

---

**ECE471-571 – Pattern Recognition**

---

**Lecture 4 – Parametric Estimation**

---

Hairong Qi, Gonzalez Family Professor  
 Electrical Engineering and Computer Science  
 University of Tennessee, Knoxville  
<http://www.eecs.utk.edu/faculty/qi>  
 Email: hqi@utk.edu

---

---

---

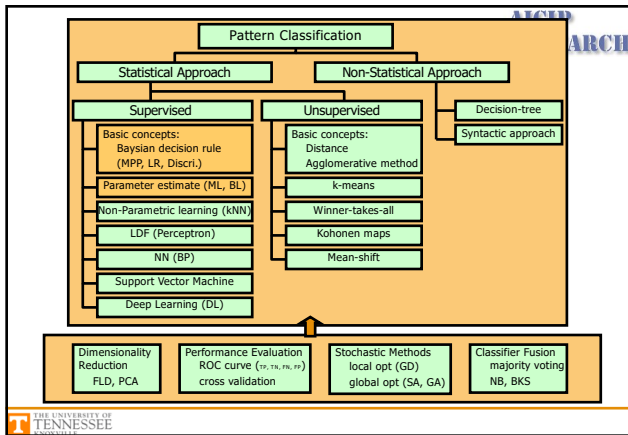
---

---

---

---

---




---

---

---


---

---

---

---

---



## Bayes Decision Rule

$$P(\omega_j | x) = \frac{p(x | \omega_j) P(\omega_j)}{p(x)}$$

**Maximum Posterior Probability** For a given  $x$ , if  $P(\omega_1 | x) > P(\omega_2 | x)$ , then  $x$  belongs to class 1, otherwise, 2.

**Likelihood Ratio** If  $\frac{p(x|\omega_1)}{p(x|\omega_2)} > \frac{P(\omega_2)}{P(\omega_1)}$ , then  $x$  belongs to class 1, otherwise, 2.

**Discriminant Function** The classifier will assign a feature vector  $x$  to class  $\omega_j$  if  $g_+(x) > g_-(x)$

Case 1: Minimum Euclidean Distance (Linear Machine),  $\Sigma_i = \sigma^2 I$

Case 2: Minimum Mahalanobis Distance (Linear Machine),  $\Sigma_i = \Sigma$

Case 3: Quadratic classifier,  $\Sigma_i = \text{arbitrary}$

---

---

---

---

---

---

---

---

## Multivariate Normal Density

AICIP  
RESEARCH

$$p(\vec{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right]$$

$\vec{x}$  : d - component column vector

$\vec{\mu}$  : d - component mean vector

$\Sigma$  : d - by - d covariance matrix

$|\Sigma|$  : determinant

$\Sigma^{-1}$  : inverse

---

---

---

---

---

---

---

---

---

---

## Estimating Normal Densities

AICIP  
RESEARCH

◆ Calculate  $\mu, \Sigma$

$$\vec{\mu}_i = \begin{bmatrix} \mu_{i1} = \frac{1}{n_i} \sum_{k=1}^{n_i} x_{k1} \\ \vdots \\ \mu_{id} = \frac{1}{n_i} \sum_{k=1}^{n_i} x_{kd} \end{bmatrix}$$

$$\Sigma_i = \begin{bmatrix} \sigma_{11} & \cdots & \sigma_{1d} \\ \vdots & \ddots & \vdots \\ \sigma_{d1} & \cdots & \sigma_{dd} \end{bmatrix} = \frac{1}{n_i - 1} \sum_{k=1}^{n_i} (\vec{x}_k - \vec{\mu}_i)(\vec{x}_k - \vec{\mu}_i)^T$$

---

---

---

---

---

---

---

---

---

---

## Covariance

AICIP  
RESEARCH

For  $d$  sets of variates denoted  $\{x_1\}, \dots, \{x_p\}, \dots, \{x_q\}, \dots, \{x_d\}$   
the covariance  $\sigma_{pq} = \text{cov}(x_p, x_q)$  of  $x_p$  and  $x_q$  is defined by

$$\begin{aligned} \text{cov}(x_p, x_q) &= E[(x_p - \mu_p)(x_q - \mu_q)] \\ &= E[x_p x_q] - E[x_p] E[x_q] + E[\mu_p \mu_q] \\ &= E[x_p x_q] - \mu_p E[x_q] + \mu_p \mu_q \\ &= E[x_p x_q] - \mu_q \mu_p - \mu_p \mu_q + \mu_p \mu_q \\ &= E[x_p x_q] - \mu_q \mu_p \end{aligned}$$

$$\begin{aligned} \text{When } p = q, \sigma_{pp} = \text{cov}(x_p, x_p) &= E[x_p x_p] - \mu_p \mu_p \\ &= E[x_p^2] - (E[x_p])^2 \\ &= \sigma_p^2 \end{aligned}$$

---

---

---

---

---

---

---

---

---

---

### \*Why? – Maximum Likelihood Estimation

$D = \{x_1, x_2, \dots, x_k, \dots, x_n\}$  is a data set of  $n$  training samples

Compare "likelihood"		$\bar{\theta} = \begin{bmatrix} \bar{\mu} \\ \Sigma \end{bmatrix}$
$p(x   \omega_i)$	$p(D   \bar{\theta})$	

$$p(D | \bar{\theta}) \xrightarrow{\text{assume samples are drawn independently}} \prod_{k=1}^n p(x_k | \bar{\theta})$$

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

$$\hat{\theta} = \arg \max_{\bar{\theta}} l(\bar{\theta})$$

---

---

---

---

---

---

---

---

---

---

### \*Derivation

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} \bar{\mu} \\ \Sigma \end{bmatrix}$$

$$p(\bar{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} (\bar{x} - \bar{\mu})^T \Sigma^{-1} (\bar{x} - \bar{\mu}) \right]$$

$$p(\bar{x}_k | \bar{\theta}) = \frac{1}{(2\pi)^{d/2} |\theta_2|^{1/2}} \exp \left[ -\frac{1}{2} (\bar{x}_k - \theta_1)^T \frac{1}{\theta_2} (\bar{x}_k - \theta_1) \right]$$

$$l(\bar{\theta}) = \sum_{k=1}^n \left( -\frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\theta_2| - \frac{1}{2} (\bar{x}_k - \theta_1)^T \frac{1}{\theta_2} (\bar{x}_k - \theta_1) \right)$$

$$\frac{\partial l}{\partial \theta_1} = 0$$

$$\frac{\partial l}{\partial \theta_2} = 0$$

---

---

---

---

---

---

---

---

---

---

### \* Derivative of a Quadratic Form

A matrix  $A$  is "positive definite" if  $x^T A x > 0 \quad \forall x \in R^d, x \neq 0$

$x^T A x$  is also called a "quadratic form".

The derivative of a quadratic form is particularly useful :

$$\frac{d}{dx} (x^T A x) = (A + A^T)x$$

---

---

---

---

---

---

---

---

---

---

## \*\*Why? – Bayesian Estimation

- Maximum likelihood estimation
  - The parameters are fixed
  - Find value for  $\theta$  that best agrees with or supports the actually observed training samples – likelihood of  $\theta$  w.r.t. the set of samples
- Bayesian estimation
  - Treat parameters as random variable themselves

$$p(D|\bar{\theta})$$

---

---

---

---

---

---

---

---

---

---

## \*\* The pdf of the parameter ( $\mu$ ) is Gaussian

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

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

$$p(\mu) = \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left[-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right]$$

---

---

---

---

---

---

---

---

---

---

## \*\* Derivation

$$p(\mu|D) = \frac{1}{C} \prod_{k=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x_k - \mu}{\sigma}\right)^2\right] \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left[-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right]$$

$$= \alpha \exp\left[-\frac{1}{2}\left(\sum_{k=1}^n \left(\frac{x_k - \mu}{\sigma}\right)^2 + \left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right)\right]$$

$$= \alpha \exp\left[-\frac{1}{2}\left(\frac{\sum_{k=1}^n x_k^2 - 2\mu \sum_{k=1}^n x_k + n\mu^2}{\sigma^2} + \frac{\mu^2 - 2\mu_0\mu + \mu_0^2}{\sigma_0^2}\right)\right]$$

$$= \beta \exp\left[-\frac{1}{2}\left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}\right)\mu^2 - 2\left(\frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2}\right)\mu\right]\right]$$

---


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

---

---

---

---

---

---

---

---

---

---

**\*\*  $\mu_n$  and  $\sigma_n$**

$$\mu_n = \left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \left( \frac{1}{n} \sum_{k=1}^n x_k \right) + \left( \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \right) \mu_0$$

Our best guess for  $\mu$  after observing  $n$  samples

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

Measures our uncertainty about this guess

◆ Behavior of Bayesian learning

- The larger the  $n$ , the smaller the  $\sigma_n$  – each additional observation decreases our uncertainty about the true value of  $\mu$
- As  $n$  approaches infinity,  $p(\mu|D)$  becomes more and more sharply peaked, approaching a Dirac delta function.
- $\mu_n$  is a linear combination between the sample mean and  $\mu_0$

---

---

---

---

---

---

---

---