

# Maximum Likelihood Parameter Estimations

---

Computer Science- Pattern Recognition

Prof. Dr. Dhahir A. Abdullah

# Parameter Estimation

- In previous chapters:
  - We could design **an optimal classifier** if we **knew** the **prior probabilities**  $P(w_i)$  and the **class-conditional probabilities**  $P(x|w_i)$
- Unfortunately, in pattern recognition applications **we rarely have this kind of complete knowledge** about the probabilistic structure of the problem.

# Parameter Estimation

- 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:
  - **Use the samples to estimate** the unknown probabilities and probability densities,
  - And then use the resulting estimates as if they **were the true values.**

# Parameter Estimation

- 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:
  - **Use the samples to estimate** the unknown probabilities and probability densities,
  - And then use the resulting estimates as if they **were the true values.**

# Parameter Estimation

- What you are going to estimate in your Homework ?

# Parameter Estimation

- If we know the **number 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.

# Parameter Estimation

- For example:
  - We can reasonably assume that the  $p(\mathbf{x}|\mathbf{w}_i)$  is a **normal density** with mean  $\boldsymbol{\mu}$  and covariance matrix  $\Sigma$ ,
  - We do not know the exact values of these quantities,
  - However, this knowledge simplifies the problem **from one of estimating an unknown function  $p(\mathbf{x}|\mathbf{w}_i)$  to one of estimating the parameters the mean  $\boldsymbol{\mu}_i$  and covariance matrix  $\Sigma_i$**
  - **Estimating  $p(\mathbf{x}|\mathbf{w}_i) \rightarrow$  estimating  $\boldsymbol{\mu}_i$  and  $\Sigma_i$**

# Parameter Estimation

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



# Parameter Estimation

- Given a bunch of data from each class how to estimate the parameters of class conditional densities,  $P(x | \omega_i)$  ?
- Ex:  $P(x | \omega_i) = N(\mu_j, \Sigma_j)$  is Normal.  
Parameters  $\theta_j = (\mu_j, \Sigma_j)$

# Two major approaches

- **Maximum-Likelihood Method**
- **Bayesian Method**
  - Use  $P(\omega_i | x)$  for our classification rule!
  - Results are nearly identical, but the approaches are different

# Maximum-Likelihood vs. Bayesian:

- Maximum Likelihood
- Parameters are fixed but unknown!
- Best parameters are obtained by maximizing the probability of obtaining the samples observed
- Bayes
- Parameters are random variables having some known distribution
- Best parameters are obtained by estimating them given the data

# Major assumptions

- A priori information  $P(\omega_i)$  for each category is available
- Samples are i.i.d. and  $P(x | \omega_i)$  is Normal

$$P(x | \omega_i) \sim N(\mu_i, \Sigma_i)$$

- Note: Characterized by 2 parameters

# Maximum-Likelihood Estimation

- Has good convergence properties as the sample size increases
- Simpler than any other alternative techniques

# Maximum-Likelihood Estimation

We assume that samples are collected randomly from a given form of distribution, whose parameters are unknown. Unknown parameters are denoted by the vector  $\theta$ .

We partition our training data  $\mathbf{X}$  into  $c$  class specific subsets  $D_1 \dots D_c$ , assuming that the data in  $D_j$  have been drawn randomly according to the distribution  $p(\mathbf{x}|\omega_j)$

If, for example, we know each distribution is normal,  $p(\mathbf{x}|\omega_j) \sim N(\mu_j, \Sigma_j)$  but we do not know the mean and covariance, then our problem is to estimate

$$\theta_j = \begin{pmatrix} \mu_j \\ \Sigma_j \end{pmatrix}$$

To make the dependence of  $p(\mathbf{x}|\omega_j)$  on  $\theta_j$  more specific, we can write  $p(\mathbf{x}|\omega_j, \theta_j)$ . Furthermore, to make our life easier – and why not – we assume that all  $\theta_j$  are independent, i.e., knowing one does not tell us anything about the others. This will allow us to work for each  $\theta_j$  separately, without worrying about interdependence of the parameters. We may therefore drop the class subscripts  $j$

# Maximum-Likelihood Estimation

- General principle
  - Assume we have  $c$  classes and
  - The samples in class  $j$  have been drawn according to the probability law  $p(\mathbf{x} \mid \omega_j)$   
 $p(\mathbf{x} \mid \omega_j) \sim \mathbf{N}(\mu_j, \Sigma_j)$   
 $p(\mathbf{x} \mid \omega_j) \equiv \mathbf{P}(\mathbf{x} \mid \omega_j, \theta_j)$  where:

# Maximum-Likelihood Estimation

- Our problem is to use the information provided by the training samples to **obtain good estimates for the unknown parameter** vectors associated with each category.



# Likelihood Function

- Use the information provided by the training samples to estimate  $\theta = (\theta_1, \theta_2, \dots, \theta_c)$ , each  $\theta_i$  ( $i = 1, 2, \dots, c$ ) is associated with each category
- Suppose that  $D$  contains  $n$  samples,  $x_1, x_2, \dots, x_n$

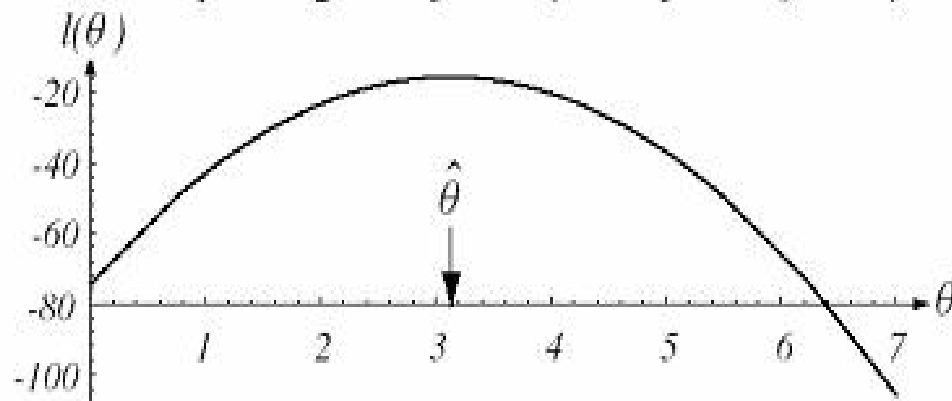
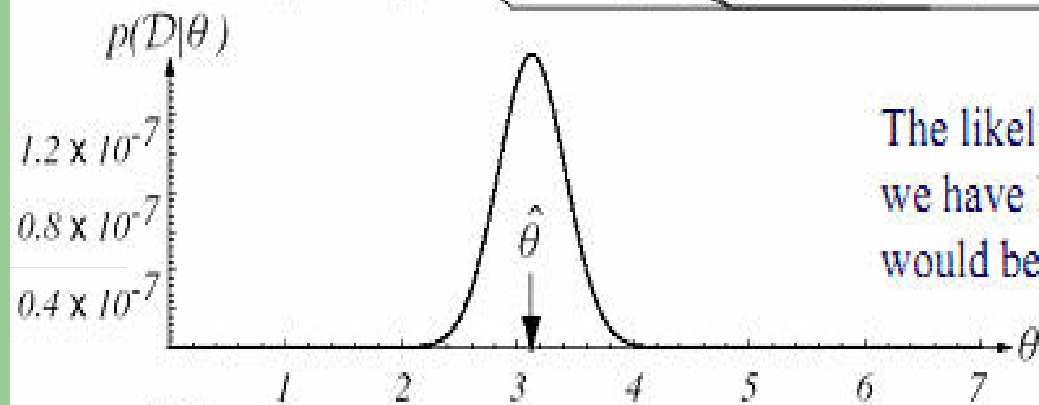
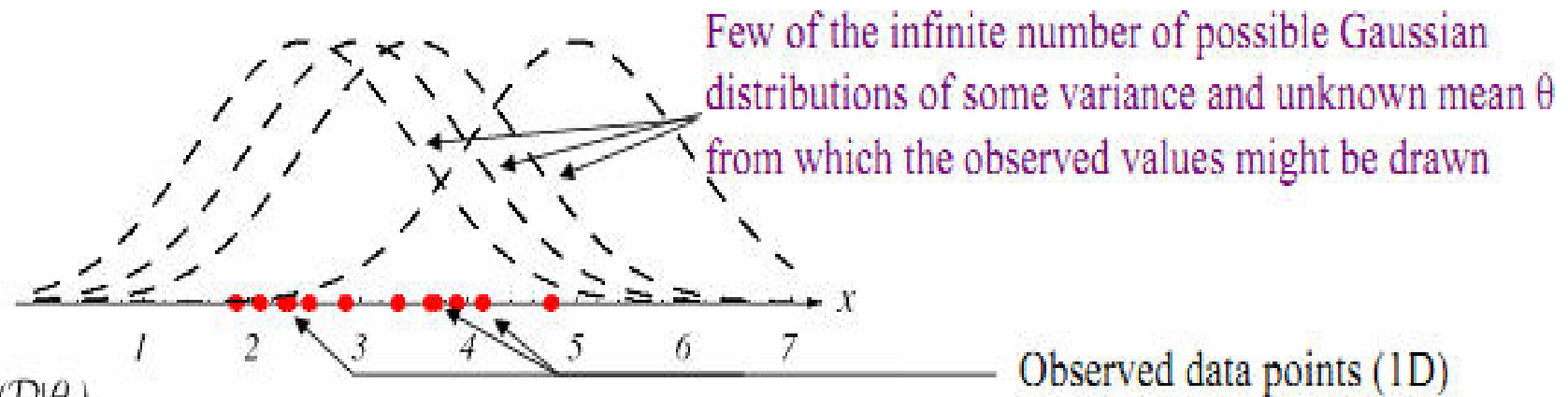
$$P(D | \theta) = \prod_{k=1}^{k=n} P(x_k | \theta) = F(\theta)$$

**$P(D | \theta)$  is called the likelihood of  $\theta$  w.r.t. the set of samples)**

Goal: find an estimate  $\hat{\theta}$

- Find  $\theta$  which maximizes  $P(D | \theta)$

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



# Maximize log likelihood function:

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

- Let  $\theta = (\theta_1, \theta_2, \dots, \theta_p)^t$  and let  $\nabla_\theta$  be the gradient operator

$$\nabla_\theta = \left[ \frac{\partial}{\partial \theta_1}, \frac{\partial}{\partial \theta_2}, \dots, \frac{\partial}{\partial \theta_p} \right]^t$$

- New problem statement: Determine  $\theta$  that maximizes the log-likelihood

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} l(\theta)$$

Set of necessary conditions for an optimum is:

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

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

- Example of a specific case: unknown  $\mu$

- $P(\mathbf{x}_i | \mu) \sim N(\mu, \Sigma)$

(Samples are drawn from a multivariate normal population)

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

$$\text{and } \nabla_{\theta\mu} \ln P(\mathbf{x}_k | \mu) = \Sigma^{-1} (\mathbf{x}_k - \mu)$$

$\theta = \mu$  therefore:

- The ML estimate for  $\mu$  must satisfy:

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

- Multiplying by  $\Sigma$  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!

**Conclusion:**

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

- ML Estimation:

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

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

$$l = \ln P(\mathbf{x}_k | \theta) = -\frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2} (\mathbf{x}_k - \theta_1)^2$$

$$\nabla_{\theta} l = \begin{pmatrix} \frac{\sigma}{\sigma\theta_1} (\ln P(\mathbf{x}_k | \theta)) \\ \frac{\sigma}{\sigma\theta_2} (\ln P(\mathbf{x}_k | \theta)) \end{pmatrix} = \mathbf{0}$$

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



Summation:

$$\left\{ \begin{array}{l} \sum_{k=1}^{k=n} \frac{1}{\hat{\theta}_2} (\mathbf{x}_k - \theta_1) = 0 \quad (1) \\ - \sum_{k=1}^{k=n} \frac{1}{\hat{\theta}_2} + \sum_{k=1}^{k=n} \frac{(\mathbf{x}_k - \hat{\theta}_1)^2}{\hat{\theta}_2^2} = 0 \quad (2) \end{array} \right.$$

Combining (1) and (2), one obtains:

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

### Example 1:

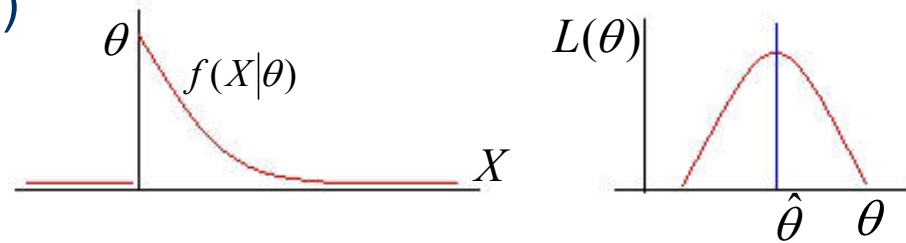
Consider an exponential distribution

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

(single feature, single parameter)

With a random sample

$$\{X_1, X_2, \dots, X_n\}$$



**Estimate  $\theta$  ?**

## Example 1:

Consider an exponential distribution

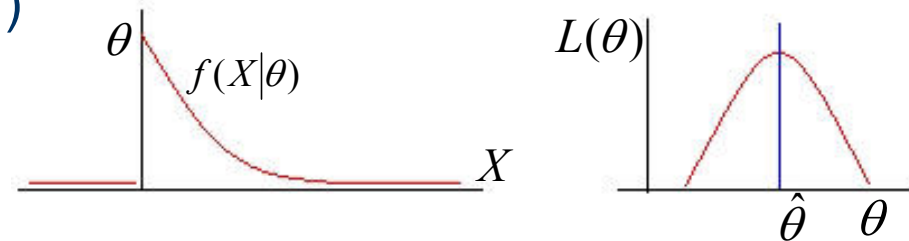
$$f(X; \theta) = \begin{cases} \theta e^{-\theta x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

(single feature, single parameter)

With a random sample

$$\{X_1, X_2, \dots, X_n\}$$

:



$$L(\theta) = f(X_1, X_2, \dots, X_n | \theta) = \prod_{i=1}^n \theta \cdot e^{-\theta \cdot x_i} \quad \text{valid for } x \geq 0$$

$$l(\theta) = \ln L(\theta) = \sum_{i=1}^n \ln \theta - \theta \sum_{i=1}^n x_i = n \ln \theta - \theta \sum_{i=1}^n x_i$$

$$\frac{dl}{d\theta} = \frac{d \ln L(\theta)}{d\theta} = \frac{n}{\theta} - \sum_{i=1}^n x_i = 0$$

$$\Rightarrow \frac{n}{\hat{\theta}} = \sum_{i=1}^n x_i \Rightarrow \hat{\theta} = \frac{1}{\frac{1}{n} \sum_{i=1}^n x_i} \quad (\text{inverse of average})$$

## Example 2:

Multivariate Gaussian with unknown mean vector  $M$ .

Assume  $\Sigma$  is known.

$k$  samples from the same distribution:

$$X_1, X_2, \dots, X_k \quad (\text{iid})$$

$$L(X | M) = \prod_{i=1}^k \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(X_i - M)^T \Sigma^{-1} (X_i - M)}$$

$$\nabla l = \nabla_M \log L = \sum_{i=1}^k \nabla_M \log \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(X_i - M)^T \Sigma^{-1} (X_i - M)}$$

$$= \sum_{i=1}^k \nabla_M \left( \frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (X_i - M)^T \Sigma^{-1} (X_i - M) \right)$$

$$= \sum_{i=1}^k (\Sigma^{-1} (X_i - \hat{M})) \quad (\text{linear algebra})$$

$$\Rightarrow 0 = \Sigma^{-1} \left( \sum_{i=1}^k X_i - k \hat{M} \right)$$

$$\hat{M} = \frac{1}{k} \sum_{i=1}^k X_i \quad \text{(sample average or sample mean)}$$

### Example 3:

Binary variables with unknown parameters  $p_i, 1 \leq i \leq n$   
(n parameters)

How to estimate these  $p_i, 1 \leq i \leq n$

Think about your homework ?

### Example 3:

Binary variables with unknown parameters  $p_i, 1 \leq i \leq n$   
(n parameters)

$$\log P(X) = \sum_{i=1}^n x_i \log p_i + \sum_{i=1}^n (1-x_i) \log(1-p_i)$$

So,

$$\begin{aligned} l = \log L &= \sum_{j=1}^k \log P(X_j) && \text{k samples} \\ &= \sum_{j=1}^k \left( \sum_{i=1}^n x_{ij} \log p_i + \sum_{i=1}^n (1-x_{ij}) \log(1-p_i) \right) \end{aligned}$$

here  $x_{ij}$  is the  $i^{th}$  element of  $j^{th}$  sample  $X_j$ .

**So,**

$$\nabla_{p_i} \log L = \begin{bmatrix} \frac{\partial}{\partial p_1} \log L \\ \frac{\partial}{\partial p_2} \log L \\ \vdots \\ \frac{\partial}{\partial p_n} \log L \end{bmatrix}$$

$$\frac{\partial}{\partial p_i} \log L = \sum_{j=1}^k \left( \frac{x_{ij}}{p_i} - ((1 - x_{ij})(1 - p_i)) \right)$$

$$\Rightarrow 0 = \frac{1}{\hat{p}_i} \sum_{j=1}^k x_{ij} - \frac{1}{1 - \hat{p}_i} \sum_{j=1}^k (1 - x_{ij})$$

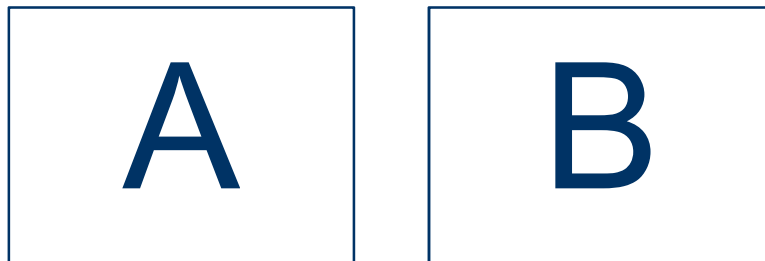
$$\Rightarrow \hat{p}_i = \frac{1}{k} \sum_{j=1}^k x_{ij}$$

**$\hat{p}_i$  is the sample average of the feature.**



Since  $X_i$  is binary,  $\sum_{j=1}^k x_{ij}$  will be the same as counting the occurrences of '1'.

Consider character recognition problem with binary matrices.



For each pixel, count the number of 1's and this is the estimate of  $p_i$  .