

Maximum-Likelihood & Bayesian Parameter Estimation

Maximum Likelihood *Versus* Bayesian Parameter Estimation

Optimal classifier can be designed knowing $P(\omega_i)$ and $p(x | \omega_i)$

Obtain them from training samples assuming known forms of pdfs, e.g., $p(x | \omega_i) \sim N(\mu_i, \Sigma_i)$ has 2 parameters

- Estimation techniques
 - Maximum-Likelihood (ML)
 - Find parameters that maximize probability of observations
 - Bayesian estimation
 - Parameters are random variables with known prior distribution sharpened by observations
 - Results nearly identical, but approaches are different

Maximum Likelihood Parameter Estimation

- Parameters are fixed but unknown
- Best parameters obtained by maximizing probability of obtaining samples observed
- Has good convergence properties as sample size increases
- Simpler than any other alternative techniques

- The General principle
Assume c classes and

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

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

$$\theta = (\mu_j, \Sigma_j)$$

Maximum Likelihood Estimation

- Use n training samples in a class to estimate θ
- If \mathcal{D} contains n independently drawn samples, x_1, x_2, \dots, x_n

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

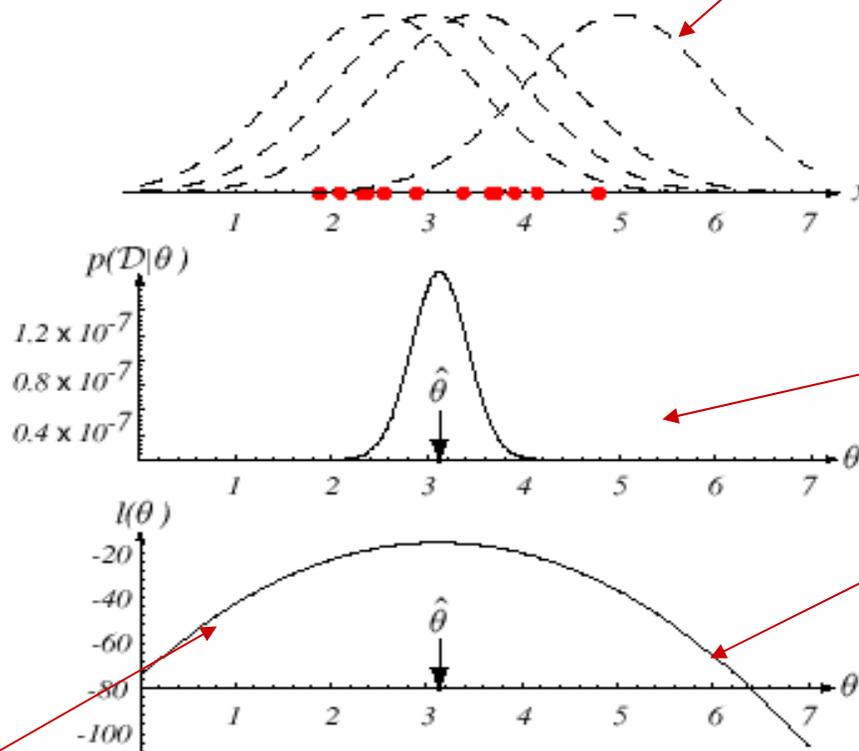
$p(\mathcal{D} | \theta)$ is called the likelihood of θ w.r.t. the set of samples

$\ell(\theta)$ is the log - likelihood of θ

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

One-Dimensional Example

Four of infinite
no. of source distributions



Likelihood as a
function of mean
(peaks at mean)
Would be sharp
peak with many
samples

Log-likelihood
function(also
peaks at mean)

Maximizing the log-likelihood function

- 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

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

Determine θ that maximizes the log-likelihood

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

Set of necessary conditions for an optimum is:

$$\nabla_{\theta} l = \sum_{k=1}^n \nabla_{\theta} \ln p(x_k | \theta) = 0$$

MLE: The Gaussian Case: unknown μ

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

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

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

$\theta = \mu$ therefore:

- The ML estimate for μ must satisfy:

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

- Multiplying by Σ

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

Just the sample mean

MLE: Gaussian Case- *unknown μ and σ*

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

$$\ln \mathbf{p}(\mathbf{x}_k | \theta) = -\frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2} (\mathbf{x}_k - \theta_1)^2$$
$$\nabla_{\theta} \mathbf{l} = \begin{pmatrix} \frac{\partial}{\partial \theta_1} (\ln \mathbf{p}(\mathbf{x}_k | \theta)) \\ \frac{\partial}{\partial \theta_2} (\ln \mathbf{p}(\mathbf{x}_k | \theta)) \end{pmatrix} = 0$$
$$\begin{cases} \frac{1}{\theta_2} (\mathbf{x}_k - \theta_1) = 0 \\ -\frac{1}{2\theta_2} + \frac{(\mathbf{x}_k - \theta_1)^2}{2\theta_2^2} = 0 \end{cases}$$

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

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

MLE Bias

- ML estimate for σ^2 is biased

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

- An elementary unbiased estimator for Σ is:

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

Sample covariance matrix