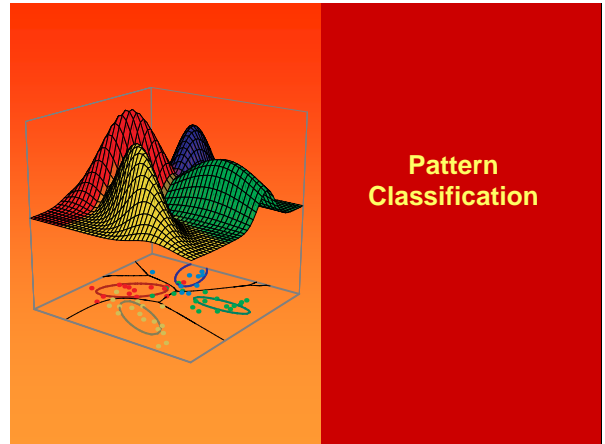




Pattern classification

Biometrics
CSE 190-a
Lecture 3

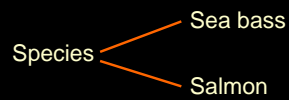
CSE 190-a, Fall 05



Pattern Classification

An Example

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

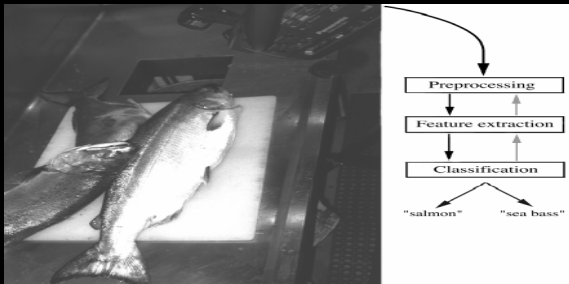


Pattern Classification, Chapter 1

• Problem Analysis

- Set up a camera and take some sample images to extract features
 - 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!

Pattern Classification, Chapter 1



Pattern Classification, Chapter 1

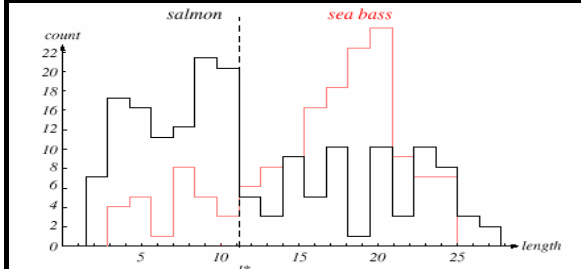
• Classification

- Select the length of the fish as a possible feature for discrimination

Pattern Classification, Chapter 1

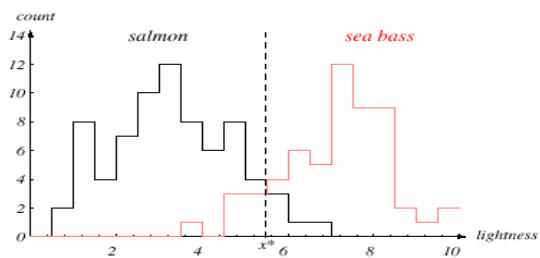
Classification

Select the length of the fish as a possible feature for discrimination

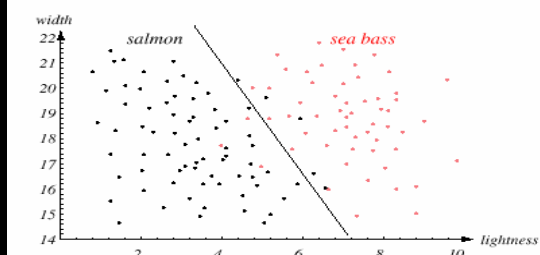
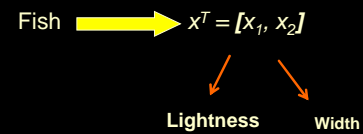


The **length** is a poor feature alone!

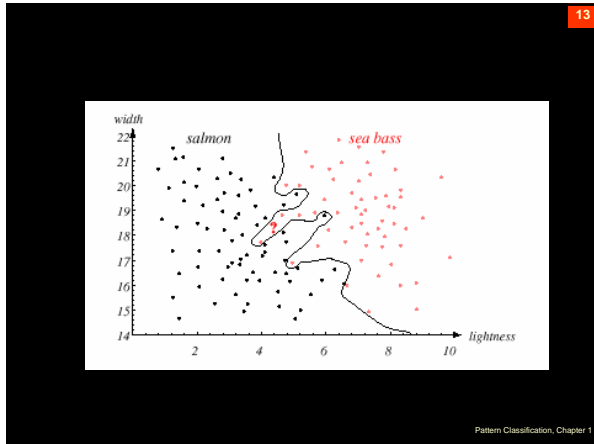
Select the **lightness** as a possible feature.



- Adopt the lightness and add the width of the fish



- We might add other features that are not correlated with the ones we already have. A precaution should be taken not to reduce the performance by adding such “noisy features”
- Ideally, the best decision boundary should be the one which provides an optimal performance such as in the following figure:



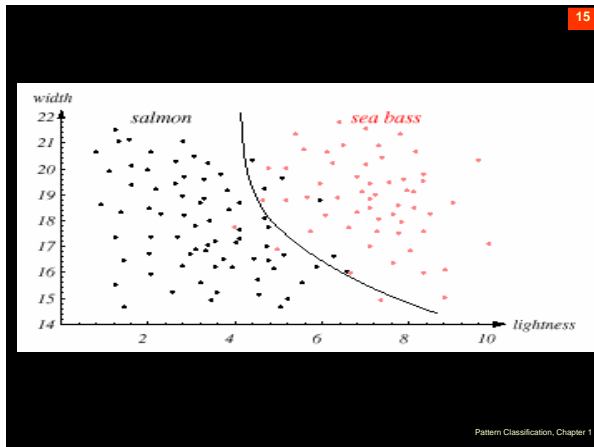
14

- However, our satisfaction is premature because the central aim of designing a classifier is to correctly classify novel input

↓

Issue of generalization!

Pattern Classification, Chapter 1



Bayesian Decision Theory
Continuous Features
(Sections 2.1-2.2)

17

Introduction

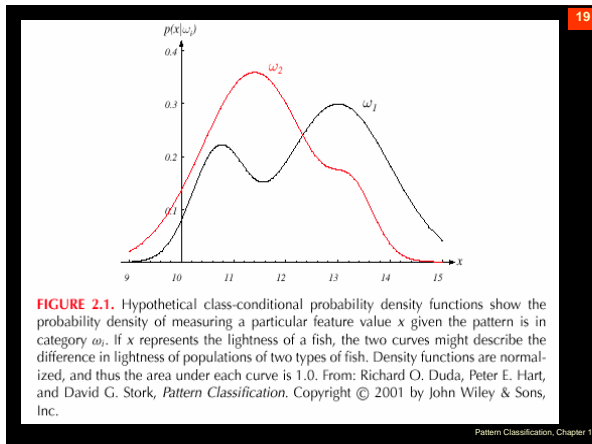
- The sea bass/salmon example
 - State of nature, prior
 - State of nature is a random variable
 - The catch of salmon and sea bass is equiprobable
 - $P(\omega_1), P(\omega_2)$ Prior probabilities
 - $P(\omega_1) = P(\omega_2)$ (uniform priors)
 - $P(\omega_1) + P(\omega_2) = 1$ (exclusivity and exhaustivity)

Pattern Classification, Chapter 1

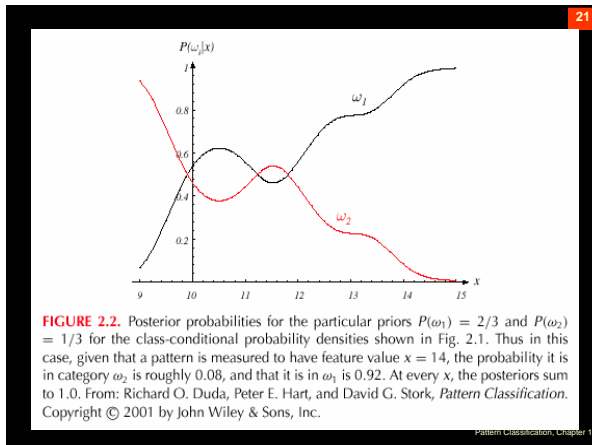
18

- Decision rule with only the prior information
 - Decide ω_1 if $P(\omega_1) > P(\omega_2)$ otherwise decide ω_2
- Use of the class-conditional information
 - $P(x | \omega_1)$ and $P(x | \omega_2)$ describe the difference in lightness between populations of sea-bass and salmon

Pattern Classification, Chapter 1



- Posterior, likelihood, evidence
 - $P(\omega_j | x) = (P(x | \omega_j) * P(\omega_j)) / P(x)$ **(BAYES RULE)**
 - In words, this can be said as:
Posterior = (Likelihood * Prior) / Evidence
 - Where in case of two categories
- $$P(x) = \sum_{j=1}^2 P(x | \omega_j) P(\omega_j)$$



-
- Intuitive decision rule given the posterior probabilities:
Given x :
if $P(\omega_1 | x) > P(\omega_2 | x)$ \rightarrow True state of nature = ω_1
if $P(\omega_1 | x) < P(\omega_2 | x)$ \rightarrow True state of nature = ω_2
- Why do this?: Whenever we observe a particular x , the probability of error is :
- $$P(\text{error} | x) = P(\omega_1 | x) \text{ if we decide } \omega_2$$
- $$P(\text{error} | x) = P(\omega_2 | x) \text{ if we decide } \omega_1$$

- Since decision rule is optimal for each feature value X , there is not better rule for all x .

- ### Bayesian Decision Theory – Continuous Features
- #### Generalization of the preceding ideas
- Use of more than one feature
 - Use more than two states of nature
 - Allowing actions and not only decide on the state of nature
 - Introduce a loss of function (more general than the probability of error)
 - Allowing actions other than classification primarily allows the possibility of rejection
 - Refusing to make a decision in close or bad cases!
 - Letting loss function state how costly each action taken is

Bayesian Decision Theory – Continuous Features 25

- Let \mathbf{X} be a vector of features.
- Let $\{\omega_1, \omega_2, \dots, \omega_c\}$ be the set of c states of nature (or "classes")
- Let $\{\alpha_1, \alpha_2, \dots, \alpha_a\}$ be the set of possible actions
- Let $\lambda(\alpha_i | \omega_j)$ be the loss for action α_i when the state of nature is ω_j

Pattern Classification, Chapter 1

What is the Expected Loss for action α_i ? 26

For any given x the expected loss is

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

$R(\alpha_i | x)$ is called the **Conditional Risk (or Expected Loss)**

Pattern Classification, Chapter 1

Overall risk 27

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

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

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

for $i = 1, \dots, a$

Pattern Classification, Chapter 1

Given a measured feature vector x , which action should we take? 28

Select the action α_i for which $R(\alpha_i | x)$ is minimum

\implies R is minimum and R in this case is called the **Bayes risk = best performance that can be achieved!**

Pattern Classification, Chapter 1

Two-Category Classification 29

α_1 : deciding ω_1

α_2 : deciding ω_2

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

loss incurred for deciding ω_j when the true state of nature is ω_i

Conditional risk:

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

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

Pattern Classification, Chapter 1

Our rule is the following: 30

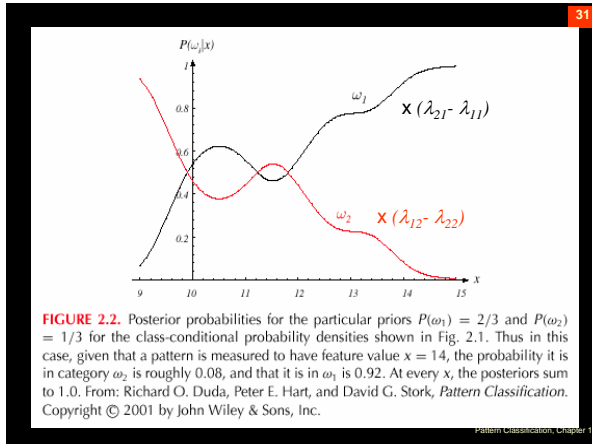
if $R(\alpha_1 | x) < R(\alpha_2 | x)$
 $\lambda_{11}P(\omega_1 | x) + \lambda_{12}P(\omega_2 | x) < \lambda_{21}P(\omega_1 | x) + \lambda_{22}P(\omega_2 | x)$
 action α_1 ; "decide ω_1 " is taken

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

Pattern Classification, Chapter 1



32

Two-Category Decision Theory: Chopping Machine

$\alpha_1 = \text{chop}$
 $\alpha_2 = \text{DO NOT chop}$
 $\omega_1 = \text{NO hand in machine}$
 $\omega_2 = \text{hand in machine}$

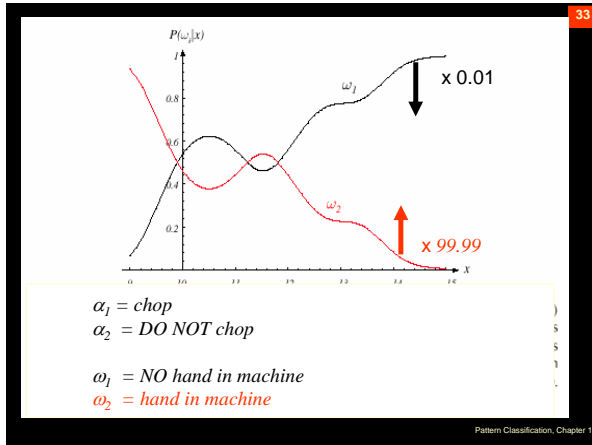
$\lambda_{11} = \lambda(\alpha_1 | \omega_1) = \$ 0.00$
 $\lambda_{12} = \lambda(\alpha_1 | \omega_2) = \$ 100.00$
 $\lambda_{21} = \lambda(\alpha_2 | \omega_1) = \$ 0.01$
 $\lambda_{22} = \lambda(\alpha_2 | \omega_2) = \$ 0.01$

Therefore our rule becomes ➔

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

$$0.01 P(x | \omega_1) P(\omega_1) > 99.99 P(x | \omega_2) P(\omega_2)$$

Pattern Classification, Chapter 1



34

Exercise to do at home!!

Select the optimal decision where:
 $\Omega = \{\omega_1, \omega_2\}$

$P(x | \omega_1) \Rightarrow N(2, 0.5)$ (Normal distribution)
 $P(x | \omega_2) \Rightarrow N(1.5, 0.2)$

$P(\omega_1) = 2/3$
 $P(\omega_2) = 1/3$

$$\lambda = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

Pattern Classification, Chapter 1

35

Minimum-Error-Rate Classification revisited

- Actions are decisions on classes
 If action α_i is taken and the true state of nature is ω_j then the decision is correct if $i = j$ and in error if $i \neq j$
- Seek a decision rule that minimizes the *probability of error* which is the *error rate*

Pattern Classification, Chapter 1

36

- Introduction of the zero-one loss function:

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

Therefore, the conditional risk for each action is:

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

$$= \sum_{j \neq i} P(\omega_j | x) = 1 - P(\omega_i | x) \quad \text{Average Prob. of Error}$$

Since $\sum P(\omega_j | x) = 1$

"The risk corresponding to this loss function is the average (or expected) probability error"

Pattern Classification, Chapter 1

- Minimize the risk requires maximize $P(\omega_i | x)$ (since $R(\alpha_i | x) = 1 - P(\omega_i | x)$)
- For Minimum error rate
 - Decide ω_i if $P(\omega_i | x) > P(\omega_j | x) \forall j \neq i$

Likelihood ratio:

The preceding rule is equivalent to the following rule:

$$\text{if } \frac{P(x | \omega_1)}{P(x | \omega_2)} > \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \cdot \frac{P(\omega_2)}{P(\omega_1)}$$

Then take action α_1 (decide ω_1)
 Otherwise take action α_2 (decide ω_2)

- Regions of decision and zero-one loss function, therefore:

$$\text{Let } \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \cdot \frac{P(\omega_2)}{P(\omega_1)} = \theta_\lambda \text{ then decide } \omega_1 \text{ if : } \frac{P(x | \omega_1)}{P(x | \omega_2)} > \theta_\lambda$$

- If λ is the zero-one loss function which means:

$$\lambda = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

then $\theta_\lambda = \frac{P(\omega_2)}{P(\omega_1)} = \theta_a$

if $\lambda = \begin{pmatrix} 0 & 2 \\ 1 & 0 \end{pmatrix}$ then $\theta_\lambda = \frac{2P(\omega_2)}{P(\omega_1)} = \theta_b$

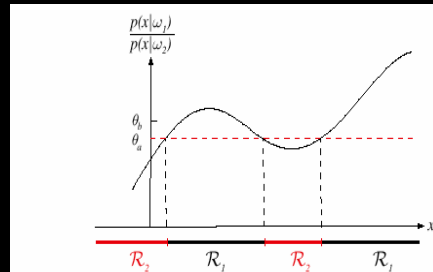


FIGURE 2.3. The likelihood ratio $p(x|\omega_1)/p(x|\omega_2)$ for the distributions shown in Fig. 2.1. If we employ a zero-one or classification loss, our decision boundaries are determined by the threshold θ_λ . If our loss function penalizes miscategorizing ω_2 as ω_1 patterns more than the converse, we get the larger threshold θ_b , and hence R_1 becomes smaller. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Classifiers, Discriminant Functions and Decision Surfaces

- Discriminant Functions: A generalization
- The multi-category case
 - Consider a set of c discriminant functions $g_i(x), i = 1, \dots, c$
 - The classifier assigns a feature vector x to class ω_i if:

$$g_i(x) > g_j(x) \forall j \neq i$$
 - Designing a classifier amounts to specifying the $g_i(x)$

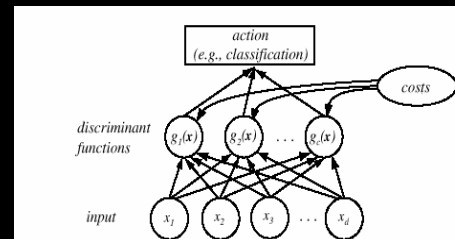


FIGURE 2.5. The functional structure of a general statistical pattern classifier which includes d inputs and c discriminant functions $g_i(x)$. A subsequent step determines which of the discriminant values is the maximum, and categorizes the input pattern accordingly. The arrows show the direction of the flow of information, though frequently the arrows are omitted when the direction of flow is self-evident. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

43

Decision Regions

- Feature space divided into c decision region
 - if $g_i(x) > g_j(x) \forall j \neq i$ then x is in \mathcal{R}_i
 - (\mathcal{R}_i means assign x to ω_i)

Decision surfaces

$\{x: \exists i, j, g_i(x) = g_j(x)\}$

Pattern Classification, Chapter 1

44

- Bayes Risk as discriminant function.**
 - Let $g_i(x) = -R(\omega_i | x)$
(max. discriminant corresponds to min. risk!)
- For the minimum error rate, discriminant function is:
 - $g_i(x) = P(\omega_i | x)$
 - (max. discrimination corresponds to max. posterior!)
 - $g_i(x) = P(x | \omega_i) P(\omega_i)$
- Any function $F(r)$ which is monotonic over $r > 0$ when applied to a set of discriminant functions, yields new discriminant function with the same decision regions/boundaries.

$g_i(x) = \ln P(x | \omega_i) + \ln P(\omega_i)$
(ln: natural logarithm!)

We'll see this form with Normal distributions

Pattern Classification, Chapter 1

47

On to higher dimensions!

Pattern Classification, Chapter 1

48

The Normal Density

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

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

Where:
 μ = mean (or expected value) of x
 σ^2 = expected squared deviation or variance

Pattern Classification, Chapter 1

49

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.

Pattern Classification, Chapter 1

50

Multivariate density

- Multivariate normal density in d dimensions is:

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

where:
 $x = (x_1, x_2, \dots, x_d)^t$ (t stands for the transpose vector form)
 $\mu = (\mu_1, \mu_2, \dots, \mu_d)^t$ mean vector
 $\Sigma = d$ by d covariance matrix
 $|\Sigma|$ and Σ^{-1} are determinant and inverse respectively

Pattern Classification, Chapter 1