

# Maximum-Likelihood & Bayesian Parameter Estimation (Sections 3.1-3.5)

Mingmin Chi

SCS Fudan University, Shanghai, China

## 1 Introduction

## 2 Maximum-Likelihood Estimation

- General Principle
- Gaussian Case: Unknown Mean  $\mu$
- Gaussian Case: Unknown  $\mu$  and  $\Sigma$
- Bias

## 3 Bayesian Estimation

- Parameter Distribution
- Gaussian Case

## 4 General Theory

- MLE vs. Bayes estimates

## 1 Introduction

## 2 Maximum-Likelihood Estimation

- General Principle
- Gaussian Case: Unknown Mean  $\mu$
- Gaussian Case: Unknown  $\mu$  and  $\Sigma$
- Bias

## 3 Bayesian Estimation

- Parameter Distribution
- Gaussian Case

## 4 General Theory

- MLE vs. Bayes estimates

- 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!

- 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 data set**

- No problem with prior estimation
- Samples are often too small for class-conditional estimation (large dimension of feature space!)

- 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 data set**
  - No problem with prior estimation
  - Samples are often too small for class-conditional estimation (large dimension of feature space!)
- If we know the number of parameters in advance, our general knowledge about the problem can be reduced by parameterizing the conditional density, e.g.,

- Data availability in a Bayesian framework
  - We could design an optimal classifier if we knew:
    - $P(\omega_i)$  (priors)
    - $P(\mathbf{x}|\omega_i)$  (class-conditional densities)

Unfortunately, we rarely have this complete information!

- Design a classifier from a **training data set**
  - No problem with prior estimation
  - Samples are often too small for class-conditional estimation (large dimension of feature space!)
- If we know the number of parameters in advance, our general knowledge about the problem can be reduced by parameterizing the conditional density, e.g.,
  - assuming that  $p(\mathbf{x}|\omega_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$
  - estimating an **unknown function**  $p(\mathbf{x}|\omega_i) \Rightarrow$

- Data availability in a Bayesian framework

- We could design an optimal classifier if we knew:

- $P(\omega_i)$  (priors)
- $P(\mathbf{x}|\omega_i)$  (class-conditional densities)

Unfortunately, we rarely have this complete information!

- Design a classifier from a **training data set**

- No problem with prior estimation
- Samples are often too small for class-conditional estimation (large dimension of feature space!)

- If we know the number of parameters in advance, our general knowledge about the problem can be reduced by parameterizing the conditional density, e.g.,

- assuming that  $p(\mathbf{x}|\omega_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$
- estimating an **unknown function**  $p(\mathbf{x}|\omega_i) \Rightarrow$  estimating the **parameters**  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\Sigma}_i$



# Ways of Parameter Estimation

- The problem of parameter estimation is a classical one in statistics and it can be approached in several ways
- Two common and reasonable procedures:

# Ways of Parameter Estimation

- The problem of parameter estimation is a classical one in statistics and it can be approached in several ways
- Two common and reasonable procedures:
  - Maximum-Likelihood (ML) estimation

# Ways of Parameter Estimation

- The problem of parameter estimation is a classical one in statistics and it can be approached in several ways
- Two common and reasonable procedures:
  - Maximum-Likelihood (ML) estimation
    - Parameters in ML estimation are **fixed but unknown**

# Ways of Parameter Estimation

- The problem of parameter estimation is a classical one in statistics and it can be approached in several ways
- Two common and reasonable procedures:
  - Maximum-Likelihood (ML) estimation
    - Parameters in ML estimation are **fixed but unknown**
    - Best parameters are obtained by **maximizing the probability of obtaining the samples available**

# Ways of Parameter Estimation

- The problem of parameter estimation is a classical one in statistics and it can be approached in several ways
- Two common and reasonable procedures:
  - Maximum-Likelihood (ML) estimation
    - Parameters in ML estimation are **fixed but unknown**
    - Best parameters are obtained by **maximizing the probability of obtaining the samples available**
  - Bayesian learning

# Ways of Parameter Estimation

- The problem of parameter estimation is a classical one in statistics and it can be approached in several ways
- Two common and reasonable procedures:
  - Maximum-Likelihood (ML) estimation
    - Parameters in ML estimation are **fixed but unknown**
    - Best parameters are obtained by **maximizing the probability of obtaining the samples available**
  - Bayesian learning
    - Parameters are **random variables** having some known distribution
    - **Training Data** convert the known distribution to a posterior density  $P(\theta|\mathbf{x}, y)$  and sharpen  $P(\hat{\theta}|\mathbf{x}, y)$
- In either approach, we use posterior density for our classification rule

# Supervised and Unsupervised Learning

## Common properties

- Samples  $\mathbf{x}$  are assumed to be obtained by

# Supervised and Unsupervised Learning

## Common properties

- Samples  $\mathbf{x}$  are assumed to be obtained by
  - selecting a state of nature  $\omega_i$  with probability  $P(\omega_i)$



# Supervised and Unsupervised Learning

## Common properties

- Samples  $\mathbf{x}$  are assumed to be obtained by
  - selecting a state of nature  $\omega_i$  with probability  $P(\omega_i)$
  - independently selecting  $\mathbf{x}$  according to the probability law  $p(\mathbf{x}|\omega_i)$

# Supervised and Unsupervised Learning

## Common properties

- Samples  $\mathbf{x}$  are assumed to be obtained by
  - selecting a state of nature  $\omega_i$  with probability  $P(\omega_i)$
  - independently selecting  $\mathbf{x}$  according to the probability law  $p(\mathbf{x}|\omega_i)$

## Distinction

- For supervised learning, we **know** the state of nature (class label) for each training sample
- For unsupervised learning, we **donot know** the state of nature (class label) for all training samples

## 1 Introduction

## 2 Maximum-Likelihood Estimation

- General Principle
- Gaussian Case: Unknown Mean  $\mu$
- Gaussian Case: Unknown  $\mu$  and  $\Sigma$
- Bias

## 3 Bayesian Estimation

- Parameter Distribution
- Gaussian Case

## 4 General Theory

- MLE vs. Bayes estimates

# Advantages

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

## 1 Introduction

## 2 Maximum-Likelihood Estimation

- General Principle

- Gaussian Case: Unknown Mean  $\mu$
- Gaussian Case: Unknown  $\mu$  and  $\Sigma$
- Bias

## 3 Bayesian Estimation

- Parameter Distribution
- Gaussian Case

## 4 General Theory

- MLE vs. Bayes estimates

# Setting

- Suppose that we separate a collection of samples  $\mathcal{D}$  according to class, so that we have  $c$  datasets,  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_c$
- The samples in  $\mathcal{D}_i$  have been drawn independently according to the probability law  $p(\mathbf{x}|\omega_i)$

# Setting

- Suppose that we separate a collection of samples  $\mathcal{D}$  according to class, so that we have  $c$  datasets,  $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_c$
- The samples in  $\mathcal{D}_i$  have been drawn independently according to the probability law  $p(\mathbf{x}|\omega_i)$
- We say such samples are **i.i.d. - independent and identically distributed** random variables

# Parameter Vector $\theta_i$ for $p(\mathbf{x}|\omega_i)$

- Assume that  $p(\mathbf{x}|\omega_i)$  has a known parametric form, determined uniquely by the value of a **parameter vector  $\theta_i$**
- E.g.,  $p(\mathbf{x}|\omega_i) \sim \mathcal{N}(\mu_i, \Sigma_i)$ , where

$$\theta_i =$$



# Parameter Vector $\theta_i$ for $p(\mathbf{x}|\omega_i)$

- Assume that  $p(\mathbf{x}|\omega_i)$  has a known parametric form, determined uniquely by the value of a **parameter vector  $\theta_i$**
- E.g.,  $p(\mathbf{x}|\omega_i) \sim \mathcal{N}(\mu_i, \Sigma_i)$ , where

$$\theta_i = (\mu_i, \Sigma_i)$$

# Parameter Vector $\theta_i$ for $p(\mathbf{x}|\omega_i)$

- Assume that  $p(\mathbf{x}|\omega_i)$  has a known parametric form, determined uniquely by the value of a **parameter vector  $\theta_i$**
- E.g.,  $p(\mathbf{x}|\omega_i) \sim \mathcal{N}(\mu_i, \Sigma_i)$ , where

$$\theta_i = (\mu_i, \Sigma_i)$$

- To show the dependence of  $p(\mathbf{x}|\omega_i)$  on  $\theta_i$

$$p(\mathbf{x}|\omega_i) \propto$$

# Parameter Vector $\theta_i$ for $p(\mathbf{x}|\omega_i)$

- Assume that  $p(\mathbf{x}|\omega_i)$  has a known parametric form, determined uniquely by the value of a **parameter vector  $\theta_i$**
- E.g.,  $p(\mathbf{x}|\omega_i) \sim \mathcal{N}(\mu_i, \Sigma_i)$ , where

$$\theta_i = (\mu_i, \Sigma_i)$$

- To show the dependence of  $p(\mathbf{x}|\omega_i)$  on  $\theta_i$

$$p(\mathbf{x}|\omega_i) \propto p(\mathbf{x}|\omega_i, \theta_i)$$

- The target of ML estimation is to

# Parameter Vector $\theta_i$ for $p(\mathbf{x}|\omega_i)$

- Assume that  $p(\mathbf{x}|\omega_i)$  has a known parametric form, determined uniquely by the value of a **parameter vector**  $\theta_i$
- E.g.,  $p(\mathbf{x}|\omega_i) \sim \mathcal{N}(\mu_i, \Sigma_i)$ , where

$$\theta_i = (\mu_i, \Sigma_i)$$

- To show the dependence of  $p(\mathbf{x}|\omega_i)$  on  $\theta_i$

$$p(\mathbf{x}|\omega_i) \propto p(\mathbf{x}|\omega_i, \theta_i)$$

- The target of ML estimation is to use the information provided by the **training samples** to obtain good estimates for the unknown parameter vectors  $\theta_1, \theta_2, \dots, \theta_c$  for each category

# Data Independent Assumption

- To simplify treatment of the problem, we shall assume that if  $i \neq j$ , the samples in  $\mathcal{D}_i$  give no information about  $\theta_j$

# Data Independent Assumption

- To simplify treatment of the problem, we shall assume that if  $i \neq j$ , the samples in  $\mathcal{D}_i$  give no information about  $\theta_j$



the parameters for the different classes are functionally independent

# Data Independent Assumption

- To simplify treatment of the problem, we shall assume that if  $i \neq j$ , the samples in  $\mathcal{D}_i$  give no information about  $\theta_j$



the parameters for the different classes are functionally independent

- We can work with each category separately and simplify our notation by deleting indications of class distinctions

Use a set  $\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_c\}$  of training samples drawn independently from the probability density  $p(\mathbf{x}|\theta)$  to estimate the unknown parameter vector  $\theta$

# Likelihood of Parameter Vector $\theta$

Suppose  $\mathcal{D} = (\mathbf{x}_i)_{i=1}^n$

Samples drawn independently

$$p(\mathcal{D}|\theta) = \prod_{i=1}^n p(\mathbf{x}_i|\theta)$$

which is called the **likelihood of  $\theta$**  with respect to the set of samples, if viewed as a function of  $\theta$

- The **maximum-likelihood estimate** of  $\theta$  is the value  $\hat{\theta}$  that maximizes  $p(\mathcal{D}|\theta)$



# Likelihood of Parameter Vector $\theta$

Suppose  $\mathcal{D} = (\mathbf{x}_i)_{i=1}^n$

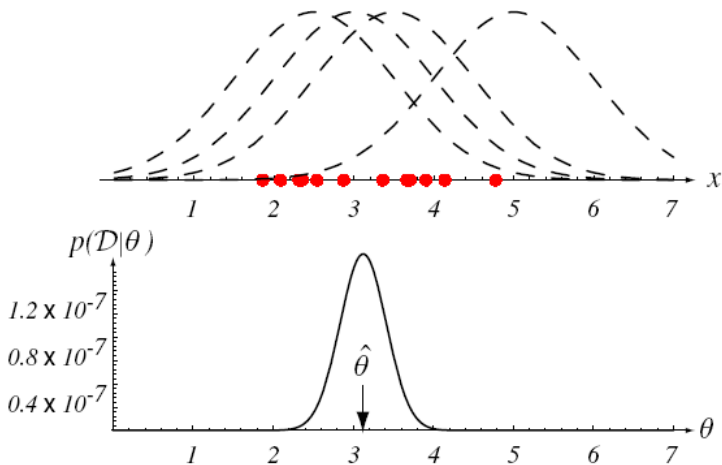
Samples drawn independently

$$p(\mathcal{D}|\theta) = \prod_{i=1}^n p(\mathbf{x}_i|\theta)$$

which is called the **likelihood of  $\theta$**  with respect to the set of samples, if viewed as a function of  $\theta$

- The **maximum-likelihood estimate** of  $\theta$  is the value  $\hat{\theta}$  that maximizes  $p(\mathcal{D}|\theta)$
- It is the value of  $\theta$  that agrees with or supports **the actually observed training samples**

$$p(\mathcal{D}|\hat{\theta}), \theta = \mu$$



# Optimal Estimation (1)

- Let  $\boldsymbol{\theta} = [\theta_1 \dots \theta_d]^\top$  and let  $\nabla_{\boldsymbol{\theta}}$  be the gradient operator

$$\nabla_{\boldsymbol{\theta}} =$$

# Optimal Estimation (1)

- Let  $\theta = [\theta_1 \dots \theta_d]^\top$  and let  $\nabla_\theta$  be the gradient operator

$$\nabla_\theta = \left[ \frac{\partial}{\partial \theta_1} \dots \frac{\partial}{\partial \theta_d} \right]^\top$$

- Define  $l(\theta)$  as the **log-likelihood** function

$$l(\theta) \equiv \ln p(\mathcal{D}|\theta)$$

- New problem statement: to determine the argument  $\theta$  that maximizes the log-likelihood

$$\hat{\theta} =$$

# Optimal Estimation (1)

- Let  $\theta = [\theta_1 \dots \theta_d]^\top$  and let  $\nabla_\theta$  be the gradient operator

$$\nabla_\theta = \left[ \frac{\partial}{\partial \theta_1} \dots \frac{\partial}{\partial \theta_d} \right]^\top$$

- Define  $l(\theta)$  as the **log-likelihood** function

$$l(\theta) \equiv \ln p(\mathcal{D}|\theta)$$

- New problem statement: to determine the argument  $\theta$  that maximizes the log-likelihood

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

# Optimal Estimation (2)

- The dependence on the dataset  $\mathcal{D}$  is implicit
- since  $l(\theta) \equiv \ln p(\mathcal{D}|\theta)$ ,

$$l(\theta) = \sum_{i=1}^n \ln p(\mathbf{x}_i|\theta)$$

and

$$\nabla_{\theta} l(\theta) = \sum_{i=1}^n \nabla_{\theta} \ln p(\mathbf{x}_i|\theta)$$

- A set of necessary conditions for the maximum likelihood estimate for  $\theta$  can be obtained from the set of  $d$  equations

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

# Optimal Estimation (3)

- Representing a true global maximum, or a local maximum

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

- second derivatives to check the global minimum



$$\lim_{n \rightarrow \infty} \hat{\theta} \equiv \theta$$

# MAP

- $p(\theta)$ , the prior probability of different parameter values
- maximum a posteriori (MAP)

$$\max_{\theta} l(\theta)p(\theta)$$

- MAP = MLE + uniform or “flat” prior
- drawback of MAP



## 1 Introduction

## 2 Maximum-Likelihood Estimation

- General Principle
- **Gaussian Case: Unknown Mean  $\mu$**
- Gaussian Case: Unknown  $\mu$  and  $\Sigma$
- Bias

## 3 Bayesian Estimation

- Parameter Distribution
- Gaussian Case

## 4 General Theory

- MLE vs. Bayes estimates

# Sample mean

- Samples are drawn from a multivariate normal population, i.e.,  
 $p(\mathbf{x}_i|\omega) \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$
- For simplicity, consider only the mean is unknown, i.e.,  $\theta =$

# Sample mean

- Samples are drawn from a multivariate normal population, i.e.,  $p(\mathbf{x}_i|\omega) \sim \mathcal{N}(\mu, \Sigma)$
- For simplicity, consider only the mean is unknown, i.e.,  $\theta = \mu$

$$\begin{aligned}\ln p(\mathbf{x}_i|\theta) &= \ln p(\mathbf{x}_i|\mu) \\ &= -\frac{1}{2}(\mathbf{x}_i - \mu)^\top \Sigma^{-1}(\mathbf{x}_i - \mu) - \frac{1}{2} \ln[(2\pi)^d |\Sigma|]\end{aligned}$$

and

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

- The ML estimate for  $\mu$  must satisfy

$$\sum_{i=1}^n \Sigma^{-1}(\mathbf{x}_i - \hat{\mu}) = 0 \Rightarrow$$

# Sample mean

- Samples are drawn from a multivariate normal population, i.e.,  $p(\mathbf{x}_i|\omega) \sim \mathcal{N}(\mu, \Sigma)$
- For simplicity, consider only the mean is unknown, i.e.,  $\theta = \mu$

$$\begin{aligned}\ln p(\mathbf{x}_i|\theta) &= \ln p(\mathbf{x}_i|\mu) \\ &= -\frac{1}{2}(\mathbf{x}_i - \mu)^\top \Sigma^{-1}(\mathbf{x}_i - \mu) - \frac{1}{2} \ln[(2\pi)^d |\Sigma|]\end{aligned}$$

and

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

- The ML estimate for  $\mu$  must satisfy

$$\sum_{i=1}^n \Sigma^{-1}(\mathbf{x}_i - \hat{\mu}) = 0 \Rightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

## 1 Introduction

## 2 Maximum-Likelihood Estimation

- General Principle
- Gaussian Case: Unknown Mean  $\mu$
- Gaussian Case: Unknown  $\mu$  and  $\Sigma$
- Bias

## 3 Bayesian Estimation

- Parameter Distribution
- Gaussian Case

## 4 General Theory

- MLE vs. Bayes estimates

# Univariate Gaussian

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

# Univariate Gaussian

- $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$
- The log-likelihood of a single point is

$$l(\theta) = \ln p(x_i|\theta) = -\frac{1}{2\theta_2}(x_i - \theta_1)^2 - \frac{1}{2} \ln 2\pi\theta_2$$

and its derivative is

$$\nabla_{\theta} l(\theta)$$

# Univariate Gaussian

- $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$
- The log-likelihood of a single point is

$$l(\theta) = \ln p(x_i|\theta) = -\frac{1}{2\theta_2}(x_i - \theta_1)^2 - \frac{1}{2} \ln 2\pi\theta_2$$

and its derivative is

$$\nabla_{\theta} l(\theta) = \begin{bmatrix} \nabla_{\theta_1} l(\theta) \\ \nabla_{\theta_2} l(\theta) \end{bmatrix} = \begin{bmatrix} \frac{1}{\theta_2}(x_i - \theta_1) \\ -\frac{1}{2\theta_2} + \frac{(x_i - \theta_1)^2}{2\theta_2^2} \end{bmatrix}$$

to the full log-likelihood



# Univariate Gaussian

- $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$
- The log-likelihood of a single point is

$$l(\theta) = \ln p(x_i|\theta) = -\frac{1}{2\theta_2}(x_i - \theta_1)^2 - \frac{1}{2} \ln 2\pi\theta_2$$

and its derivative is

$$\nabla_{\theta} l(\theta) = \begin{bmatrix} \nabla_{\theta_1} l(\theta) \\ \nabla_{\theta_2} l(\theta) \end{bmatrix} = \begin{bmatrix} \frac{1}{\theta_2}(x_i - \theta_1) \\ -\frac{1}{2\theta_2} + \frac{(x_i - \theta_1)^2}{2\theta_2^2} \end{bmatrix}$$

to the full log-likelihood

$$\sum_{i=1}^n \nabla_{\theta} \ln p(x_i|\theta) = 0 \Rightarrow \begin{cases} \sum_{i=1}^n \frac{1}{\hat{\theta}_2}(x_i - \hat{\theta}_1) = 0 \\ \sum_{i=1}^n -\frac{1}{2\hat{\theta}_2} + \frac{(x_i - \hat{\theta}_1)^2}{2\hat{\theta}_2^2} = 0 \end{cases}$$

## Univariate normal case

- we can obtain the following maximum-likelihood estimates for  $\mu$  and  $\sigma^2$

$$\begin{cases} \hat{\mu} = \hat{\theta}_1 = \frac{1}{n} \sum_{i=1}^n x_i \\ \hat{\sigma}^2 = \hat{\theta}_2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 \end{cases}$$

## Multivariate normal case

- With similar analysis, we can obtain the following maximum-likelihood estimates for  $\mu$  and  $\Sigma^2$

$$\begin{cases} \hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \\ \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu})(\mathbf{x}_i - \hat{\mu})^\top \end{cases}$$

## 1 Introduction

## 2 Maximum-Likelihood Estimation

- General Principle
- Gaussian Case: Unknown Mean  $\mu$
- Gaussian Case: Unknown  $\mu$  and  $\Sigma$
- Bias

## 3 Bayesian Estimation

- Parameter Distribution
- Gaussian Case

## 4 General Theory

- MLE vs. Bayes estimates

# Mean

- **Bias** of an estimator  $\hat{\theta}$ :

$$\mathcal{E}_{\theta} \left[ (f(\hat{\theta}; \mathcal{D}) - f(\theta; \mathcal{D}))^2 \right]$$

- ML estimate  $\hat{\mu}$  of the mean  $\mu$  is **unbiased**

$$\mathcal{E}[\hat{\mu}] = \mathcal{E} \left[ \frac{1}{n} \sum_{i=1}^n x_i \right] = \mu$$

# Variance

- ML estimate for the variance  $\sigma^2$  is **biased**

$$\mathcal{E} \left[ \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 \right] = \frac{n-1}{n} \sigma^2 \neq \sigma^2$$

which is asymptotically unbiased

- An elementary unbiased estimator for  $\hat{\Sigma}$  is given by

$$\mathbf{C} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top$$

Sample covariance matrix

which is absolutely unbiased

## 1 Introduction

## 2 Maximum-Likelihood Estimation

- General Principle
- Gaussian Case: Unknown Mean  $\mu$
- Gaussian Case: Unknown  $\mu$  and  $\Sigma$
- Bias

## 3 Bayesian Estimation

- Parameter Distribution
- Gaussian Case

## 4 General Theory

- MLE vs. Bayes estimates

# Class-Conditional Densities

- In MLE,  $\theta$  was supposed fixed
- In BE,  $\theta$  is a random variable
- The computation of posterior probabilities  $P(\omega_i|\mathbf{x})$  lies at the heart of Bayesian classification
- if prior probabilities and class densities are unknown, how to compute a posteriori?
- Goal: Given the sample set  $\mathcal{D}$ , compute  $P(\omega_i|\mathbf{x}, \mathcal{D})$
- Bayes formula can be written

$$P(\omega_i|\mathbf{x}, \mathcal{D}) =$$

# Class-Conditional Densities

- In MLE,  $\theta$  was supposed fixed
- In BE,  $\theta$  is a random variable
- The computation of posterior probabilities  $P(\omega_i|\mathbf{x})$  lies at the heart of Bayesian classification
- if prior probabilities and class densities are unknown, how to compute a posteriori?
- Goal: Given the sample set  $\mathcal{D}$ , compute  $P(\omega_i|\mathbf{x}, \mathcal{D})$
- Bayes formula can be written

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



# Class-Conditional Densities

- In MLE,  $\theta$  was supposed fixed
- In BE,  $\theta$  is a random variable
- The computation of posterior probabilities  $P(\omega_i|\mathbf{x})$  lies at the heart of Bayesian classification
- if prior probabilities and class densities are unknown, how to compute a posteriori?
- Goal: Given the sample set  $\mathcal{D}$ , compute  $P(\omega_i|\mathbf{x}, \mathcal{D})$
- Bayes formula can be written

$$\begin{aligned} P(\omega_i|\mathbf{x}, \mathcal{D}) &= \frac{p(\mathbf{x}|\omega_i, \mathcal{D})P(\omega_i|\mathcal{D})}{\sum_{j=1}^C p(\mathbf{x}|\omega_j, \mathcal{D})P(\omega_j|\mathcal{D})} \\ &\propto \frac{p(\mathbf{x}|\omega_i, \mathcal{D}_i)P(\omega_i)}{\sum_{j=1}^C p(\mathbf{x}|\omega_j, \mathcal{D}_j)P(\omega_j)} \end{aligned}$$

## 1 Introduction

## 2 Maximum-Likelihood Estimation

- General Principle
- Gaussian Case: Unknown Mean  $\mu$
- Gaussian Case: Unknown  $\mu$  and  $\Sigma$
- Bias

## 3 Bayesian Estimation

- Parameter Distribution
- Gaussian Case

## 4 General Theory

- MLE vs. Bayes estimates

# Central problem of Bayesian learning

- Desired probability density  $p(\mathbf{x})$  is unknown, but has a known parametric form  $p(\mathbf{x}|\theta)$
  - The only thing assumed unknown is the value of a parameter vector  $\theta$
- 
- known prior density  $p(\theta)$
  - $\mathcal{D}$  converts the prior to a posterior density  $p(\theta|\mathcal{D})$

Learning a probability density function  $\Rightarrow$  estimating a parameter vector

$p(\mathbf{x}),$

$$p(\mathbf{x}), p(\mathbf{x}|\mathcal{D}),$$

$$p(\mathbf{x}), p(\mathbf{x}|\mathcal{D}), p(\mathbf{x}, \theta|\mathcal{D})$$

$$p(\mathbf{x}), p(\mathbf{x}|\mathcal{D}), p(\mathbf{x}, \theta|\mathcal{D})$$

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

$$p(\mathbf{x}), p(\mathbf{x}|\mathcal{D}), p(\mathbf{x}, \theta|\mathcal{D})$$

$$\begin{aligned} p(\mathbf{x}|\mathcal{D}) &= \int p(\mathbf{x}, \theta|\mathcal{D}) d\theta \\ &= \int p(\mathbf{x}|\theta, \mathcal{D}) p(\theta|\mathcal{D}) d\theta \end{aligned}$$



$$p(\mathbf{x}), p(\mathbf{x}|\mathcal{D}), p(\mathbf{x}, \theta|\mathcal{D})$$

$$\begin{aligned} p(\mathbf{x}|\mathcal{D}) &= \int p(\mathbf{x}, \theta|\mathcal{D}) d\theta \\ &= \int p(\mathbf{x}|\theta, \mathcal{D}) p(\theta|\mathcal{D}) d\theta \\ &= \int p(\mathbf{x}|\theta) p(\theta|\mathcal{D}) d\theta \end{aligned}$$

$$p(\mathbf{x}), p(\mathbf{x}|\mathcal{D}), p(\mathbf{x}, \theta|\mathcal{D})$$

$$\begin{aligned} p(\mathbf{x}|\mathcal{D}) &= \int p(\mathbf{x}, \theta|\mathcal{D}) d\theta \\ &= \int p(\mathbf{x}|\theta, \mathcal{D}) p(\theta|\mathcal{D}) d\theta \\ &= \int p(\mathbf{x}|\theta) p(\theta|\mathcal{D}) d\theta \end{aligned}$$

- which links the desired class-conditional density  $p(\mathbf{x}|\mathcal{D})$  to the posterior density  $p(\theta|\mathcal{D})$  for the unknown parameter
- which can be performed numerically, e.g., by Monte-Carlo simulation
- if  $p(\theta|\mathcal{D})$  peaks very sharply about some value  $\hat{\theta}$ ,  $p(\mathbf{x}|\mathcal{D}) \simeq p(\mathbf{x}|\hat{\theta})$

## 1 Introduction

## 2 Maximum-Likelihood Estimation

- General Principle
- Gaussian Case: Unknown Mean  $\mu$
- Gaussian Case: Unknown  $\mu$  and  $\Sigma$
- Bias

## 3 Bayesian Estimation

- Parameter Distribution
- Gaussian Case

## 4 General Theory

- MLE vs. Bayes estimates

Use the BE to calculate the *a posteriori* density  $p(\theta|\mathcal{D})$  and the desired probability density  $p(\mathbf{x}|\mathcal{D})$  for the case where  $p(\mathbf{x}|\theta) \sim \mathcal{N}(\mu, \Sigma)$

# Univariate Case: $p(\mu|\mathcal{D})$ - Prior $p(\mu)$

Consider the case where  $\mu$  is the only unknown parameter. For the simplicity, we treat first the univariate case, i.e.,

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

where the only unknown quantity is the mean  $\mu$

- Assume that prior knowledge we might have about  $\mu$  can be expressed by a known prior density  $p(\mu)$
- Assume that

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

where  $\mu_0$  and  $\sigma_0^2$  are **known**

# Univariate Case: $p(\mu|\mathcal{D})$

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

# Univariate Case: $p(\mu|\mathcal{D})$

$$\begin{aligned} p(\mu|\mathcal{D}) &= \frac{p(\mathcal{D}|\mu)p(\mu)}{\int p(\mathcal{D}|\mu)p(\mu)d\mu} \\ &= \alpha \prod_{i=1}^n p(\mathbf{x}_i|\mu)p(\mu) \end{aligned}$$

where  $\alpha$  is a normalization factor that depends on  $\mathcal{D}$  but is independent of  $\mu$

$$p(x|\mu) \sim \mathcal{N}(\mu, \sigma^2) \text{ and } p(\mu) \sim \mathcal{N}(\mu_0, \sigma_0^2)$$

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

# Bayesian Learning

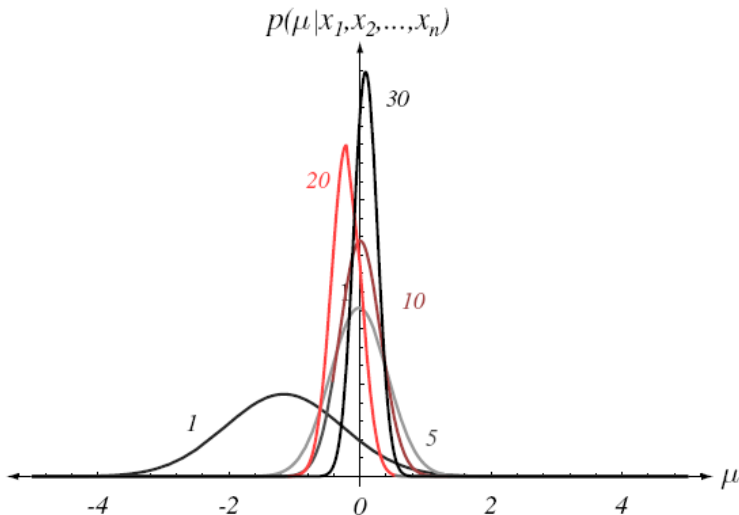
$$p(\mu|\mathcal{D}) = \frac{1}{\sqrt{2\pi}\sigma_n} \exp \left[ -\frac{1}{2} \left( \frac{\mu - \mu_n}{\sigma_n} \right)^2 \right]$$
$$\mu_n = \left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \hat{\mu}_n \right) + \left( \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0 \right)$$
$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}$$

- $\mu_n$  represents our best guess for  $\mu$  after observing  $n$  samples
- $\sigma_n^2$  measures our uncertainty about this guess
- As  $n$  increases,  $p(\mu|\mathcal{D})$  becomes more and more sharply peaked, approaching a Dirac delta function as  $n$  approaches infinity

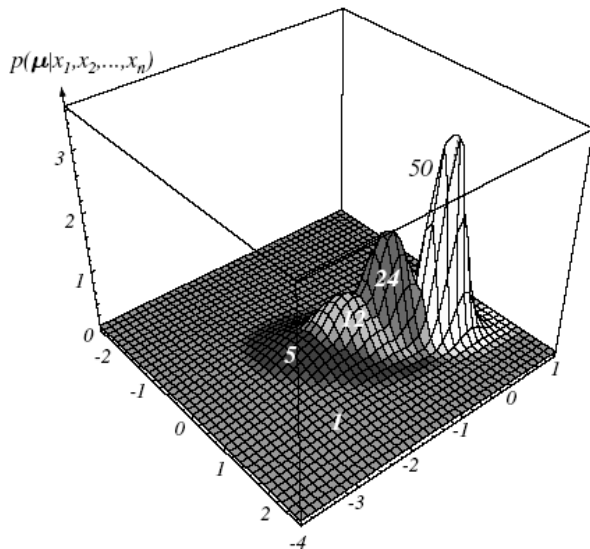
This is known as **Bayesian learning**



# Bayesian Learning: Visualization in 1-D



# Bayesian Learning: Visualization in 2-D



# Univariate Case: $p(x|\mathcal{D})$

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

Finally, we can obtain

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

Therefore:  $P(x|\omega_i, \mathcal{D}_i)$  together with  $P(\omega_i)$  and using Bayes formula, we obtain the Bayesian classification rule:

$$\max_{\omega_i} [P(\omega_i|x, \mathcal{D}_i)] = \max_{\omega_i} [P(x|\omega_i, \mathcal{D}_i) \cdot P(\omega_i)]$$

# Univariate Case: $p(x|\mathcal{D})$

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

Finally, we can obtain

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

Therefore:  $P(x|\omega_i, \mathcal{D}_i)$  together with  $P(\omega_i)$  and using Bayes formula, we obtain the Bayesian classification rule:

$$\max_{\omega_i} [P(\omega_i|x, \mathcal{D}_i)] = \max_{\omega_i} [P(x|\omega_i, \mathcal{D}_i) \cdot P(\omega_i)]$$

# Conjugate prior

## definition

- for a given probability distribution  $p(\mathbf{x}|\theta)$ , we can seek a prior  $p(\theta)$  that is conjugate to the likelihood function, so that the posterior distribution has the same functional form as the prior
- e.g., Gaussian case, for unknown  $\theta = \mu$

$$p(x_i|\mu) \sim \mathcal{N}(\mu, \sigma^2)$$

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

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

pdf	conjugate prior
Multinomial	Dirichlet
Binomial	Beta

# Bernoulli Case

- consider  $P(\mathbf{x}|\theta) = \text{Bern}(\mathbf{x}|\mu) = \mu^{\mathbf{x}}(1 - \mu)^{1-\mathbf{x}}$ , where  $\theta = \mu$  is the unknown parameter
- a conjugate prior distribution  $p(\mu) = \text{Beta}(\alpha, \beta)$  where  $\alpha, \beta$  are both known
- $p(\mu|\mathcal{D}) = \text{Beta}(\alpha + \sum_{i=1}^n \mathbf{x}_i, \beta + n - \sum_{i=1}^n \mathbf{x}_i)$

## 1 Introduction

## 2 Maximum-Likelihood Estimation

- General Principle
- Gaussian Case: Unknown Mean  $\mu$
- Gaussian Case: Unknown  $\mu$  and  $\Sigma$
- Bias

## 3 Bayesian Estimation

- Parameter Distribution
- Gaussian Case

## 4 General Theory

- MLE vs. Bayes estimates

- So far, we see how the BE can be used to obtain the desired density  $p(\mathbf{x}|\mathcal{D})$  in a special case-uni/multi- variate Gaussian
- This approach can be generalized to apply to any situation in which unknown density can be parameterized

## Basic Assumption

- Form of the density  $p(\mathbf{x}|\theta)$  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  $\mathcal{D}$  of  $n$  random variables  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  that follows  $P(\mathbf{x})$



The basic problem is to compute the posterior density  $P(\theta|\mathcal{D})$  to derive  $P(\mathbf{x}|\mathcal{D})$

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

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

$$p(\mathcal{D}|\theta) = \prod_{i=1}^n p(\mathbf{x}_i|\theta)$$

## Problems

- Computation difficulties
- convergence of  $p(\mathbf{x}|\mathcal{D}) \rightarrow p(\mathbf{x})$ ?

# Recursive Bayes learning

## Problems

- Provided  $\mathcal{D}^n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , if  $n > 1$

$$p(\mathcal{D}^n|\theta) = p(\mathbf{x}_n|\theta)p(\mathcal{D}^{n-1}|\theta)$$
$$\rightarrow p(\theta|\mathcal{D}^n) = \frac{p(\mathbf{x}_n|\theta)p(\theta|\mathcal{D}^{n-1})}{\int p(\mathbf{x}_n|\theta)p(\theta|\mathcal{D}^{n-1})d\theta}$$

- $p(\theta|\mathcal{D}^0) = p(\theta)$  and  $p(\theta|\mathbf{x}_1)$ ,  $p(\theta|\mathbf{x}_1, \mathbf{x}_2)$ ,  $\dots$
- To estimate  $p(\theta|\mathcal{D}^n)$ , all the training data in  $\mathcal{D}^{n-1}$  should be kept
- sufficient statistics, where distributions can be represented using only a few parameters

# Recursive Bayes learning

## Problems

- Provided  $\mathcal{D}^n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , if  $n > 1$

$$p(\mathcal{D}^n|\theta) = p(\mathbf{x}_n|\theta)p(\mathcal{D}^{n-1}|\theta)$$
$$\rightarrow p(\theta|\mathcal{D}^n) = \frac{p(\mathbf{x}_n|\theta)p(\theta|\mathcal{D}^{n-1})}{\int p(\mathbf{x}_n|\theta)p(\theta|\mathcal{D}^{n-1})d\theta}$$

- $p(\theta|\mathcal{D}^0) = p(\theta)$  and  $p(\theta|\mathbf{x}_1), p(\theta|\mathbf{x}_1, \mathbf{x}_2), \dots$
- To estimate  $p(\theta|\mathcal{D}^n)$ , all the training data in  $\mathcal{D}^{n-1}$  should be kept
- sufficient statistics, where distributions can be represented using only a few parameters

## 1 Introduction

## 2 Maximum-Likelihood Estimation

- General Principle
- Gaussian Case: Unknown Mean  $\mu$
- Gaussian Case: Unknown  $\mu$  and  $\Sigma$
- Bias

## 3 Bayesian Estimation

- Parameter Distribution
- Gaussian Case

## 4 General Theory

- MLE vs. Bayes estimates

# Comparison

## Different behaviors

	MLE	Bayes
computation	differential, gradient	multidimensional integration
Interpretability	single best model $\hat{\theta}$	weighted average of models
prior	$p(\mathbf{x} \hat{\theta}) = p(\mathbf{x} \theta)$	$p(\mathbf{x} \mathcal{D}) \neq p(\mathbf{x} \hat{\theta})$

## Similar behavior

- asymptotic limit of infinite training data
- strongly peaked  $p(\mathbf{x}|\hat{\theta})$  and the prior  $p(\theta)$  is uniform or flat

# Classification error

## Three sources

- **Bayes error**: overlapping densities for different category  $p(\mathbf{x}|\omega_i)$ , intrinsic and cannot be eliminated
- **Model error**: domain knowledge dependent
- **Estimation error**: parameters estimated from a finite set of samples, error can be reduced by increasing the number of training data