

Chapter 3

Maximum-Likelihood and Bayesian Parameter Estimation (3,4,5)

- Bayesian Estimation (BE)
- Bayesian Parameter Estimation: Gaussian Case
- Bayesian Parameter Estimation: General Estimation

3.3 Bayesian Estimation

- ▶ In MLE θ was supposed to be a fixed value
- ▶ In BE θ is a random variable
- ▶ The computation of posterior probabilities $P(\omega_i | x)$ lies at the heart of Bayesian classification
- ▶ Goal: compute $P(\omega_i | x, D)$

Given the sample D , Bayes formula can be written

$$P(\omega_i | x, D) = \frac{p(x | \omega_i, D) \cdot P(\omega_i | D)}{\sum_{j=1}^c p(x | \omega_j, D) \cdot P(\omega_j | D)}$$

Sample D \longrightarrow likelihood (conditional probability)

\downarrow
posterior probabilities

- To demonstrate the preceding equation, we use:

$$D = D_1 \cup D_2 \dots \cup D_c \quad x \in D_i \rightarrow x \text{ is } \omega_i$$

D_i has no influence on $p(x | \omega_j, D_j)$ if $i \neq j$

$P(\omega_i) = P(\omega_i | D)$ (Training sample provides this!)

Thus :

$$P(\omega_i | x, D) = \frac{P(x | \omega_i, D_i).P(\omega_i)}{\sum_{j=1}^c P(x | \omega_j, D_j).P(\omega_j)}$$

Further illustration

likelihood (conditional probability)

$$\begin{aligned} p(x | \omega_i, D) &= p(x) \cong p(x | D) \\ &= \int p(x, \theta | D) d\theta \\ &= \int p(x | \theta) p(\theta | D) d\theta \end{aligned}$$

posterior $p(\theta | D)$ **Key**

Unknown θ and known prior density $p(\theta)$

Description of the above illustration

► Parameter Distribution

- $p(x)$ is unknown, we assume it has a known parametric form $p(x | \theta)$, and value of parameter θ is unknown

- Known prior density $p(\theta)$
- Training data convert $p(\theta)$ to a posterior $p(\theta | D)$

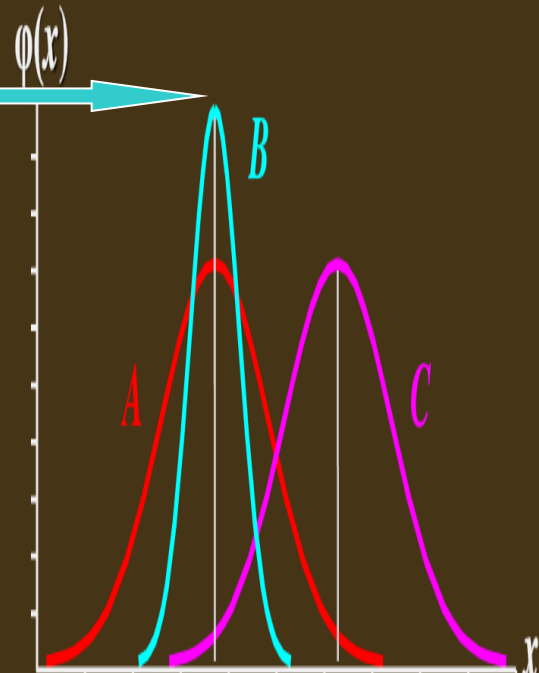
- Our path:

$$\begin{aligned} p(x | \omega_i, D) &= p(x) \cong p(x | D) \\ &= \int p(x, \theta | D) d\theta \\ &= \int p(x | \theta) p(\theta | D) d\theta \end{aligned}$$

- If $p(\theta | D)$ peaks very sharply about parameter $\hat{\theta}$ and $p(x | \theta)$ is smooth, and the tails of the integral are not important, then

$$p(x | D) = \int p(x | \theta) p(\theta | D) d\theta \quad \Rightarrow \quad p(x | D) \cong p(x | \hat{\theta})$$

e.g. the green curve



3.4 Bayesian Parameter Estimation: Gaussian Case

- ▶ **Goal:** Estimate θ using the a-posteriori density $P(\theta \mid D)$
- ▶ The univariate case: $P(\mu \mid D)$
 μ is the only unknown parameter

$$P(x \mid \mu) \sim N(\mu, \sigma^2)$$

$$P(\mu) \sim N(\mu_0, \sigma_0^2)$$

(μ_0 and σ_0 are known!)



$$P(\mu | D) = \frac{P(D | \mu) \cdot P(\mu)}{\int P(D | \mu) \cdot P(\mu) d\mu} \quad (1)$$

$$= \alpha \prod_{k=1}^{k=n} P(x_k | \mu) \cdot P(\mu)$$

- Reproducing density

$$P(\mu | D) \sim N(\mu_n, \sigma_n^2) \quad (2)$$

Identifying (1) and (2)
yields:

$$\mu_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \cdot \mu_0$$

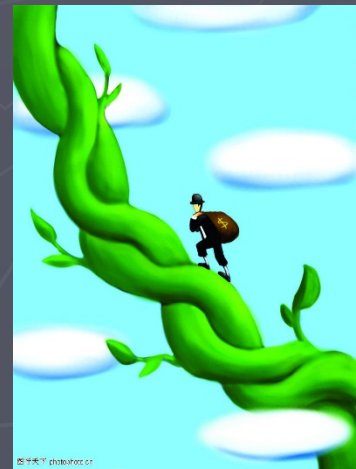
and $\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2} \quad \hat{\mu}_n = \frac{1}{n} (x_1 + \dots + x_n)$



► Understanding

- μ_n represents our best guess for μ after observing n samples
- σ_n^2 measures our uncertainty about this guess
- Add samples to decrease the uncertainty
- Bayse Learning: as n increase, $p(\mu | D)$ becomes more and more sharply peaked, approaching a Dirac delta function as n approaches infinity

$$\frac{\sigma_n^2}{\sigma_0^2}$$



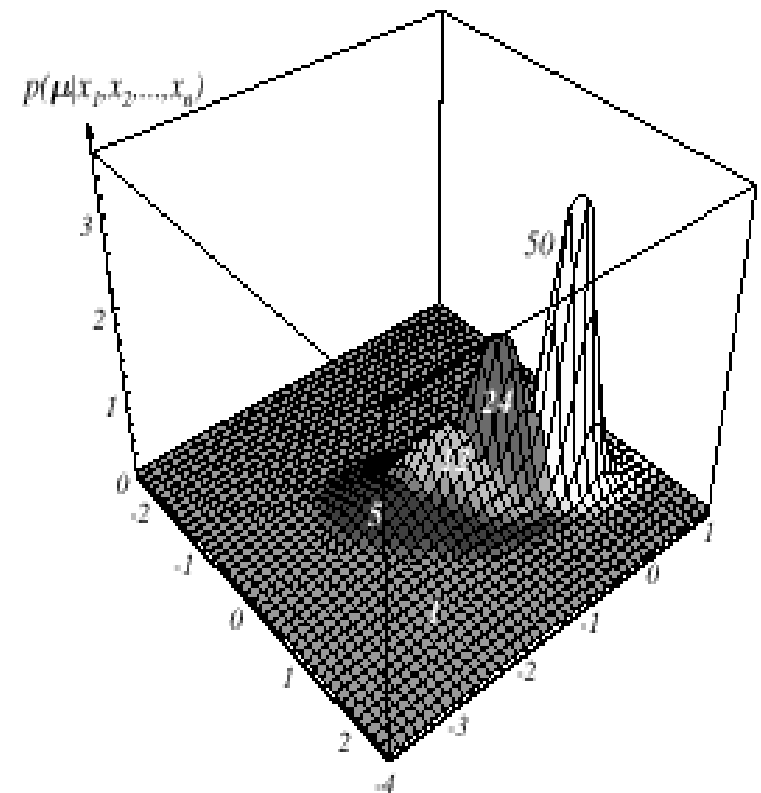
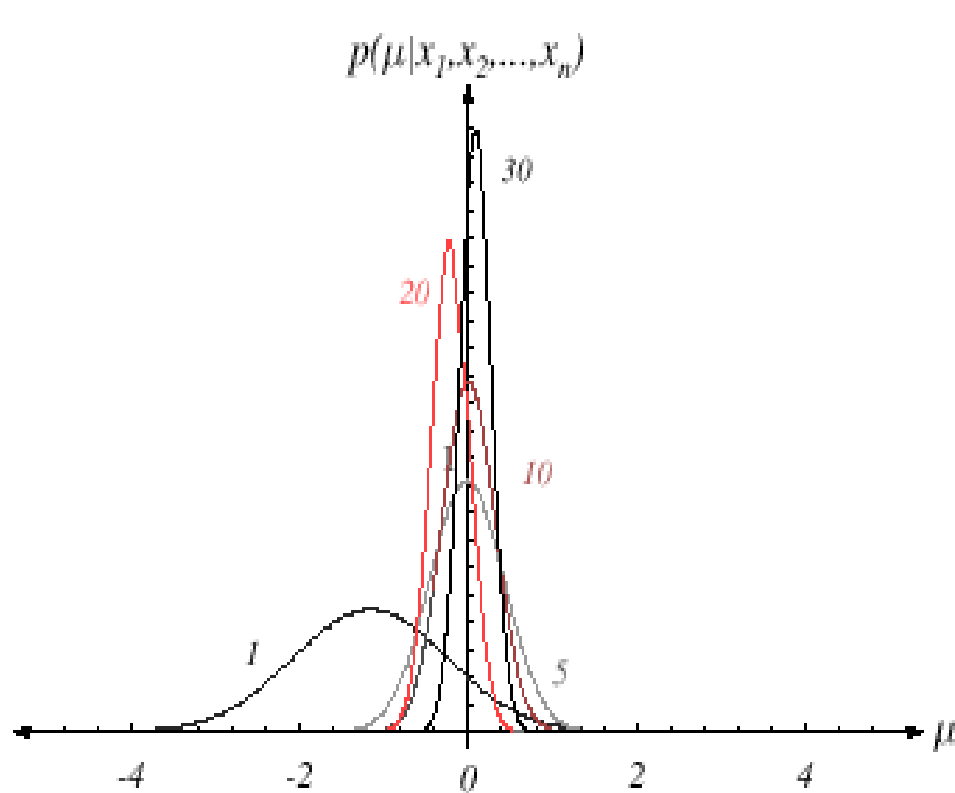


FIGURE 3.2. Bayesian learning of the mean of normal distributions in one and two dimensions. The posterior distribution estimates are labeled by the number of training samples used in the estimation. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



► The univariate case $P(x | D)$

- $P(\mu | D)$ computed as above
- $P(x | D)$ remains to be computed!

$P(x | D) = \int P(x | \mu) \cdot P(\mu | D) d\mu$ is Gaussian

- It provides:

$$P(x | D) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$$

(Desired class-conditional density $P(x | D_j, \omega_j)$)

Therefore: $P(x | D_j, \omega_j)$ together with $P(\omega_j)$ are known

And using Bayes formula, we obtain the Bayesian classification rule:

$$\underset{\omega_j}{Max} [P(\omega_j | x, D)] \equiv \underset{\omega_j}{Max} [P(x | \omega_j, D_j) \cdot P(\omega_j)]$$



3.5 Bayesian Parameter Estimation: General Theory

► $P(x \mid D)$ computation can be applied to any situation in which the unknown density can be parametrized. the basic assumptions are:

- The form of $P(x \mid \theta)$ is assumed known, but the value of θ is not known exactly
- Our knowledge about θ is assumed to be contained in a known prior density $P(\theta)$
- The rest of our knowledge θ is contained in a set D of n random variables x_1, x_2, \dots, x_n that follows unknown $P(x)$



► The basic problem is:

“Compute the posterior density $P(\theta \mid D)$ ”

then “Derive

$$p(x \mid D) = \int p(x \mid \theta) p(\theta \mid D) d\theta$$

”

Using Bayes formula, we have:

$$P(\theta \mid D) = \frac{P(D \mid \theta) \cdot P(\theta)}{\int P(D \mid \theta) \cdot P(\theta) d\theta},$$

And by independence assumption:

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



► Bayse incremental learning

$$D^n = \{x_1, \dots, x_n\}$$

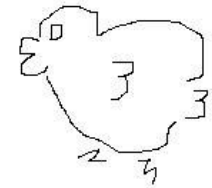
$$p(D^n | \theta) = p(x_n | \theta) p(D^{n-1} | \theta)$$

$$p(\theta | D^n) = \frac{p(D^n | \theta) p(\theta)}{\int p(D^n | \theta) p(\theta) d\theta} = \frac{p(x_n | \theta) p(D^{n-1} | \theta) p(\theta)}{\int p(x_n | \theta) p(D^{n-1} | \theta) p(\theta) d\theta}$$

$$= \frac{p(x_n | \theta) \frac{p(D^{n-1} | \theta) p(\theta)}{p(D^{n-1})}}{\int p(x_n | \theta) \frac{p(D^{n-1} | \theta) p(\theta)}{p(D^{n-1})} d\theta}$$

$$= \frac{p(x_n | \theta) p(\theta | D^{n-1})}{\int p(x_n | \theta) p(\theta | D^{n-1}) d\theta}$$

$$p(\theta | D^0) = p(\theta)$$



Maximum Likelihood vs Bayse Estimation

► Which is better ?

► Maximum Likelihood vs Bayse Estimation

- Computational complexity: ML
- Interpretability : ML
- Confidence in prior information

► Source of classification error

- Bayes Error
- Model Error
- Estimation Error



- Drawbacks of ML estimation
 - ▶ Some observations may be not consistent with the fact

