

Chapter 3 (part 1): Maximum-Likelihood & Bayesian Parameter Estimation

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



All materials used in this course were taken from the textbook "*Pattern Classification*" by Duda et al., John Wiley & Sons, 2001 with the permission of the authors and the publisher

✿ 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)

Unfortunately, we rarely have this complete information!

✿ Design a classifier from a training sample

- ✿ No problem with prior estimation
- ✿ 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!

✱ Maximum-Likelihood Estimation

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

✱ General principle

- ✱ Assume we have c classes and

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

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

$$\theta = (\mu_j, \Sigma_j) = (\mu_j^1, \mu_j^2, \dots, \sigma_j^{11}, \sigma_j^{22}, \text{cov}(x_j^m, x_j^n) \dots)$$

2

- ✱ 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

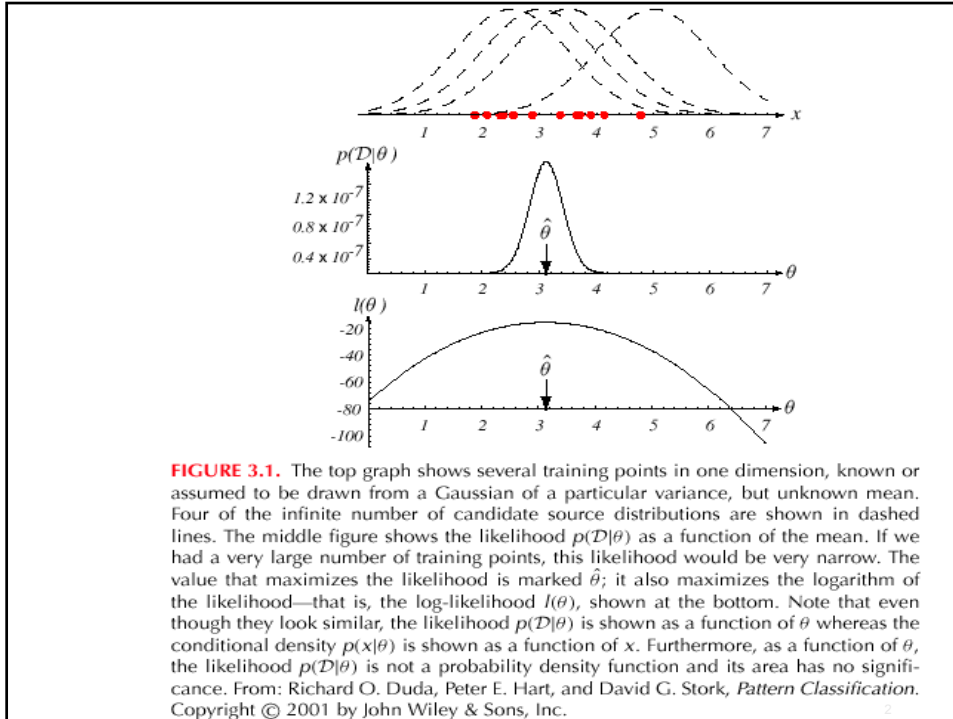
$$P(D | \theta) = \prod_{k=1}^{k=n} P(x_k | \theta) = F(\theta)$$

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

- ✱ ML estimate of θ is, by definition the value that $\hat{\theta}$ maximizes $P(D | \theta)$

“It is the value of θ that best agrees with the actually observed training sample”

2



✱ **Optimal estimation**

- Let $\theta = (\theta_1, \theta_2, \dots, \theta_p)^t$ and let ∇_{θ} 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

$$l(\theta) = \ln P(\mathcal{D} | \theta)$$

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

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

Set of necessary conditions for an optimum is:

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

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

2

✱ Example of a specific case: unknown μ

✱ $P(x_i | \mu) \sim N(\mu, \Sigma)$

(Samples are drawn from a multivariate normal population)

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

and $\nabla_{\theta\mu} \ln P(\mathbf{x}_k | \mu) = \Sigma^{-1} (\mathbf{x}_k - \mu)$

$\theta = \mu$ therefore:

• The ML estimate for μ must satisfy:

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

2

- Multiplying by Σ and rearranging, we obtain:

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

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

2

✱ ML Estimation:

✱ Gaussian Case: *unknown μ and σ*

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

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

$$\nabla_{\theta} l = \begin{pmatrix} \frac{\sigma}{\sigma\theta_1} (\ln P(\mathbf{x}_k | \theta)) \\ \frac{\sigma}{\sigma\theta_2} (\ln P(\mathbf{x}_k | \theta)) \end{pmatrix} = 0$$

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

2

Summation:

$$\begin{cases} \sum_{k=1}^{k=n} \frac{1}{\hat{\theta}_2} (\mathbf{x}_k - \theta_1) = 0 & (1) \end{cases}$$

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

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

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

✱ Bias

- ✱ ML estimate for σ^2 is biased

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

- ✱ An elementary unbiased estimator for Σ is:

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

Sample covariance matrix

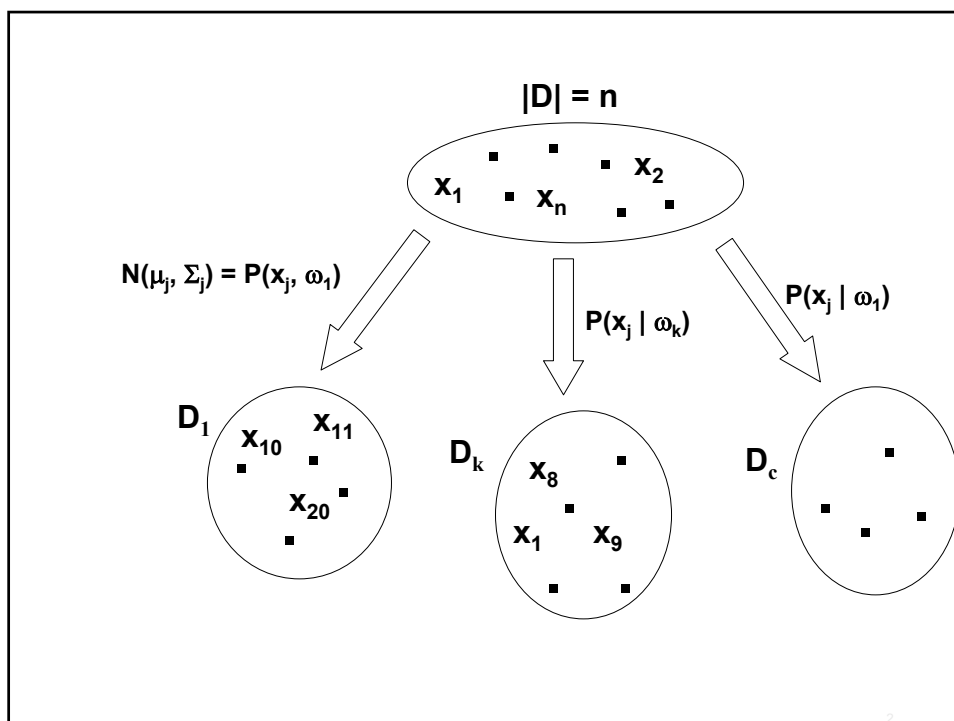
✱ Appendix: ML Problem Statement

✱ Let $D = \{x_1, x_2, \dots, x_n\}$

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

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

2



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

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

$$\begin{aligned} \text{Max}_{\theta} P(\mathbf{D} | \theta) &= \text{Max} P(\mathbf{x}_1, \dots, \mathbf{x}_n | \theta) \\ &= \text{Max} \prod_{k=1}^n P(\mathbf{x}_k | \theta) \end{aligned}$$

2