

Chapter 3:

Maximum-Likelihood & Bayesian Parameter Estimation (part 1)

- | Introduction
- | Maximum-Likelihood Estimation
 - | Example of a Specific Case
 - | The Gaussian Case: unknown μ and σ
 - | Bias
- | Appendix: ML Problem Statement

I Introduction

I Data availability in a Bayesian framework

- I We could design an optimal classifier if we knew:
 - I $P(\omega_i)$ (priors)
 - I $P(x | \omega_i)$ (class-conditional densities)

Unfortunately, we rarely have this complete information!

I Design a classifier from a training sample

- I No problem with prior estimation
- I Samples are often too small for class-conditional estimation (large dimension of feature space!)

- | A priori information about the problem
- | Normality of $P(x | \omega_i)$

$$P(x | \omega_i) \sim N(\mu_i, \Sigma_i)$$

- | Characterized by 2 parameters
- | Estimation techniques
 - | Maximum-Likelihood (ML) and the Bayesian estimations
 - | Results are nearly identical, but the approaches are 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 distribution
- | In either approach, we use $P(\omega_i | x)$ for our classification rule!

I Maximum-Likelihood Estimation

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

I General principle

- I Assume we have c classes and

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

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

$$q = (m_j, S_j) = (m_j^1, m_j^2, \dots, S_j^{11}, S_j^{22}, \text{cov}(x_j^m, x_j^n) \dots)$$

- Use the information provided by the training samples to estimate $\theta = (\theta_1, \theta_2, \dots, \theta_c)$, each θ_i ($i = 1, 2, \dots, c$) is associated with each category
- Suppose that D contains n samples, x_1, x_2, \dots, x_n (samples drawn independently)

$$P(D | q) = \prod_{k=1}^n P(x_k | q) = F(q)$$

$P(D | q)$ is called the likelihood of q w.r.t. the set of samples)

- ML estimate of θ is, by definition the value that \hat{q} maximizes $P(D | \theta)$
 “It is the value of θ that best agrees with the actually observed training sample”

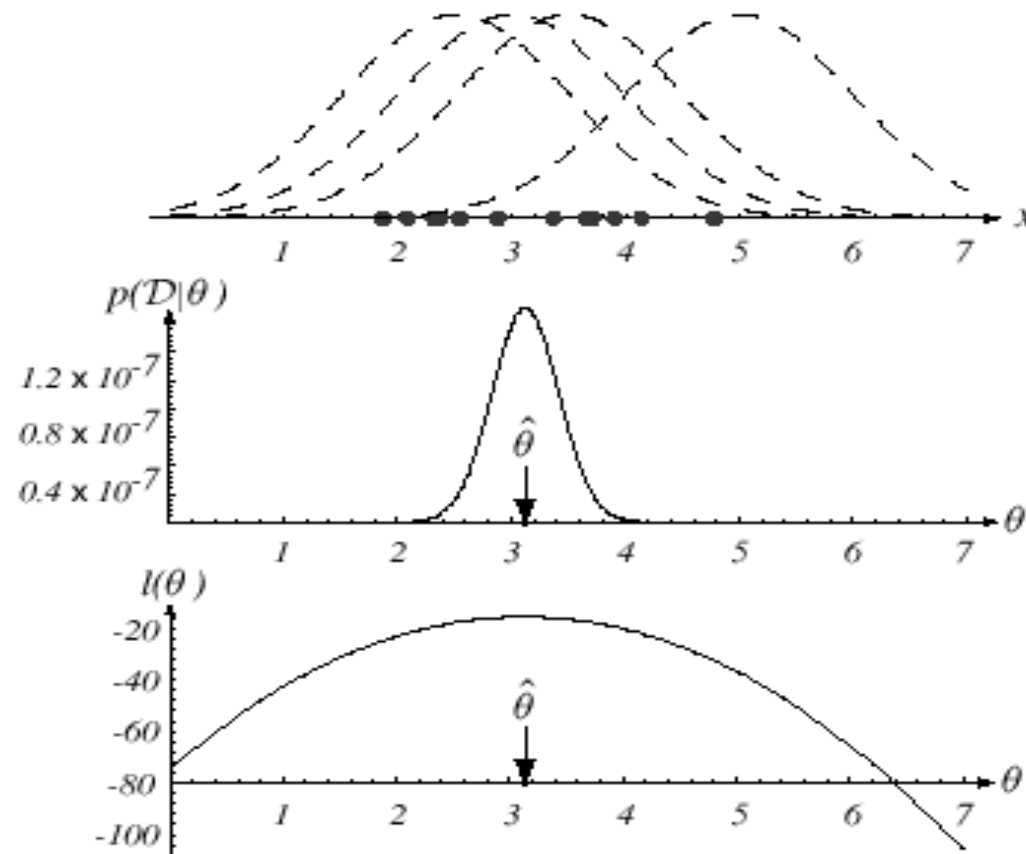


FIGURE 3.1. The top graph shows several training points in one dimension, known or assumed to be drawn from a Gaussian of a particular variance, but unknown mean. Four of the infinite number of candidate source distributions are shown in dashed lines. The middle figure shows the likelihood $p(\mathcal{D}|\theta)$ 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 θ whereas the conditional density $p(x|\theta)$ is shown as a function of x . Furthermore, as a function of θ , 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 ∇_θ be the gradient operator

$$\tilde{N}_q = \left(\frac{\partial \ell}{\partial \theta_1}, \frac{\partial \ell}{\partial \theta_2}, \dots, \frac{\partial \ell}{\partial \theta_p} \right)^t$$

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

$$\ell(\theta) = \ln P(D | \theta)$$

- New problem statement:
determine θ that maximizes the log-likelihood

$$\hat{q} = \arg \max_q \ell(q)$$

Set of necessary conditions for an optimum is:

$$(\nabla_q \mathbf{l} = \sum_{k=1}^{k=n} \nabla_q \ln P(x_k | \mathbf{q}))$$

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

I Example of a specific case: unknown μ

I $P(x_i | \mu) \sim N(\mu, \Sigma)$

(Samples are drawn from a multivariate normal population)

$$\ln P(\mathbf{x}_k | \mathbf{m}) = -\frac{1}{2} \ln[(2\pi)^d |S|] - \frac{1}{2} (\mathbf{x}_k - \mathbf{m})^t S^{-1} (\mathbf{x}_k - \mathbf{m})$$

$$\text{and } \tilde{N}_{qm} \ln P(\mathbf{x}_k | \mathbf{m}) = S^{-1} (\mathbf{x}_k - \mathbf{m})$$

$\theta = \mu$ therefore:

- The ML estimate for μ must satisfy:

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

- Multiplying by Σ and rearranging, we obtain:

$$\hat{\mathbf{m}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

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

Conclusion:

If $P(\mathbf{x}_k | \omega_j)$ ($j = 1, 2, \dots, c$) 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!

I ML Estimation:

- I Gaussian Case: *unknown m and s*
(one dimension) $q = (q_1, q_2) = (\mu, \sigma^2)$

$$\mathbf{l} = \ln P(x_k | q) = -\frac{1}{2} \ln(2\pi q_2) - \frac{1}{2q_2} (x_k - q_1)^2$$

$$\nabla_{\theta} \mathbf{l} = \begin{pmatrix} \frac{\partial}{\partial q_1} (\ln P(x_k | q)) \\ \frac{\partial}{\partial q_2} (\ln P(x_k | q)) \end{pmatrix} = 0$$

$$\begin{cases} \frac{1}{q_2} (x_k - q_1) = 0 \\ -\frac{1}{2q_2} + \frac{(x_k - q_1)^2}{2q_2^2} = 0 \end{cases}$$

Summation:

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

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

$$\bar{\mathbf{m}} = \frac{\sum_{k=1}^{k=n} \mathbf{x}_k}{n} \quad ; \quad S^2 = \frac{\sum_{k=1}^{k=n} (\mathbf{x}_k - \bar{\mathbf{m}})^2}{n}$$

I Bias

- I ML estimate for σ^2 is biased

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

- I An elementary unbiased estimator for Σ is:

$$\mathbf{C} = \frac{1}{n-1} \sum_{k=1}^n (\mathbf{x}_k - \bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}})^t$$

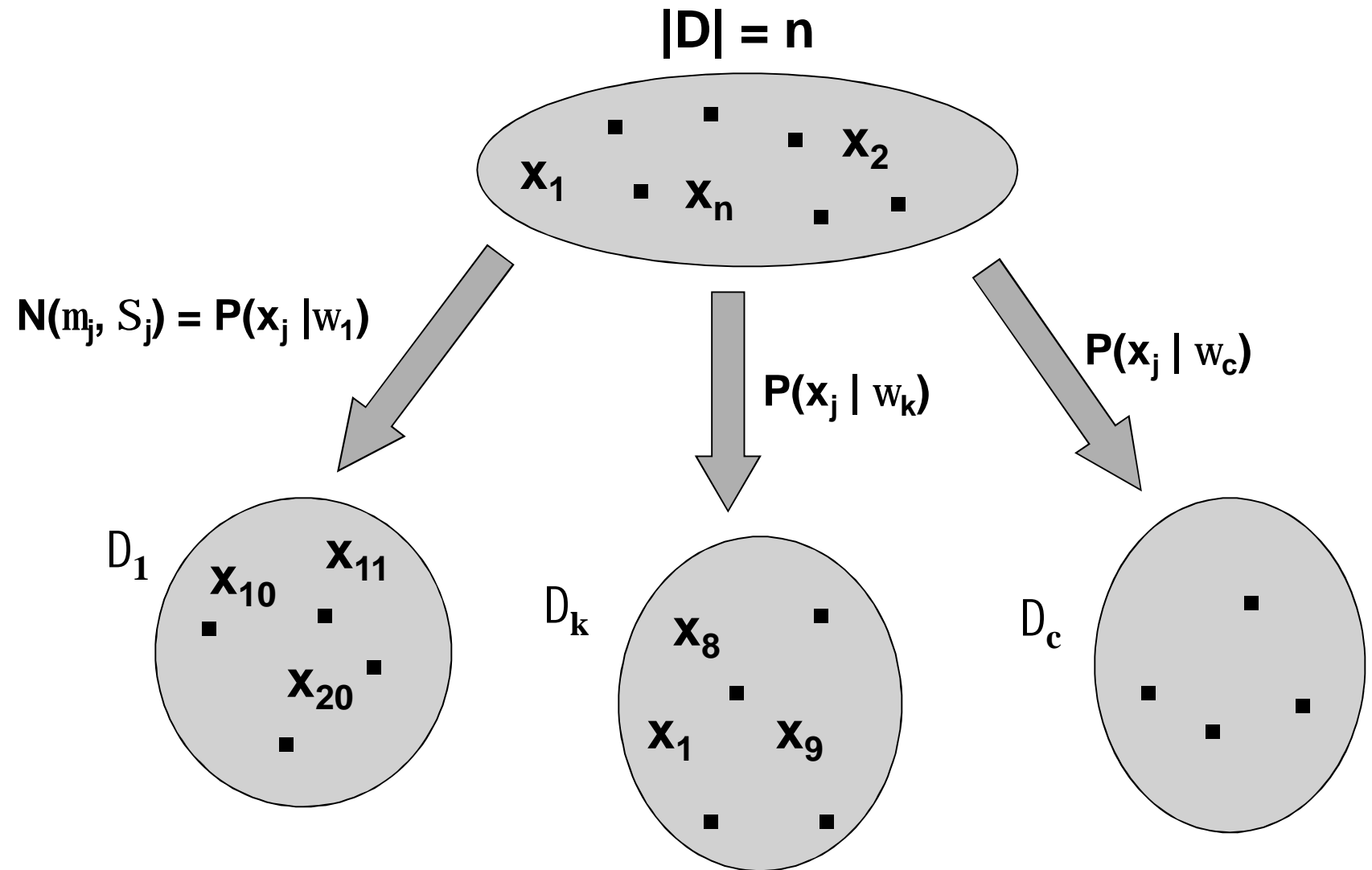
Sample covariance matrix

I Appendix: ML Problem Statement

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

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

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



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

Problem: find \hat{q} such that:

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

summary

