

Parameter Estimation

- Data availability in a Bayesian framework
 - We could design an optimal classifier if we knew:
 - $P(\omega_i)$ (priors)
 - $p(\mathbf{x} | \omega_i)$ (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 $p(\mathbf{x} \mid \omega_i)$

$$p(\mathbf{x} \mid \omega_i) \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

Characterized by 2 parameters

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
- Parameters are chosen in a way that they best support/ describe the training data.

- Bayesian methods view the parameters as random variables having some known distribution
- Observation of the samples converts this to a posterior density, thereby revising our opinion about the true values of the parameters.
- In the Bayesian case, we shall see that a typical effect of observing additional samples is to sharpen the a posteriori density function, causing it to peak near the true values of the parameters. This phenomenon is known as *Bayesian learning*.

- In either approach, we use $P(\omega_i | \mathbf{x})$ for our classification rule!
- Bayes Theorem is the key...!

Maximum-Likelihood Estimation

- 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
$$p(\mathbf{x} \mid \omega_j) \sim N(\mu_j, \Sigma_j)$$
$$p(\mathbf{x} \mid \omega_j) \equiv P(\mathbf{x} \mid \omega_j, \theta_j) \text{ where:}$$

- 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, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$

$$p(\mathcal{D}|\theta) = \prod_{k=1}^n p(\mathbf{x}_k|\theta).$$

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

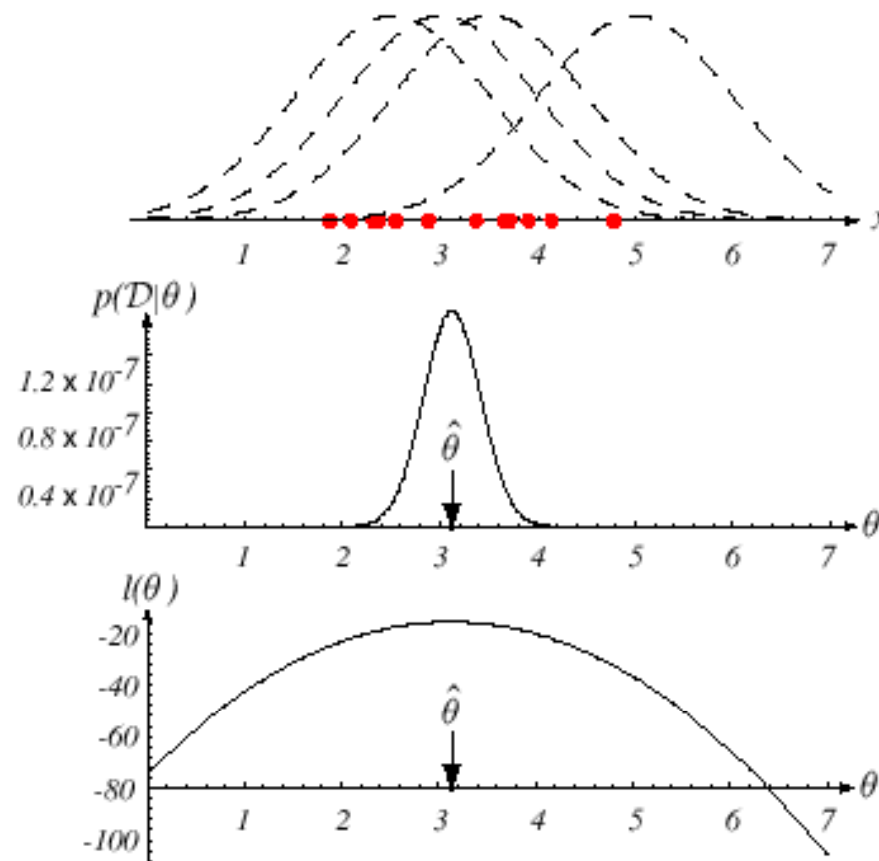


FIGURE 3.1. The top graph shows several training points in one dimension, known or assumed to be drawn from a Gaussian of a particular variance, but unknown mean. Four of the infinite number of candidate source distributions are shown in dashed lines. The middle figure shows the likelihood $p(\mathcal{D}|\theta)$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

Maximum likelihood estimation

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

$$\nabla_{\theta} \equiv \begin{bmatrix} \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_p} \end{bmatrix}$$

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

$$l(\theta) \equiv \ln p(\mathcal{D}|\theta)$$

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

Maximum likelihood estimation

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

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

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

Set of necessary conditions :

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

Gaussian Case : unknown μ , known Σ

$$p(\mathbf{x} \mid \mu) \sim N(\mu, \Sigma)$$

(Samples are drawn from a multivariate normal population)

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

and

$$\nabla_{\theta} \ln p(\mathbf{x}_k \mid \mu) = \Sigma^{-1} (\mathbf{x}_k - \mu).$$

$\theta = \mu$ therefore:

- The ML estimate for μ must satisfy:

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

- 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!

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

Gaussian Case: *unknown μ and σ*

– Gaussian Case: *unknown μ and σ*

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

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

$$\nabla_{\theta} l = \nabla_{\theta} \ln p(x_k|\theta) = \begin{bmatrix} \frac{1}{\theta_2}(x_k - \theta_1) \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} \end{bmatrix}.$$

Solving :

$$\sum_{k=1}^n \frac{1}{\hat{\theta}_2} (x_k - \hat{\theta}_1) = 0$$

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

Combining above equations, one obtains:

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

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

The Gaussian Case: Unknown μ and Σ

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

and

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

– Let $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$

$$p(x_1, \dots, x_n \mid \theta) = \prod p(\mathbf{x}_k \mid \theta); \quad |D| = n$$

Our goal is to determine (value of θ that makes this sample the most representative!)

