

# ICS663: Pattern Recognition

Department of Information and Computer Sciences  
University of Hawai'i at Manoa

Kyungim Baek

ICS663 (Fall 2015)

1

## Announcement

- Homework assignment #1
  - Due: **Monday September 21, by 5:00 PM**
- Project proposal
  - Due: **Wednesday Septmber 23, by 5:00 PM**
  - Presentation (~ 10 minutes)
    - Monday (9/28)
      - BJ, Thomas, Tyson, Sharif, Danny, Jeremy
    - Wednesday (9/30)
      - Tetsuya, Kelly, Nurit
- **Exam I: Wednesday, October 7**

ICS663 (Fall 2015)

2

## Previously...

- **Parametric density estimation**
  - Estimating  $P(\omega_i)$  is not very difficult in the context of supervised learning
  - Estimating  $p(\mathbf{x} | \omega_i)$  is hard, especially in high-dimensional feature spaces and when the number of training samples available seems to be too small
    - Assumption:  $p(\mathbf{x} | \omega_i) = p(\mathbf{x} | \omega_i, \theta) - \text{known}$  parametric form,  $\exists$  a 'true'  $\theta_0$
  - Two approaches
    - Maximum-Likelihood estimation
    - Bayesian estimation

ICS663 (Fall 2015)

3

## Lecture 6

- Parametric Density Estimation
  - Maximum-Likelihood Estimation

ICS663 (Fall 2015)

4

## Previously...

- **Maximum-Likelihood (ML) estimation:** *parameters of probabilistic distributions are fixed but unknown values*
  - Parameters are unknown constants to be identified through training
  - Best estimate of parameter values is achieved by maximizing the probability of obtaining the samples observed

ICS663 (Fall 2015)

5

## Maximum-Likelihood Estimation

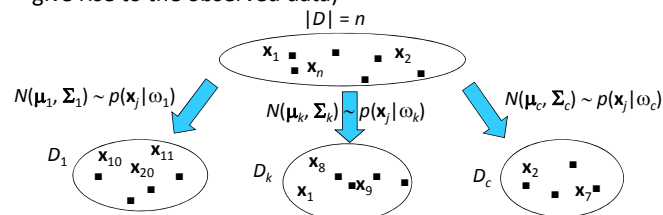
- General principle
  - Assume we have  $c$  classes and  $p(\mathbf{x} | \omega_i) \equiv p(\mathbf{x} | \omega_i, \theta_i)$ 
    - For example, if  $p(\mathbf{x} | \omega_i, \theta_i) \sim N(\mu_i, \Sigma_i)$  then
 
$$\theta_i = (\mu_i, \Sigma_i) = (\mu_i^1, \mu_i^2, \dots, \sigma_i^{11}, \sigma_i^{22}, \dots, \text{cov}(x_i^m, x_i^n), \dots)$$
  - 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
  - Assumption: independent and identically distributed (i.i.d.) samples

ICS663 (Fall 2015)

6

## ML Problem Statement

- Let  $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  and divided into  $c$  data sets according to the class as  $D_1, D_2, \dots, D_c$
- Our goal is to determine  $\theta_{ML}$  (value of  $\theta$  that is most likely to give rise to the observed data)



- Find  $\theta_i^{ML}$  such that  $\theta_i^{ML} = \max_{\theta_i} p(D_i | \theta_i) = \max_{\theta_i} \prod_{\mathbf{x}_k \in D_i} p(\mathbf{x}_k | \theta_i)$

ICS663 (Fall 2015)

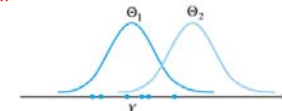
7

## ML Estimation (cont'd)

- Suppose that  $D$  contains  $n$  samples,  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ . The likelihood of  $\theta$  with respect to  $D$  is:

$$p(D | \theta) = \prod_{k=1}^n p(\mathbf{x}_k | \theta) = L(\theta)$$

- ML estimate of  $\theta$  is, by definition the value  $\theta_{ML}$  that maximizes  $p(D | \theta)$
- "It is the value of  $\theta$  that best agrees with the actually observed training sample"



ICS663 (Fall 2015)

8

## Log-Likelihood

- The value  $\theta$  of maximizing the likelihood function also maximizes its logarithm, known as the **log-likelihood** function

$$l(\theta) = \ln[p(D|\theta)] = \sum_{k=1}^n \ln[p(\mathbf{x}_k|\theta)]$$

- ML estimate

- 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  $\theta_{\text{ML}} = \arg \max_{\theta} l(\theta)$

- ML condition:

$$\nabla_{\theta} l(\theta) = \sum_{k=1}^n \nabla_{\theta} [\ln p(\mathbf{x}_k|\theta)] = 0$$

ICS663 (Fall 2015)

9

## ML Estimation Summary

$$P(\omega_i|\mathbf{x}) = P(\omega_i|\mathbf{x}, \theta_i) = \frac{p(\mathbf{x}|\omega_i, \theta_i)P(\omega_i|\theta_i)}{p(\mathbf{x}|\theta_i)}$$

$$p(\mathbf{x}|\hat{\theta})$$

maximize the likelihood:

$$\arg \max_{\theta} p(D|\theta)$$

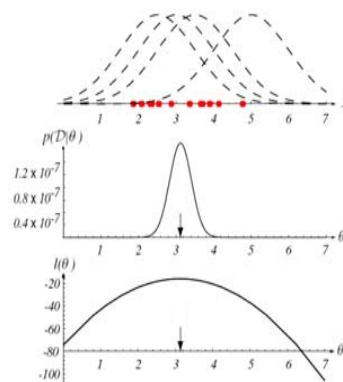
$$\prod_{k=1}^n p(\mathbf{x}_k|\theta) \Rightarrow \sum_{k=1}^n \ln p(\mathbf{x}_k|\theta) : \text{log-likelihood}$$

ICS663 (Fall 2015)

10

## Graphical View of ML Estimation

- Top graph shows several training points in one dimension, known or assumed to be drawn from a Gaussian of a particular variance, but unknown mean. Four of the infinite number of candidate source distributions are shown in dashed lines.
- The middle figure shows the likelihood  $p(D|\theta)$  as a function of the mean. If we had a very large number of training points, this likelihood would be very narrow. The value that maximizes the likelihood is marked by  $\theta_{\text{ML}}$ ; it also maximizes the logarithm of the likelihood shown at the bottom.
- Note: even though they look similar, the likelihood  $p(D|\theta)$  is shown as a function of  $\theta$  whereas the conditional density  $p(\mathbf{x}|\theta)$  is shown as a function of  $\mathbf{x}$ .



ICS663 (Fall 2015)

11

## The Gaussian Case: Unknown $\mu$

- Samples are drawn from a multivariate normal population:  
 $p(\mathbf{x}_i|\mu) \sim \mathcal{N}(\mu, \Sigma)$

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

$$\nabla_{\mu} \ln p(\mathbf{x}_k|\mu) = \Sigma^{-1} (\mathbf{x}_k - \mu)$$

$$\frac{\partial}{\partial \mathbf{x}} [\mathbf{y}'\mathbf{x}] = \mathbf{y}$$

$$\frac{\partial}{\partial \mathbf{x}} [\mathbf{x}'\mathbf{M}\mathbf{x}] = (\mathbf{M} + \mathbf{M}')\mathbf{x}$$

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

$$\sum_{k=1}^n \Sigma^{-1} (\mathbf{x}_k - \mu_{\text{ML}}) = 0$$

sample mean

- Multiplying by  $\Sigma$  and rearranging, we obtain:  $\mu_{\text{ML}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$

Just the arithmetic average of the training samples!

ICS663 (Fall 2015)

12

## Unknown $\mu$ and $\sigma$

- Univariate case:  $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$

$$\ln p(x_k | \theta) = -\frac{1}{2} \ln 2\pi - \frac{1}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} (x_k - \mu)^2$$

$$\nabla_{\theta} \ln p(x_k | \theta) = \begin{pmatrix} \frac{\partial}{\partial \mu} \ln p(x_k | \theta) \\ \frac{\partial}{\partial \sigma^2} \ln p(x_k | \theta) \end{pmatrix} = \begin{pmatrix} \frac{1}{\sigma^2} (x_k - \mu) \\ -\frac{1}{2\sigma^2} + \frac{(x_k - \mu)^2}{2(\sigma^2)^2} \end{pmatrix}$$

ICS663 (Fall 2015)

13

## Unknown $\mu$ and $\sigma$ (cont'd)

- The ML estimate for  $\mu$  and  $\sigma^2$  must satisfy:

$$\sum_{k=1}^n \frac{1}{\sigma_{ML}^2} (x_k - \mu_{ML}) = 0 \quad (1)$$

$$-\sum_{k=1}^n \frac{1}{\sigma_{ML}^2} + \sum_{k=1}^n \frac{(x_k - \mu_{ML})^2}{(\sigma_{ML}^2)^2} = 0 \quad (2)$$

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

$$\mu_{ML} = \frac{1}{n} \sum_{k=1}^n x_k \quad ; \quad \sigma_{ML}^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \mu_{ML})^2$$

ICS663 (Fall 2015)

14

## MLE: Multivariate Gaussians

$$\ln p(\mathbf{x}_k | \boldsymbol{\mu}) = -\frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{x}_k - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_k - \boldsymbol{\mu})$$

The log-likelihood function of  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  is:

$$\begin{aligned} l(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \sum_{k=1}^n \ln p(\mathbf{x}_k | \boldsymbol{\mu}) \\ &= -\frac{nd}{2} \ln 2\pi - \frac{n}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_k - \boldsymbol{\mu}) \\ &= -\frac{nd}{2} \ln 2\pi - \frac{n}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{k=1}^n (\mathbf{x}_k' \boldsymbol{\Sigma}^{-1} \mathbf{x}_k - 2\boldsymbol{\mu}' \boldsymbol{\Sigma}^{-1} \mathbf{x}_k) + n\boldsymbol{\mu}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \end{aligned}$$

ICS663 (Fall 2015)

15

## MLE: Multivariate Gaussians (cont'd)

$$\begin{aligned} \frac{\partial l(\boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}} &= \frac{\partial}{\partial \boldsymbol{\mu}} \left( -\frac{nd}{2} \ln 2\pi - \frac{n}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{k=1}^n (\mathbf{x}_k' \boldsymbol{\Sigma}^{-1} \mathbf{x}_k - 2\boldsymbol{\mu}' \boldsymbol{\Sigma}^{-1} \mathbf{x}_k) + n\boldsymbol{\mu}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right) \\ &= -\frac{1}{2} \left( -2\boldsymbol{\Sigma}^{-1} \sum_{k=1}^n \mathbf{x}_k + 2n\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right) \\ &= \mathbf{0} \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}} [\mathbf{y}' \mathbf{x}] &= \frac{\partial}{\partial \mathbf{x}} [\mathbf{x}' \mathbf{y}] = \mathbf{y} \\ \frac{\partial}{\partial \mathbf{x}} [\mathbf{x}' \mathbf{M} \mathbf{x}] &= (\mathbf{M} + \mathbf{M}') \mathbf{x} \end{aligned}$$

Therefore,  $\boldsymbol{\Sigma}^{-1} \sum_{k=1}^n \mathbf{x}_k = n\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$

This gives the ML solution,

$$\boldsymbol{\mu}_{ML} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

ICS663 (Fall 2015)

16

## MLE: Multivariate Gaussians (cont'd)

- For  $\Sigma$ , we need *trace* operator:  $tr(\mathbf{A}) = \sum_{i=1 \dots n} a_{ii}$
- Trace is invariant under cyclical permutations of matrix product:  $tr(\mathbf{ABC}) = tr(\mathbf{CAB}) = tr(\mathbf{BCA})$
- Derivatives of quadratic forms:
  - $\mathbf{x}^t \mathbf{A} \mathbf{x} = tr(\mathbf{x}^t \mathbf{A} \mathbf{x}) = tr(\mathbf{x} \mathbf{x}^t \mathbf{A})$   $\because \mathbf{x}^t \mathbf{A} \mathbf{x}$  is a scalar
  - Derivative of  $tr(\mathbf{BA}) = \mathbf{B}^t$

$$\frac{\partial}{\partial a_{ij}} tr(\mathbf{BA}) = \frac{\partial}{\partial a_{ij}} \sum_k \sum_l b_{kl} a_{lk} = b_{ji}$$

– Therefore

$$\frac{\partial}{\partial \mathbf{A}} \mathbf{x}^t \mathbf{A} \mathbf{x} = \frac{\partial}{\partial \mathbf{A}} tr(\mathbf{x} \mathbf{x}^t \mathbf{A}) = (\mathbf{x} \mathbf{x}^t)^t = \mathbf{x} \mathbf{x}^t$$

ICS663 (Fall 2015)

17

## MLE: Multivariate Gaussians (cont'd)

$$\begin{aligned} l(\boldsymbol{\mu}, \Sigma) &= -\frac{nd}{2} \ln 2\pi - \frac{n}{2} \ln |\Sigma| - \frac{1}{2} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu})^t \Sigma^{-1} (\mathbf{x}_k - \boldsymbol{\mu}) \\ &= -\frac{nd}{2} \ln 2\pi + \frac{n}{2} \ln |\Sigma^{-1}| - \frac{1}{2} \sum_{k=1}^n tr[(\mathbf{x}_k - \boldsymbol{\mu})^t \Sigma^{-1} (\mathbf{x}_k - \boldsymbol{\mu})] \\ &= -\frac{nd}{2} \ln 2\pi + \frac{n}{2} \ln |\Sigma^{-1}| - \frac{1}{2} \sum_{k=1}^n tr[(\mathbf{x}_k - \boldsymbol{\mu})(\mathbf{x}_k - \boldsymbol{\mu})^t \Sigma^{-1}] \end{aligned}$$

$|\Sigma^{-1}| = |\Sigma^{-1}|$

$$\frac{\partial l(\boldsymbol{\mu}, \Sigma)}{\partial \Sigma^{-1}} = \frac{n}{2} \Sigma - \frac{1}{2} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu})(\mathbf{x}_k - \boldsymbol{\mu})^t = \mathbf{0}$$

$\frac{\partial}{\partial \mathbf{A}} \ln |\mathbf{A}| = \mathbf{A}^{-t}$

This give the ML solution,

$$\Sigma_{ML} = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu}_{ML})(\mathbf{x}_k - \boldsymbol{\mu}_{ML})^t$$

ICS663 (Fall 2015)

18

## ML Analysis

- Geometrical view:** samples are a cloud of points whose center is the sample mean (employed as ML estimator)
- ML Bias:** ML estimator  $\sigma_{ML}^2$  for  $\sigma^2$  is **biased** over a small population of samples, but it is **asymptotically unbiased** as the number of training samples becomes very large

$$E[\sigma_{ML}^2] = E\left[\frac{1}{n} \sum_{k=1}^n (x_k - \mu_{ML})^2\right] = \frac{n-1}{n} \sigma^2 \neq \sigma^2$$

See the notes.

- The variance over a population of  $n$  training samples is:

$$\hat{\sigma}^2 = \frac{n}{n-1} \sigma_{ML}^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \mu_{ML})^2$$

unbiased estimator for  $\sigma^2$

- An elementary unbiased estimator for  $\Sigma$  is the **sample covariance matrix**:

$$\mathbf{C} = \frac{1}{n-1} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu}_{ML})(\mathbf{x}_k - \boldsymbol{\mu}_{ML})^t$$

ICS663 (Fall 2015)

19

## MLE Example Problem

- Suppose we toss a coin  $n$  times and observe the outcomes. We model the coin by a single parameter  $\theta$  that represents the probability of tossing heads. Given  $n$  independent observed tosses  $D = \{x_1, x_2, \dots, x_n\}$ , where  $x_i = 0$  if tail and  $x_i = 1$  if head, let  $n_H$  be the number of times the coin turned up heads. Then, the likelihood function is:

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

Show that the maximum likelihood estimate of  $\theta$  is

$$\hat{\theta}_{ML} = \frac{n_H}{n}.$$

ICS663 (Fall 2015)

20