from  $\mathbf{x}$  to  $\mu_i$
- ▶ what are the shape of the decision boundaries in feature space?

# White Gaussian Vectors

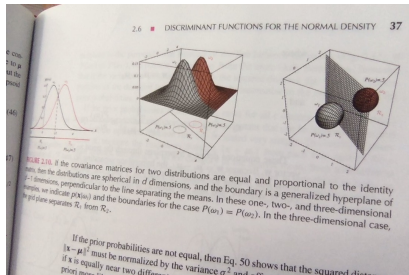
- ▶ spherical Gaussian vectors:  
 $\Sigma_i = \sigma^2 I$

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

- ▶ remove terms & factors common to all  $i$

- ▶  $g_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x} + w_{i0}$ ,

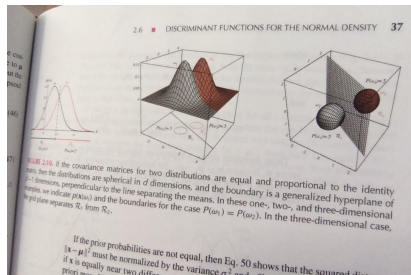
- ▶  $\mathbf{w}_i = \frac{1}{\sigma^2} \mu_i$ ,  $w_{i0} = -\frac{1}{2\sigma^2} \mu_i^t \mu_i + \ln P[\omega_i]$



- ▶  $w_{i0}$  is the  $i$ th *bias* or *threshold*
- ▶  $\mathbf{w}_i^t \mathbf{x}$  is a linear operator on the feature vector
- ▶ a *linear machine* is a classifier which uses such  $g_i(\mathbf{x})$



# Decision Boundaries for White Gaussian Vectors



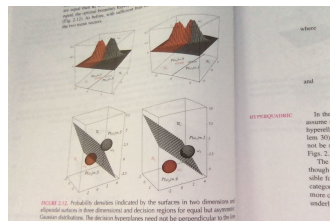
- ▶  $g_1(\mathbf{x}) = g_2(\mathbf{x})$  specifies  $\mathbf{x}$  on a decision boundary, or
- ▶  $\mathbf{w}_{12}^t (\mathbf{x} - \mathbf{x}_0) = 0$ , with
- ▶  $\mathbf{w}_{12} = \mu_1 - \mu_2$ , and
- ▶  $\mathbf{x}_0 = \frac{1}{2} (\mu_1 + \mu_2) - \frac{\sigma^2}{\|\mu_1 - \mu_2\|^2} \ln \frac{P[\omega_1]}{P[\omega_2]} (\mu_1 - \mu_2)$

- ▶ boundary is a hyperplane
  - ▶ orthogonal to  $\mathbf{w}_{12}$ ,
  - ▶ passes through  $\mathbf{x}_0$
- ▶ decision regions are contiguous

# Commonly-colored Gaussian Vectors

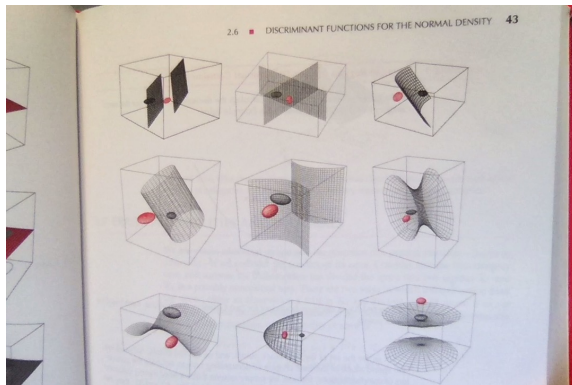
$$\mathbf{x} \sim \mathcal{N}(\mu_i, \Sigma)$$

- ▶  $\Sigma_i = \Sigma$
- ▶  $\mathbf{y} = \mathbf{A}_w \mathbf{x}$ , then colored Gaussian becomes white in  $\mathbf{y}$ -feature space
- ▶ in  $\mathbf{x}$ -feature space, the discriminant functions are:
- ▶  $g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mu_i)^t \Sigma^{-1}(\mathbf{x} - \mu_i) + \ln P[\omega_i]$  is equivalent to:
- ▶  $g_i(\mathbf{x}) = \mathbf{q}_i^t \mathbf{x} + c_{i0}$ , where
  - ▶  $\mathbf{q}_i = \Sigma^{-1} \mu_i$
  - ▶  $c_{i0} = -\frac{1}{2} \mu_i^t \Sigma^{-1} \mu_i + \ln P[\omega_i]$



- ▶ boundary is a hyperplane
- ▶ not orthogonal to mean difference

## Decision Boundaries - Arbitrary Gaussian Vectors



- ▶  $g_i(\mathbf{x}) = \mathbf{x}^t \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^t \mathbf{x} + w_{i0}$
- ▶  $\mathbf{W}_i = -\frac{1}{2} \boldsymbol{\Sigma}_i^{-1}$
- ▶  $\mathbf{w}_i = \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i$

$$w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^t \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P[\omega_i]$$

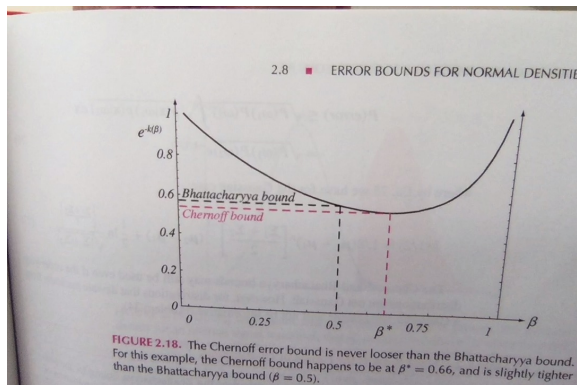
## Binary Error Probability Calculation

- ▶ Bayes Formula:
- ▶  $P[\omega_j|\mathbf{x}] = p(\mathbf{x}|\omega_j)P[\omega_j]/p(\mathbf{x})$
- ▶  $p(\mathbf{x}) = \sum_{j=1}^2 p(\mathbf{x}|\omega_j)P[\omega_j]$
- ▶ Bayes Decision Rule (uniform costs):
- ▶ decide  $\omega_1$  if  $P[\omega_1|\mathbf{x}] > P[\omega_2|\mathbf{x}]$
- ▶  $P(\text{error}|\mathbf{x}) = \min(P(\omega_1|\mathbf{x}), P(\omega_2|\mathbf{x}))$
- ▶  $P(\text{error}) = \int P(\text{error}|\mathbf{x})p(\mathbf{x})d\mathbf{x}$
- ▶  $\min[a, b] \leq a^\beta b^{1-\beta}$ , for  $a, b \geq 0, 0 \leq \beta \leq 1$
- ▶ Chernoff Bound
- ▶  $P(\text{error}) \leq P^\beta(\omega_1)P^{1-\beta}(\omega_2) \int p^\beta(\mathbf{x}|\omega_1)p^{1-\beta}(\mathbf{x}|\omega_2)$
- ▶ Bhattacharyya Bound: set  $\beta = 1/2$
- ▶ No integration over decision regions!

## Gaussian Chernoff Bounds

$$\triangleright P(\text{error}) \leq P^\beta(\omega_1)P^{1-\beta}(\omega_2)e^{-k(\beta)},$$

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



- $\triangleright$  1 parameter to optimize!

## Next:

Intuition

Bayes Decision Rule Intuition

Minimum Risk Decisions

Minimax Bayes Decision Rule

Neyman-Pearson Detection Rule

Classifiers, Discriminant Functions, Decision Surfaces

The Normal Distribution

Discriminant Functions for the Normal Density

Signal Detection Theory and Operating Characteristics

# Signal Detection Theory

- ▶ two classes:  $\omega_1$  = null class,  $\omega_2$  = alternative class
- ▶ observe  $\mathbf{x}$ , decide for  $\omega_1$  (negative)
  - ▶ or decide for  $\omega_2$  (positive)

Four possible events:

- ▶ *Hit*: true positive
- ▶ *False alarm*: false positive
- ▶ *Miss*: false negative
- ▶ *Correct Rejection*: true negative

# Decision Events

		Condition (as determined by "Gold Standard" or "Ground Truth")		
		Positive	Negative	
Test outcome	Positive	True Positive $T_p$	False Positive (Type I Error) $F_p$	→ Positive Predicted Value
	Negative	False Negative (Type II Error) $F_n$	True Negative $T_n$	→ Negative Predicted Value
		↓ Sensitivity	↓ Specificity	

$$= \frac{T_p}{T_p + F_n}$$

$$= \frac{T_n}{F_p + T_n}$$

Sensitivity  
1 - Specificity



## Alternative Event Labels

### Medical Diagnoses:

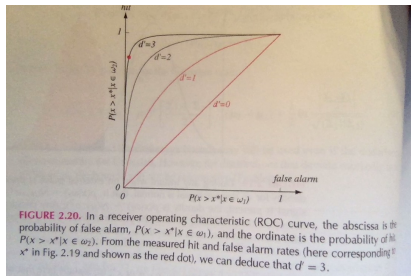
- *Sensitivity:*  
 $1 - P[\text{error} | \omega_2] = T_p / (T_p + F_n)$
- *Specificity:*  
 $1 - P[\text{error} | \omega_1] = T_n / (T_n + F_p)$

### Information Retrieval:

- *Precision:*  $T_p / (T_p + F_n)$
- *Recall:*  $T_p / (T_p + F_n) = ?$

		correct result /	classification
		C1	C2
obtained result / classification	C1	tp (true positive)	fp (false positive)
	C2	fn (false negative)	tn (true negative)

# Receiver Operating Characteristic Curves



- ▶ ROCs display Type I error rate (false alarm prob) vs...
- ▶ 1-Type II error rate (hit prob)
- ▶ shown for scalar Gaussians, common variance
- ▶ discriminability  

$$d' = |\mu_1 - \mu_2| / \sigma$$

## 2-Column Template