

# Maximum-Likelihood and Bayesian Parameter Estimation (part 2)

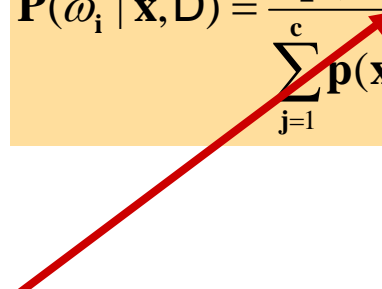
Bayesian Estimation

Bayesian Parameter Estimation: Gaussian Case

Bayesian Parameter Estimation: General Estimation

# Bayesian Estimation

- The parameter  $\theta$  is a random variable
  - Computation of posterior probabilities  $P(\omega_i | \mathbf{x})$  lies at the heart of Bayesian classification
  - Goal: compute  $P(\omega_i | \mathbf{x}, D)$
  - Given the sample  $D$ , Bayes formula is written

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


- $c$  separate parameter estimation problems  $p(\mathbf{x}|D)$

# Parameter Distribution

- Although desired pdf  $p(x)$  is unknown, we assume that it has a known parametric form.
- Only thing unknown is the value of parameter vector  $\theta$

$$p(x | D) = \int p(x, \theta | D) d\theta$$

or

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



Core Task remaining

# Bayesian Parameter Estimation: Univariate Gaussian Case

Assume  $\mu$  is only unknown,  $\sigma$ ,  $\mu_0$  and  $\sigma_0$  known, ie,

$$p(x | \mu) \sim N(\mu, \sigma^2)$$

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

We need

$$\begin{aligned} p(\mu | \mathbf{D}) &= \frac{p(\mathbf{D} | \mu) \cdot p(\mu)}{\int p(\mathbf{D} | \mu) \cdot p(\mu) d\mu} & (1) \\ &= \alpha \prod_{k=1}^{k=n} p(x_k | \mu) \cdot p(\mu) \end{aligned}$$

$$p(\mu | \mathbf{D}) \sim N(\mu_n, \sigma_n^2) \quad (2)$$

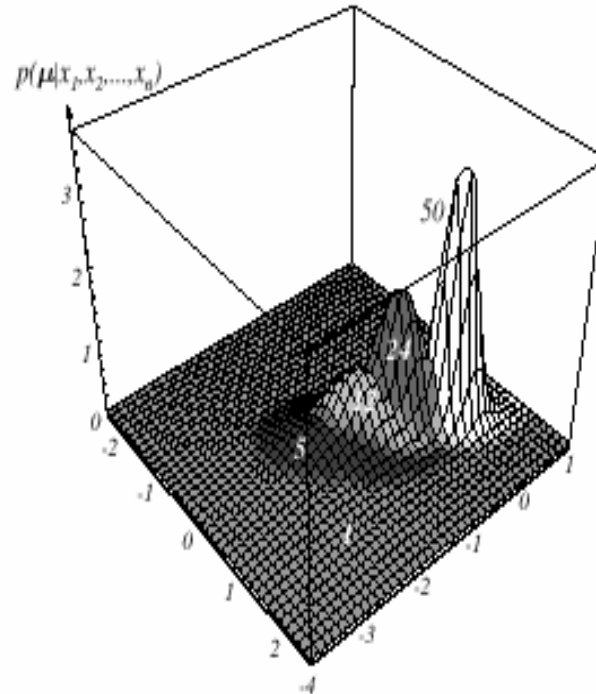
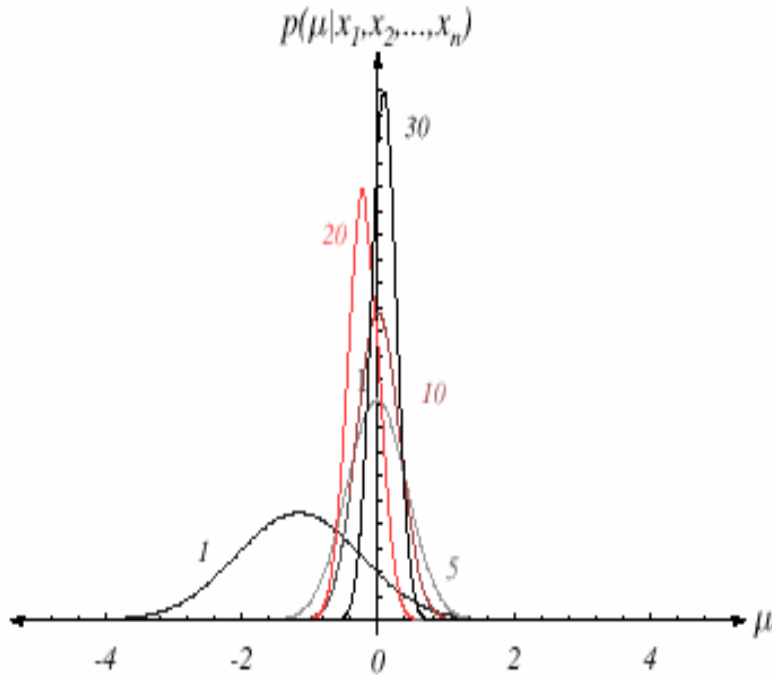
Reproducing density

$$\mu_n = \left( \frac{n\sigma_0^2}{n_0\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \cdot \mu_0$$

$$\text{and } \sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}$$

Prior information is combined with the empirical information in the samples to obtain the *a posteriori* density

# Bayesian Learning



$$\mu_n = \left( \frac{n\sigma_0^2}{n_0\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n_0\sigma_0^2 + \sigma^2} \cdot \mu_0$$

$$\text{and } \sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n_0\sigma_0^2 + \sigma^2}$$

$$p(\mu | \mathbf{D}) \sim N(\mu_n, \sigma_n^2)$$

A linear combination of assumed means

# Computing $p(x | D)$

- $p(\mu | D)$  is computed
- $p(x | D)$  remains

$$\begin{aligned} p(x | D) &= \int p(x | \mu) \cdot p(\mu | D) d\mu \\ &= \int \frac{1}{\sqrt{2\pi\sigma}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] \frac{1}{\sqrt{2\pi\sigma_n}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu_n}{\sigma_n}\right)^2\right] d\mu \\ &= \frac{1}{2\pi\sigma\sigma_n} \exp\left[-\frac{1}{2}\frac{(x-\mu_n)^2}{\sigma^2 + \sigma_n^2}\right] f(\sigma, \sigma_n) \end{aligned}$$

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

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

Multivariate case is similar:

$$p(x | D) \sim N(\mu_n, \Sigma + \Sigma_n)$$

# Bayesian Parameter Estimation: General Theory

$p(\mathbf{x} | D)$  computation can be applied to any situation in which unknown density can be parameterized

Basic assumptions:

- Form of  $p(\mathbf{x} | \theta)$  known, value of  $\theta$  not known exactly
- Initial knowledge of  $\theta$  in known prior density  $p(\theta)$
- Rest of knowledge about  $\theta$  is contained in a set  $D$  of  $n$  random variables  $x_1, x_2, \dots, x_n$  that follows  $p(\mathbf{x})$

# General Bayesian Parameter Estimation

Compute posterior density  $p(\theta | D)$  then  $p(x | D)$  using

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

Using Bayes formula:

$$p(\theta | D) = \frac{p(D | \theta) \cdot p(\theta)}{\int p(D | \theta) \cdot p(\theta) d\theta},$$

By independence assumption:

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



# Recursive Bayes Incremental Learning

- Explicitly indicate number of samples in a set for a given category as  $D^n = \{x_1, \dots, x_n\}$

- Then from  $p(D | \theta) = \prod_{k=1}^{k=n} p(x_k | \theta)$  we can write

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

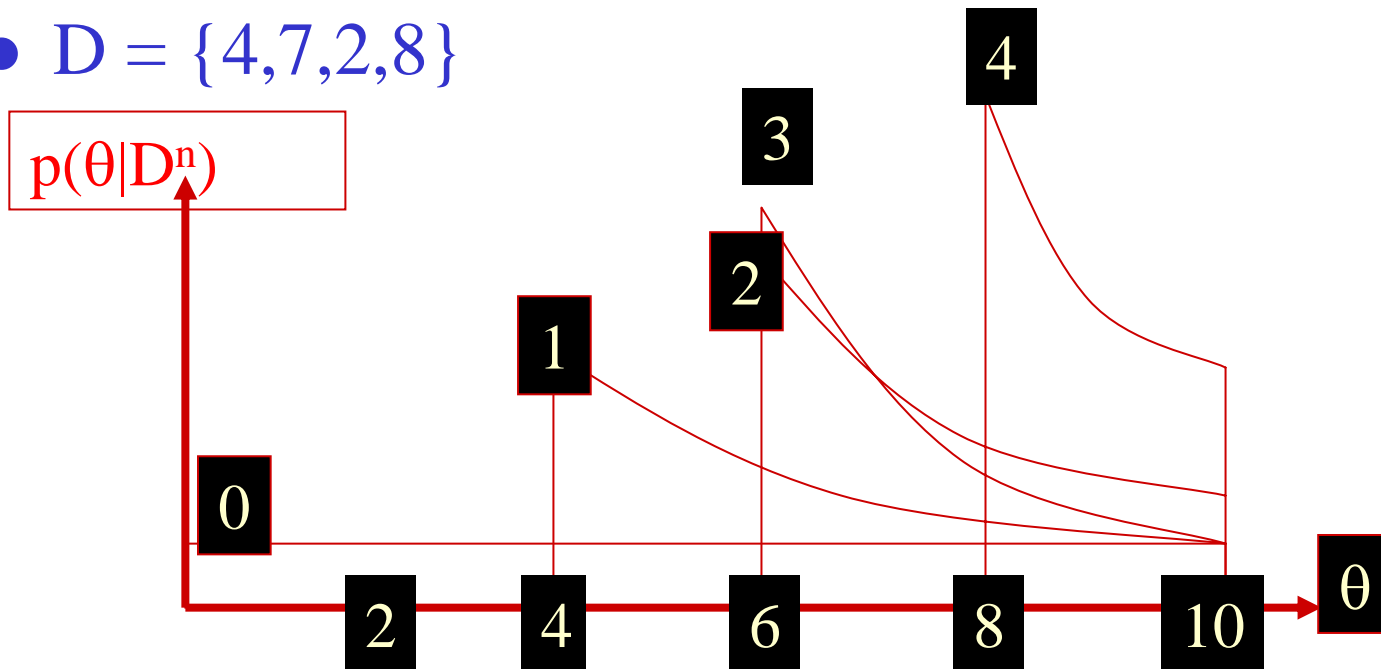
$$p(\theta | D^n) = \frac{p(x_n | \theta) \cdot p(\theta | D^{n-1})}{\int p(x_n | \theta) \cdot p(\theta | D^{n-1}) d\theta},$$

# Example of Recursive Bayes Learning

- One-dim samples from uniform distribution

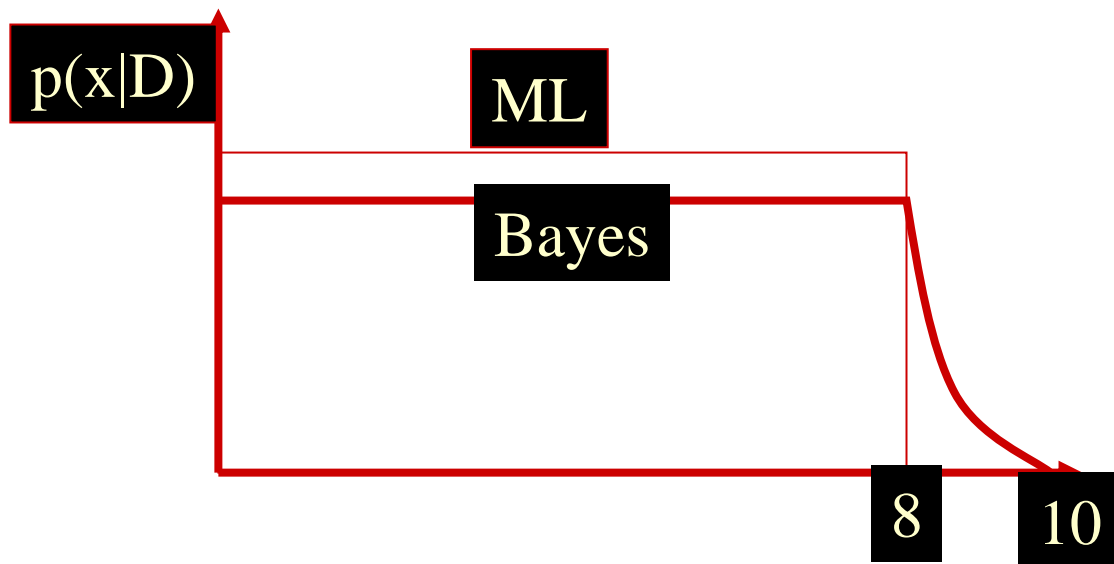
$$p(x | \theta) \sim U(0, \theta) = \begin{cases} \frac{1}{\theta} & 0 \leq x \leq \theta \\ 0 & \text{otherwise} \end{cases}$$

- Parameter distribution: Uniform over  $0 \leq \theta \leq 10$
- $D = \{4, 7, 2, 8\}$



# Learning the parameter of Uniform Distribution using Recursive Bayes

- As more points are incorporated
- Bayes has a tail above 8 reflecting prior information



# Three Sources of Error in Classification

- Bayes or Indistinguishability Error
  - Due to overlapping densities. Inherent property of given feature set
- Model Error
  - Due to an incorrect model. Can only be eliminated if designer specifies true model that generated the data.
- Estimation Error
  - Parameters are estimated from a finite sample.