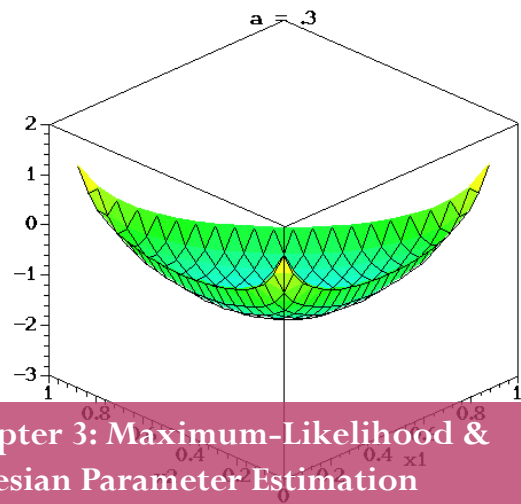


# PATTERN RECOGNITION



## Chapter 3: Maximum-Likelihood & Bayesian Parameter Estimation

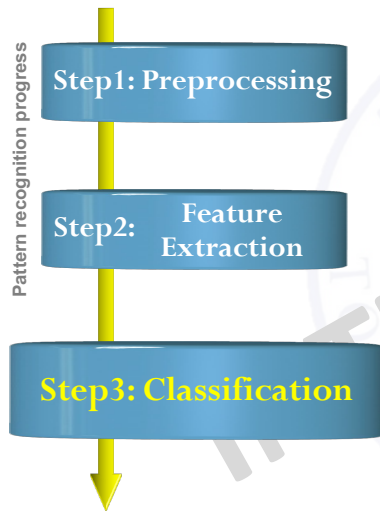
### Ch.3 Content



- 1 **Introduction**
- 2 Maximum-Likelihood Estimation
- 3 Bayesian Estimation

# 1. Introduction

- Pattern Recognition System



## Bayesian Classification

$$P(\omega_i) = ?$$

$$p(x|\omega_i) = ?$$

### Minimize error rate decision

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

### Minimum-Risk Decision decision rule

$$- R(\alpha_i|x) = \sum_{j=1}^c \lambda(\alpha_i|\omega_j)P(\omega_j|x), i = 1, 2, \dots, a$$

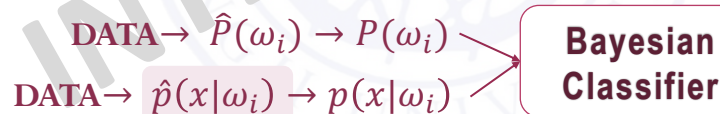
# 1. Introduction

- **Parameter Estimation** — Data availability in a Bayesian framework

$$P(\omega_i) = ?$$

$$p(x|\omega_i) = ?$$

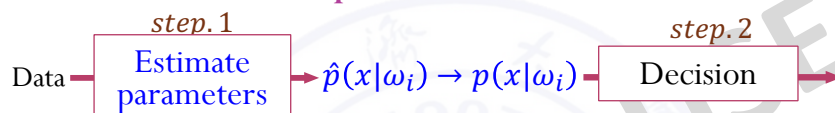
- In pattern recognition applications we rarely have this kind of complete knowledge about the probabilistic structure of the problem.
  - We have a number of design samples or training data.
  - The problem is to find some way to use this information to design or train the classifier.
  - One approach is to use the samples to estimate the unknown probabilities and probability densities, and to use the resulting estimates as if they were the true values.



# 1. Introduction



- Bayes Decision based on samples



**Example:**  $p(x|\omega_i) \sim N(\mu_i, \Sigma_i)$   
Samples are independent and identically distributed (i.i.d.)



- Class-conditional densities  $p(x|\omega_i)$

- Nonparameter Estimation
- **Parameter Estimation:** If we know the **number of parameters** in advance and our general knowledge about the problem permits us to **parameterize the conditional densities** then severity of the problem can be reduced significantly.

# 1. Introduction



- **Parameter Estimation techniques**

- Maximum-Likelihood Estimations(MLE)
- Bayesian Estimations(BE)

## MLE

- Parameters are **fixed** but unknown.
- Best parameters are obtained by maximizing the probability of obtaining the samples observed.

## BE

- Parameters are **random variables** having some known distribution.
- Best parameters are obtained by estimating them given the data.

- Results are nearly identical, but the approaches are different.

- 1 Introduction
- 2 **Maximum-Likelihood Estimation**
- 3 Bayesian Estimation

## 2. Maximum-Likelihood Estimation

- **Maximum-likelihood Estimation (MLE)**

- **GOAL:** determine the most likely values of the population parameter value (e.g,  $\mu$ ,  $\sigma$ ,  $\beta$ ,  $\rho$ , ... ) given an observed sample value (e.g,  $\bar{x}$ ,  $s$ ,  $b$ ,  $r$ , ....).



- **Simpler** than any other alternative techniques.

- General principle
- Gaussian Case: unknown  $\mu$
- Gaussian Case: unknown  $\mu$  and  $\sigma$

## 2. Maximum-Likelihood Estimation

### • General principle

#### • Preliminaries and Notations

$$p(x|\omega_i) = p(x|\omega_i, \theta_i)$$

$\theta = (\theta_1, \theta_2 \dots \theta_c)$ : unknown parameters

$i = 1, 2, \dots c$

$\theta_i = (\theta_1, \theta_2 \dots \theta_p)^t$ : Each  $\theta_i$  is associated with each category.

$\hat{\theta}$ : estimated parameter

$D_i = \{x_1, \dots, x_n\}$ :  $n$  samples (i.i.d.)

■ **Task**: find  $\hat{\theta}$  of parameters  $\theta$

■ **Method**: ML estimation of  $\theta$  is the value  $\hat{\theta}$  that maximizes  $p(D|\theta)$ .

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

—  $p(D|\theta)$  is called the likelihood of  $\theta$  with respect to the set of samples.

**“It is the value of  $\hat{\theta}$  that best agrees with the actually observed training sample”**

$$\begin{aligned} \text{Max}_{\theta} p(D|\theta) &\Rightarrow \hat{\theta} = \underset{\theta}{\operatorname{argmax}} l(\theta) = \underset{\theta}{\operatorname{argmax}} p(D|\theta) \\ &= \underset{\theta}{\operatorname{argmax}} p(x_1, \dots, x_n|\theta) = \underset{\theta}{\operatorname{argmax}} \prod_{k=1}^n p(x_k|\theta) \end{aligned}$$

## 2. Maximum-Likelihood Estimation

### • General principle

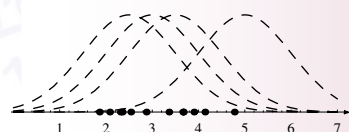
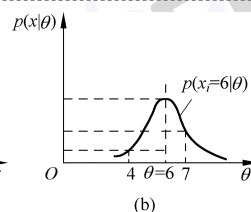
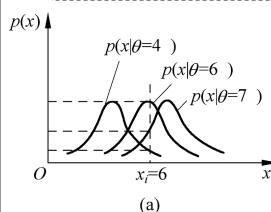
**Example**  $p(x|\omega_i) \sim N(\mu_i, \Sigma_i)$

$$p(x|\omega_i) = p(x|\omega_i, \theta_i)$$

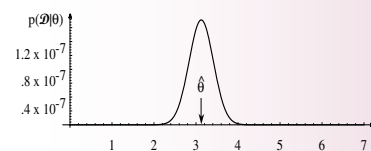
—  $\theta_i = (\mu_i) = (\mu_i^1, \mu_i^2, \dots)$ ,  $i = 1, 2, \dots c$ ,  $\mu_i$  is unknown.

— Use the information provided by the training samples to estimate  $\theta = (\theta_1, \theta_2 \dots \theta_c)$ , each  $\theta_i$  is associated with each category.

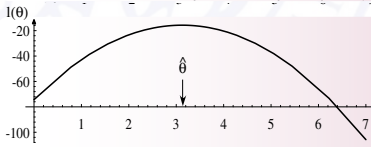
— log-likelihood function.



Shows several training points in one dimension, known variance, but unknown mean. Four of the infinite number of candidate source distributions are shown in dashed lines.



Shows the likelihood  $p(D|\theta)$ . If we had a very large number of training points, this likelihood would be very narrow. The value that maximizes the likelihood is marked  $\hat{\theta}$ .



Shows the log-likelihood  $l(\theta)$  ( $\ln l(\theta)$ ) which marks the maximum  $\hat{\theta}$ .

## 2. Maximum-Likelihood Estimation

- **General principle**

- **Method: MLE log-likelihood function**

- $D$  contains  $n$  samples:  $\{x_1, x_2 \dots x_n\}$
- Unknown parameters:  $\theta = (\theta_1, \theta_2 \dots \theta_p)^t$

### STEP

- ① We define  $l(\theta)$  as the log-likelihood function:  

$$l(\theta) = \ln p(D|\theta) = \sum_{k=1}^n \ln p(x_k|\theta)$$
- ② determine  $\hat{\theta}$  that maximizes the log-likelihood  $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} l(\theta)$
- ③ Set of necessary conditions for an optimum is:  

$$\nabla_{\theta} l = \sum_{k=1}^n \nabla_{\theta} \ln p(x_k|\theta)$$
- ④ Set the derivative of  $l(\theta)$  equal to zero and solving for  $\hat{\theta}$ .

$$\nabla_{\theta} l = 0 \Rightarrow \frac{d[l(\theta)]}{d\theta} = \frac{d[\ln p(x_k|\theta)]}{d\theta} = 0$$

$\nabla_{\theta}$  is the gradient operator

$$\nabla_{\theta} = \begin{bmatrix} \frac{\partial}{\partial \theta_1} \\ \frac{\partial}{\partial \theta_2} \\ \vdots \\ \frac{\partial}{\partial \theta_p} \end{bmatrix}$$

## 2. Maximum-Likelihood Estimation

- **Gaussian Case: unknown  $\mu$**

- Multivariate normal population:  $p(x_k|\theta) \sim N(\mu, \Sigma)$   
 (Samples are drawn from a multivariate normal population)

— Unknown parameter:  $\theta = \mu$

— Sample set  $D: \{x_1, x_2 \dots x_n\}$ , a sample point:  $x_k$

### MLE

- ① **log-likelihood function:**  $l(\theta) = \ln p(x_k|\mu) = -\frac{1}{2} \ln[(2\pi)^d |\Sigma|] - \frac{1}{2} (x_k - \mu)^t \Sigma^{-1} (x_k - \mu)$
- ② gradient operation:  $\nabla_{\theta} l(\theta) = \nabla_{\mu} \ln p(x_k|\mu) = \Sigma^{-1} (x_k - \mu) = 0$
- ③ ML estimate for  $\mu$ :  $\sum_{k=1}^n \Sigma^{-1} (x_k - \hat{\mu}) = 0$ ,  $D: \{x_1, x_2 \dots x_n\}$

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

Just the arithmetic average of the samples of the training samples.

## 2. Maximum-Likelihood Estimation



- **Gaussian Case: unknown  $\mu$  and  $\sigma$**

- Univariate case:  $p(x_k|\theta) \sim N(\mu, \sigma^2)$

(Samples are drawn from a univariate normal population)

– Unknown parameters:  $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$

– Sample set  $D: \{x_1, x_2 \dots x_n\}$ , a sample point  $x_k$

### MLE

$$\begin{aligned} \textcircled{1} \quad l &= \ln p(x_k|\theta) = -\frac{1}{2} \ln 2\pi \theta_2 - \frac{1}{2\theta_2} (x_k - \theta_1)^2 \\ \textcircled{2} \quad \nabla_{\theta} l &= \nabla_{\theta} \ln P(x_k|\theta) = \begin{bmatrix} \frac{\partial}{\partial \theta_1} \ln p(x_k|\theta) \\ \frac{\partial}{\partial \theta_2} \ln p(x_k|\theta) \end{bmatrix} = 0 \end{aligned}$$

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

$$\begin{cases} \sum_{k=1}^n \frac{1}{\theta_2} (x_k - \hat{\theta}_1) = 0 & (1) \\ -\sum_{k=1}^n \frac{1}{\theta_2} + \sum_{k=1}^n \frac{(x_k - \hat{\theta}_1)^2}{\theta_2^2} = 0 & (2) \end{cases}$$

$$\begin{aligned} \textcircled{3} \quad & \hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k \\ \textcircled{4} \quad & \hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})^2 \end{aligned}$$

## 2. Maximum-Likelihood Estimation



- **Gaussian Case: unknown  $\mu$  and  $\Sigma$**

- Multivariate case:  $p(x_k|\theta) \sim N(\mu, \Sigma)$

(Samples are drawn from a multivariate normal population)

– Unknown parameters  $\theta: \theta_1 = \mu, \theta_2 = \Sigma$

– Sample set  $D: \{x_1, x_2 \dots x_n\}$ , a sample point  $x_k$

### MLE

$$\begin{aligned} \hat{\mu} &= \frac{1}{n} \sum_{k=1}^n x_k \\ \hat{\Sigma} &= \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})(x_k - \hat{\mu})^t \end{aligned}$$



## 2. Maximum-Likelihood Estimation



- **Gaussian Case:** Bias

	Biased estimator	unbiased estimator
$\sigma^2$	$E \left[ \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})^2 \right] = \frac{n-1}{n} \sigma^2 \neq \sigma^2$	$E \left[ \frac{1}{n-1} \sum_{k=1}^n (x_k - \hat{\mu})^2 \right] = \sigma^2$
covariance matrix	$\hat{\Sigma} = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})(x_k - \hat{\mu})^t = \frac{n-1}{n} \Sigma$	$\Sigma = \frac{1}{n-1} \sum_{k=1}^n (x_k - \hat{\mu})(x_k - \hat{\mu})^t$

asymptotically unbiased  
when  $n$  is very large

## 2. Maximum-Likelihood Estimation



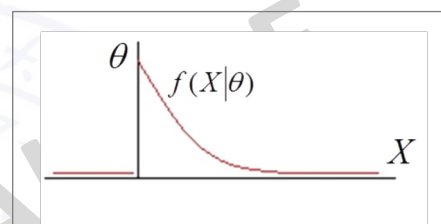
### Exercise

Consider an exponential distribution

$$f(x; \theta) = \begin{cases} \theta e^{-\theta x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

(single feature, single parameter),

With a random sample set  $\{x_1, x_2 \dots x_n\}$ , i. d. d.,  
**estimate  $\theta$  ?**



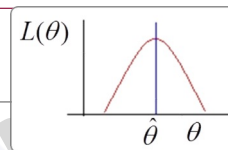


## 2. Maximum-Likelihood Estimation



Exercise

$$f(x; \theta) = \begin{cases} \theta e^{-\theta x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$



- ①  $L(\theta) = f(x_1, x_2, \dots, x_n | \theta) = \prod_{k=1}^n \theta e^{-\theta x_k}$
- ②  $l(\theta) = \ln L(\theta) = \sum_{k=1}^n \ln \theta - \theta \sum_{k=1}^n x_k$   
 $= n \ln \theta - \theta \sum_{k=1}^n x_k, \quad \text{valid for } x \geq 0$
- ③  $\frac{dl}{d\theta} = \frac{d \ln L(\theta)}{d\theta} = \frac{n}{\theta} - \sum_{k=1}^n x_k = 0$   
 $\Rightarrow \frac{n}{\hat{\theta}} = \sum_{k=1}^n x_k$   
 $\Rightarrow \hat{\theta} = \frac{1}{\frac{1}{n} \sum_{k=1}^n x_k}$

## 2. Maximum-Likelihood Estimation



### • General principle: MLE log-likelihood function

- **Note1:** the likelihood function **must be differentiable** and then this method can be used.
- **Note2:**  $\hat{\theta}$  could represent a true **global maximum**, a **local maximum** or **minimum**, or an inflection point of  $l(\theta)$ . We should **check** each solution to identify which is the **global optimum**.

Example

$D = \{x_1, x_2 \dots x_n\}$ , iid.

$$p(x|\theta) = \begin{cases} \frac{1}{\theta_2 - \theta_1} & \theta_1 < x < \theta_2 \\ 0 & \text{others} \end{cases}$$

Estimate  $\theta_1, \theta_2$ ?

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

$$= \begin{cases} \ln p(x_1, x_2 \dots x_n | \theta) = -n \ln(\theta_2 - \theta_1) & \theta_1 < x < \theta_2 \\ 0 & \text{others} \end{cases}$$

$$= \begin{cases} \frac{\partial}{\partial \theta_1} \ln p(x_1, x_2 \dots x_n | \theta) = n \frac{1}{\theta_2 - \theta_1} \\ \frac{\partial}{\partial \theta_2} \ln p(x_1, x_2 \dots x_n | \theta) = -n \frac{1}{\theta_2 - \theta_1} \end{cases}$$

It is uninteresting. Need to find a new method.

- ① Introduction
- ② Maximum-Likelihood Estimation
- ③ **Bayesian Estimation**

### 3. Bayesian Estimation

- **MLE & BE**

- In MLE,  $\theta$  was supposed fix.
- In BE,  $\theta$  is a random variable.

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

### 3. Bayesian Estimation

#### • Bayesian Estimation

- 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  $D = D_1 \cup D_2 \dots \cup D_c$ ,
- $x \in D_i$ ,  $D_i$  has no influence on  $p(x|\omega_j, D_j)$  if  $i \neq j$ .  $p(x|\omega_i, D) = p(x|\omega_i, D_i)$ .
- $P(\omega_i) = P(\omega_i|D)$

- Bayes formula can be written:

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

1 The *functional forms* for unknown densities;

2 Ranges for the values of unknown parameters;

3 Training samples  $D$ .

Prior information

### 3. Bayesian Estimation

#### • Bayesian Estimation

- $p(x)$  is unknown
  - We assume it has a *known parametric form*  $p(x|\theta)$ ;
  - Value of parameter  $\theta$  is unknown,  $p(\theta) \rightarrow p(\theta|D)$ .



Bayesian Estimation

$p(x)$   
Probability density

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

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

- If  $p(\theta|D)$  peaks very sharply about parameter  $\hat{\theta}$  and  $p(x|\theta)$  is smooth, then  $p(x|D) \approx p(x|\hat{\theta})$ .

### 3. Bayesian Estimation



- **Bayesian Parameter Estimation: Gaussian Case**

- **GOAL:**

- Calculate the a-posteriori density  $p(\theta|D)$
    - Calculate the desired probability density  $p(x|D)$ .

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

- Assume:  $p(x|\mu) \sim N(\mu, \sigma^2)$ , the only unknown parameter  $\mu$ .

- The univariate case:  $p(\mu|D)$ ,  $p(x|D)$
    - The multivariate case:  $p(\mu|D)$ ,  $p(x|D)$ .

### 3. Bayesian Estimation



- **Bayesian Parameter Estimation: Gaussian Case**

- The univariate case:  $p(\mu|D)$

- Known prior density  $p(\mu) \sim N(\mu_0, \sigma_0^2)$ 
      - Assume  $\mu_0$  and  $\sigma_0$  are known.
      - $\mu_0$  represents our best a priori guess for  $\mu$ .
      - $\sigma_0^2$  measures our uncertainty about this guess.

$$p(\mu|D) = \frac{p(D|\mu)p(\mu)}{\int p(D|\mu)p(\mu)d\mu}$$

$$D = \{x_1, x_2, \dots, x_n\}$$

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

$$\begin{aligned} p(x_k|\mu) &\sim N(\mu, \sigma^2) \\ p(\mu) &\sim N(\mu_0, \sigma_0^2) \end{aligned}$$

$$\begin{aligned} p(\mu|D) &= \alpha \prod_{k=1}^n \overbrace{\frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x_k - \mu}{\sigma}\right)^2\right]}^{p(x_k|\mu)} \overbrace{\frac{1}{\sqrt{2\pi}\sigma_0} \exp\left[-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right]}^{p(\mu)} \\ &= \alpha' \exp\left[-\frac{1}{2}\left(\sum_{k=1}^n \left(\frac{\mu - x_k}{\sigma}\right)^2 + \left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right)\right] \\ &= \alpha'' \exp\left[-\frac{1}{2}\left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}\right)\mu^2 - 2\left(\frac{1}{\sigma^2}\sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2}\right)\mu\right]\right] \end{aligned}$$

### 3. Bayesian Estimation

#### • Bayesian Parameter Estimation: Gaussian Case

##### □ The univariate case: $p(\mu|D)$

$$p(\mu|D) = \frac{1}{\sqrt{2\pi}\sigma_n} \exp\left[-\frac{1}{2}\left(\frac{\mu - \mu_n}{\sigma_n}\right)^2\right]$$

$$\frac{1}{\sigma_n^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}$$

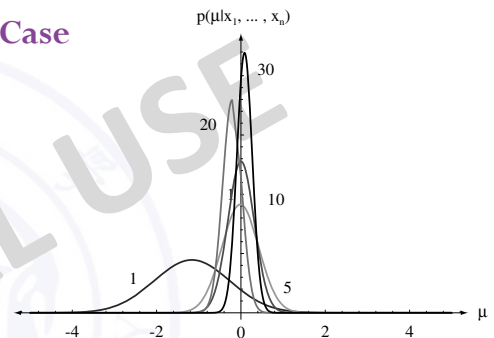
$$\mu_n = \frac{n}{\sigma^2} x_n + \frac{\mu_0}{\sigma_0^2}$$

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

$$\mu_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2}\right) x_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0$$

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

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



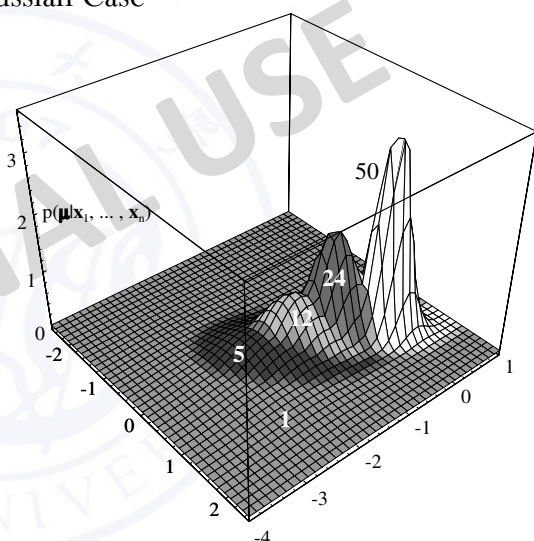
- $\mu_n$  represents our best guess for  $\mu$  after observing  $n$  samples.
- $\sigma_n^2$  measures our uncertainty about this guess.  $\sigma_n^2$  decreases monotonically as  $n$  increases. As the same time,  $p(\mu|D)$  becomes more and more sharply peaked.
- This behavior is commonly known as **Bayesian learning**.

### 3. Bayesian Estimation

#### • Bayesian Parameter Estimation: Gaussian Case

##### • The 2D case: $p(\mu|D)$

- Bayesian learning of the mean of normal distributions in two dimensions. The posterior distribution estimates are labeled by the number of training samples used in the estimation.



### 3. Bayesian Estimation

- **Bayesian Parameter Estimation: Gaussian Case**

- The univariate case:  $p(x|D)$

- $p(\mu|D)$  has been computed  $p(\mu|D) \sim N(\mu_n, \sigma_n^2)$
    - $p(\mu|D) \rightarrow p(x|D)$

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

$$\begin{aligned}
 \mu_n &= \left(\frac{n\sigma_0^2}{n\sigma_0^2+\sigma^2}\right)\hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2+\sigma^2}\mu_0 \\
 \sigma_n^2 &= \frac{\sigma_0^2\sigma^2}{n\sigma_0^2+\sigma^2}
 \end{aligned}$$

$$f(\sigma, \sigma_n) = \int \exp\left[-\frac{1}{2}\frac{\sigma^2+\sigma_n^2}{\sigma^2\sigma_n^2}\left(\mu - \frac{\sigma_n^2x + \sigma^2\mu_n}{\sigma^2+\sigma_n^2}\right)^2\right] d\mu$$

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

— The density  $p(x|D)$  is the desired class-conditional density  $p(x|D_j, \omega_j)$ , and together with the prior probabilities  $P(\omega_j)$  it gives us the probabilistic information needed to design the classifier.

### 3. Bayesian Estimation

- **Bayesian Parameter Estimation: Gaussian Case**

- The multivariate case:  $p(x|D) \sim N(\mu, \Sigma)$

- $\mu$  is the only unknown parameter
    - $\Sigma$  is known
    - Known prior density  $p(\mu) \sim N(\mu_0, \Sigma_0^2)$
    - $\mu_0$  and  $\Sigma_0$  are known!

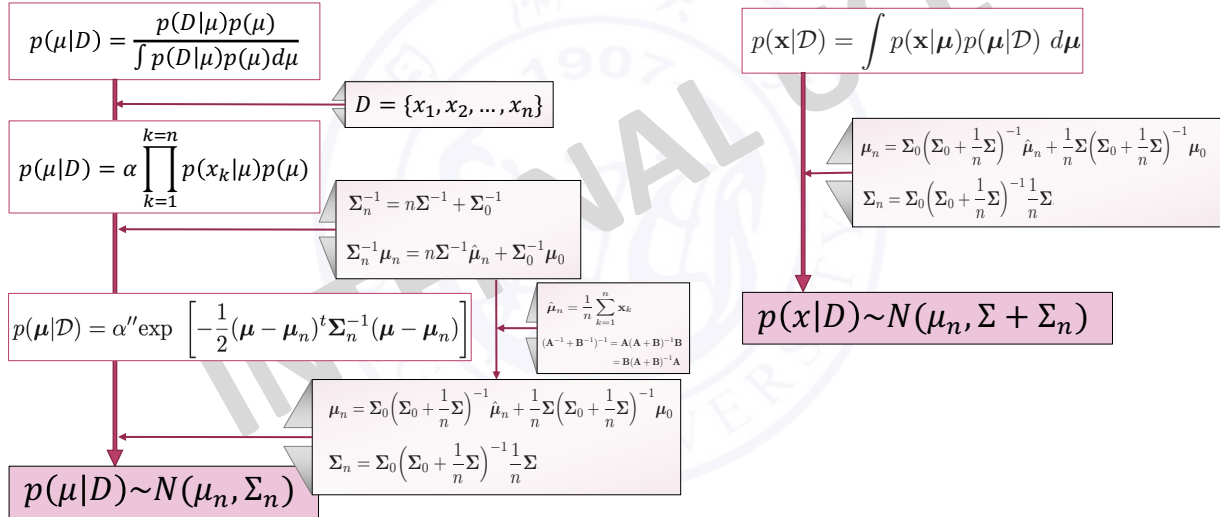
- $p(\mu|D)$

- $p(x|D)$

### 3. Bayesian Estimation

#### • Bayesian Parameter Estimation: Gaussian Case

##### □ The multivariate case: $p(\mu|D)$



### 3. Bayesian Estimation

#### • Bayesian Parameter Estimation:

##### General Theory

- This approach can be generalized to apply to any situation in which the unknown density can be parameterized. The basic assumptions are summarized as follows:
  - The form of the density  $p(x|\theta)$  is assumed to be known, but the value of the parameter vector  $\theta$  is not known exactly.
  - Our initial knowledge about  $\theta$  is assumed to be contained in a known prior density  $p(\theta)$ .
  - The rest of our knowledge about  $\theta$  is contained in a set  $D$  of  $n$  samples  $x_1, x_2, \dots, x_n$  drawn independently according to the unknown probability density  $p(x)$ .

$$p(x|D) = \int p(x|\theta) p(\theta|D) d\theta \approx p(x|\hat{\theta})$$

—If  $p(\theta|D)$  peaks very sharply about parameter  $\hat{\theta}$  and  $p(x|\theta)$  is smooth, then  $p(x|D) \approx p(x|\hat{\theta})$ .

$$p(\theta|D) \rightarrow p(x|D)$$

$$\textcircled{1} p(\theta|D) = \frac{p(D|\theta)p(\theta)}{\int p(D|\theta)p(\theta)d\theta}$$

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

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

$$\textcircled{2} p(x|D) = \int p(x|\theta) p(\theta|D) d\theta$$



### 3. Bayesian Estimation

- Recursive Bayes approach

$$\begin{aligned}
 p(\theta|D^n) &= \frac{p(D^n|\theta)p(\theta)}{\int p(D^n|\theta)p(\theta)d\theta} \\
 &\quad \leftarrow p(D^n|\theta) = p(x_n|\theta)p(D^{n-1}|\theta) \\
 &= \frac{p(x_n|\theta)p(D^{n-1}|\theta)p(\theta)}{\int p(x_n|\theta)p(D^{n-1}|\theta)p(\theta)d\theta} \\
 &= \frac{p(x_n|\theta)p(\theta|D^{n-1})}{\int p(x_n|\theta)p(\theta|D^{n-1})d\theta}
 \end{aligned}$$

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

#### Incremental learning:

$$\begin{aligned}
 p(\theta) &\xrightarrow{D^0 = \{\}} p(\theta|D^0) = p(\theta) \\
 p(\theta|x_1) &\xrightarrow{D^1 = \{x_1\}} p(\theta|D^1) \propto p(x_1|\theta)p(\theta|D^0) \\
 &\vdots \\
 p(\theta|x_1, x_1 \dots x_n) &\xrightarrow{D^n = \{x_1, x_2, \dots x_n\}} p(\theta|D^n) \propto p(x_n|\theta)p(\theta|D^{n-1})
 \end{aligned}$$

### 3. Bayesian Estimation

- Recursive Bayes Learning

#### Example

Suppose we believe our one dimensional samples come from a uniform distribution

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

But initially we know only that our parameter is bounded  $0 < \theta \leq 10$  (uniform distribution). Use recursive Bayes methods to estimate  $\theta$  and the underlying densities from the data  $D = \{4, 7, 2, 8\}$ .

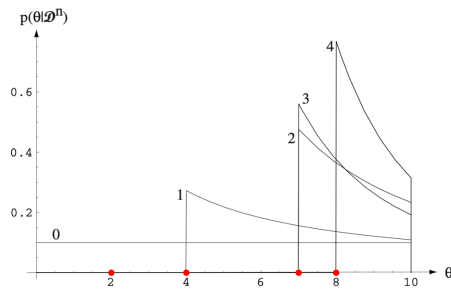
$$\begin{aligned}
 D^0 = \{\} &\quad \dots \quad p(\theta|D^0) = p(\theta) = U(0, 10) \\
 x_1 = 4 &\quad \dots \quad p(\theta|D^1) \propto p(x_1|\theta)p(\theta|D^0) = \begin{cases} 1/\theta & 4 \leq \theta \leq 10 \\ 0 & \text{otherwise} \end{cases} \\
 x_2 = 7 &\quad \dots \quad p(\theta|D^2) \propto p(x_2|\theta)p(\theta|D^1) = \begin{cases} 1/\theta^2 & 7 \leq \theta \leq 10 \\ 0 & \text{otherwise} \end{cases} \\
 &\quad \dots \dots \quad p(\theta|D^n) \propto \frac{1}{\theta^n}, \quad \max_x[D^n] \leq \theta \leq 10
 \end{aligned}$$

### 3. Bayesian Estimation

#### • Recursive Bayes Learning

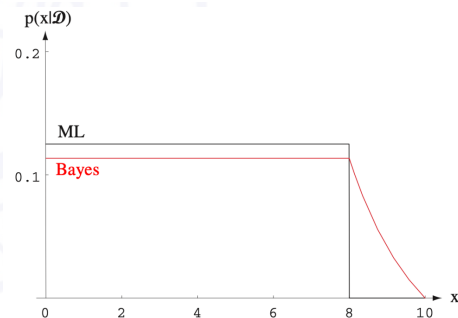
**Example:**  $p(x|\theta) \sim U(0, \theta) = \begin{cases} 1/\theta & 0 \leq x \leq \theta \\ 0 & \text{otherwise} \end{cases}$

The posterior  $p(\theta|\mathcal{D}^n)$  for the model and  $n$  points in the data set in this example. For  $n = 0$ , the posterior starts out as a flat, uniform density from 0 to 10, denote  $p(\theta) \sim U(1, 10)$ . As more points are incorporated, it becomes increasingly peaked at the value of the highest data point.  $\mathcal{D} = \{4, 7, 2, 8\}$



- Given full data set, the maximum-likelihood solution here is clearly  $\hat{\theta} = 8$ , and this implies a uniform  $p(x|\mathcal{D}) \sim U(0, 8)$ .
- According to Bayesian methodology, which requires the integration, the density is uniform up to  $x = 8$ , but has a tail at higher values—an indication that the influence of our prior  $p(\theta)$  has not yet been swamped by the information in the training data.

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



### 3. Bayesian Estimation

#### Maximum-Likelihood Estimation and Bayesian Estimation

	MLE	BE
Method	Maximum likelihood approach estimates a point in $\theta$ space	Bayesian approach estimates a distribution
Computational complexity	Simpler	Bayesian method has strong theoretical and methodological arguments supporting it.
results		Better
	When used for classifiers, they mostly give same result.	

## Chapter 3

### PART I

