

# Pattern Classification

# Chapter 3:

## Maximum-Likelihood & Bayesian Parameter Estimation (3.1,3.2)

- ▶ Introduction
- ▶ Maximum-Likelihood Estimation
  - The Gaussian Case 1: unknown  $\mu$
  - The Gaussian Case 2: unknown  $\mu$  and  $\sigma$
  - Bias



# 3.1 Introduction

- ▶ Data availability in a Bayesian framework
  - We could design an optimal classifier if we knew:
    - ▶  $P(\omega_i)$  (priors)
    - ▶  $P(x | \omega_i)$  (class-conditional densities)—Unknown parameters
    - ▶ Unfortunately, we rarely have both complete information!
- ▶ Design a classifier from a training sample
  - No problem with the estimation of prior probabilities
  - Samples are often too few for the estimation of class-conditional densities
  - Complexity for large dimension of feature space

- To simplify above problem

- ▶ Normality of  $P(x \mid \omega_i)$

- ▶  $P(x \mid \omega_i) \sim N(\mu_i, \Sigma_i)$ : Characterized by 2 parameters

- ▶ The problem is changed from estimating  $P(x \mid \omega_i)$  to estimating  $\mu_i, \Sigma_i$

- ▶ Estimation techniques

- Maximum-Likelihood (ML) and the Bayesian estimations

- Results are nearly identical, but the approaches are conceptually different



- Parameters in ML estimation are fixed but unknown!

Best parameters are obtained by maximizing the probability of obtaining the samples observed



- Bayesian methods view the parameters as random variables having some known prior distribution. Training data allow us to convert a distribution on this variable into a posterior probability density

In either approach, we use  $P(\omega_i | x)$  for our classification rule!



# 3.2 Maximum-Likelihood Estimation

## ► M-L Estimation

- Has good convergence properties as the sample size increases
- Simpler than any other alternative techniques

## ► General principle

- Assume we have  $c$  classes and

$$p(x | \omega_j) \sim N(\mu_j, \Sigma_j)$$

$$p(x | \omega_j) \equiv p(x | \omega_j, \theta_j) \text{ where:}$$

$$\theta_j = (\mu_j, \Sigma_j) = (\mu_j^1, \mu_j^2, \dots, \sigma_j^{11}, \sigma_j^{22} \dots)$$





- Use the information provided by the training samples  $D = (D_1, D_2, \dots, D_c)$  to estimate

$\theta = (\theta_1, \theta_2, \dots, \theta_c)$ , each  $\theta_i$  ( $i = 1, 2, \dots, c$ ) is associated with each category.

assume  $D_i$  give no information about  $\theta_j$  if  $i \neq j$

So we use each class separately to simplify our notation

Suppose that  $D$  contains  $n$  samples,  $x_1, x_2, \dots, x_n$

$$p(D | \theta) = \prod_{k=1}^n p(x_k | w, \theta) = \prod_{k=1}^n p(x_k | \theta) = F(\theta)$$

$p(D | \theta)$  is called the likelihood of  $\theta$



- ML estimation of  $\theta$  is, by definition, the value  $\hat{\theta}$  that maximizes  $p(D | \theta)$

“It is the value of  $\theta$  that best agrees with the actually observed training sample”

# ► ML Problem Statement

- Let  $D = \{x_1, x_2, \dots, x_n\}$

$$p(x_1, \dots, x_n \mid \theta) = \prod_{k=1}^n p(x_k \mid \theta); \quad |D| = n$$

Our goal is to determine  $\hat{\theta}$  (value of  $\theta$  that makes this sample the most representative!)





$$|D| = n$$

$$N(\mu_j, \Sigma_j) = P(x_j | \omega_c)$$

$x_1$   $x_2$   $x_n$

$$P(x_j | \omega_k)$$

$$P(x_j | \omega_1)$$

$D_c$

$x_{10}$

$x_{11}$

$x_{20}$

$D_k$

$x_8$

$x_1$

$x_9$

$D_1$

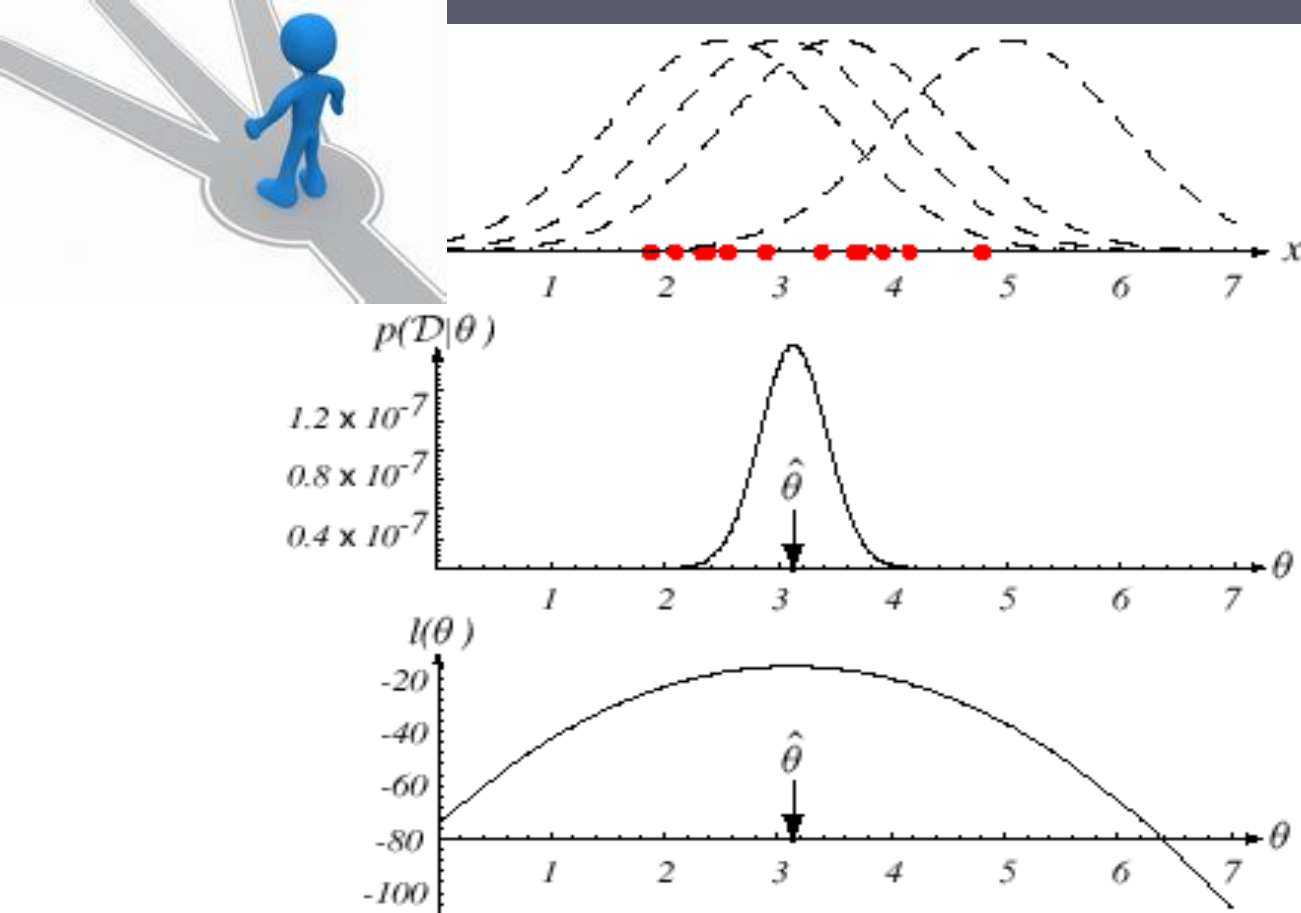


$$\theta = (\theta_1, \theta_2, \dots, \theta_c)$$

Problem: find  $\hat{\theta}$  such that:

$$\begin{aligned} \text{Max}_{\theta} P(D \mid \theta) &= \text{Max} P(x_1, \dots, x_n \mid \theta) \\ &= \text{Max} \prod_{k=1}^n P(x_k \mid \theta) \end{aligned}$$





**FIGURE 3.1.** The top graph shows several training points in one dimension, known or assumed to be drawn from a Gaussian of a particular variance, but unknown mean. Four of the infinite number of candidate source distributions are shown in dashed lines. The middle figure shows the likelihood  $p(\mathcal{D}|\theta)$  as a function of the mean. If we had a very large number of training points, this likelihood would be very narrow. The value that maximizes the likelihood is marked  $\hat{\theta}$ ; it also maximizes the logarithm of the likelihood—that is, the log-likelihood  $l(\theta)$ , shown at the bottom. Note that even though they look similar, the likelihood  $p(\mathcal{D}|\theta)$  is shown as a function of  $\theta$  whereas the conditional density  $p(x|\theta)$  is shown as a function of  $x$ . Furthermore, as a function of  $\theta$ , the likelihood  $p(\mathcal{D}|\theta)$  is not a probability density function and its area has no significance. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

## ■ Optimal estimation

- Let  $\theta = (\theta_1, \theta_2, \dots, \theta_p)^t$  and let  $\nabla_{\theta}$  be the gradient operator

$$\nabla_{\theta} = \left[ \frac{\partial}{\partial \theta_1}, \frac{\partial}{\partial \theta_2}, \dots, \frac{\partial}{\partial \theta_p} \right]^t$$

- We define  $l(\theta)$  as the log-likelihood function

$$l(\theta) = \ln p(D \mid \theta)$$

- New problem statement:

determine  $\theta$  that maximizes the log-likelihood

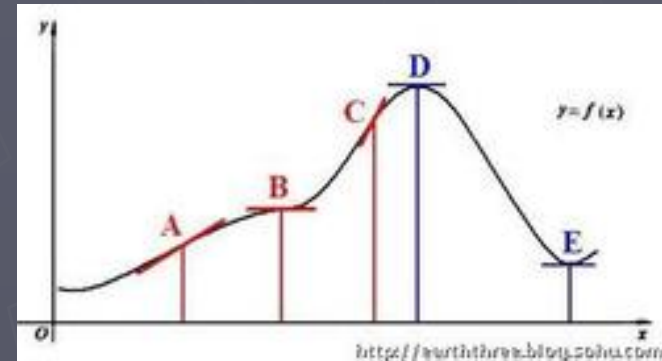
$$\hat{\theta} = \arg \max_{\theta} l(\theta)$$



- Set of necessary conditions for an optimum is:

$$(\nabla_{\theta} l = \sum_{k=1}^{k=n} \nabla_{\theta} \ln P(x_k | \theta))$$

$$\nabla_{\theta} l = 0$$



- Global maximum, local maximum or minimum, inflection point
- For reference: MAP estimators (Max a posteriori)

$$l(\theta) p(\theta)$$

## ► Example of a specific case 1: unknown $\mu$

- $p(x_i | \mu) \sim N(\mu, \Sigma)$

(Samples are drawn from a multivariate normal population)

$$\ln p(x_k | \mu) = -\frac{1}{2} \ln[(2\pi)^d |\Sigma|] - \frac{1}{2} (x_k - \mu)^t (\Sigma)^{-1} (x_k - \mu)$$

and  $\nabla_{\mu} \ln p(x_k | \mu) = (\Sigma)^{-1} (x_k - \mu)$

- The ML estimate for  $\mu$  must satisfy:

$$\sum_{k=1}^{k=n} \Sigma^{-1} (\mathbf{x}_k - \hat{\mu}) = \mathbf{0}$$



- Multiplying by  $\Sigma$  and rearranging, we obtain:

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^{k=n} x_k$$



Just the arithmetic average of the samples of the training samples!

## Conclusion:

If  $P(x_k | \omega_j)$  ( $j = 1, 2, \dots, c$ ) is supposed to be Gaussian in a  $d$ -dimensional feature space; then we can estimate the vector

$\theta = (\theta_1, \theta_2, \dots, \theta_c)^t$  and perform an optimal classification!

## ► Example of a specific case 2

- Gaussian Case: *unknown  $\mu$  and  $\sigma$*

$$\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$$

$$l = \ln P(x_k | \theta) = -\frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2} (x_k - \theta_1)^2$$

$$\nabla_{\theta} l = \begin{pmatrix} \frac{\partial}{\partial \theta_1} (\ln P(x_k | \theta)) \\ \frac{\partial}{\partial \theta_2} (\ln P(x_k | \theta)) \end{pmatrix} = 0$$

$$\begin{cases} \frac{1}{\theta_2} (x_k - \theta_1) = 0 \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} = 0 \end{cases}$$

Summation:

$$\begin{cases} \sum_{k=1}^{k=n} \frac{1}{\hat{\theta}_2} (x_k - \hat{\theta}_1) = 0 & (1) \\ -\sum_{k=1}^{k=n} \frac{1}{\hat{\theta}_2} + \sum_{k=1}^{k=n} \frac{(x_k - \hat{\theta}_1)^2}{\hat{\theta}_2^2} = 0 & (2) \end{cases}$$

Combining (1) and (2), one obtains:

$$\mu = \sum_{k=1}^{k=n} \frac{x_k}{n} \quad ; \quad \sigma^2 = \frac{\sum_{k=1}^{k=n} (x_k - \mu)^2}{n}$$



## ► Bias

- ML estimate for  $\sigma^2$  is biased

$$E[\hat{\sigma}^2] = E\left[\frac{1}{n} \sum (x_i - \bar{x})^2\right] = \frac{n-1}{n} \cdot \sigma^2 \neq \sigma^2$$

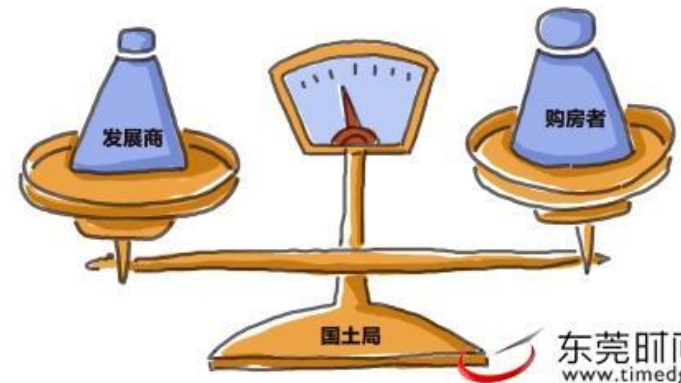
- An elementary unbiased estimator for  $\Sigma$  is:

$$C = \frac{1}{n-1} \sum_{k=1}^{k=n} (x_k - \hat{\mu})(x_k - \hat{\mu})^t$$

Sample covariance matrix

- ML estimate for  $\Sigma$  is biased

$$\hat{\Sigma} = \frac{n-1}{n} C$$



- Absolutely unbiased, asymptotically unbiased
- Prove ML estimate for  $\sigma^2$  is biased

$$E[x^2] = D[x] + E[x]^2$$

$$E\left[\sum_{i=1}^n x_i^2\right] = n(\sigma^2 + \mu^2)$$

$$E[\bar{x}^2] = D[\bar{x}] + E(\bar{x})^2 = \frac{1}{n}\sigma^2 + \mu^2$$

$$\begin{aligned} E\left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\right] &= \frac{1}{n} E\left[\sum_{i=1}^n x_i^2 - n\bar{x}^2\right] \\ &= \frac{1}{n} \left[ n(\sigma^2 + \mu^2) - n\left(\frac{1}{n}\sigma^2 + \mu^2\right) \right] = \frac{n-1}{n} \sigma^2 \end{aligned}$$