# Review of Estimation Theory

Berlin 2003

References:

1. X. Huang et. al., Spoken Language Processing, Chapter 3

# Introduction

- Estimation theory is the most important theory and method in $\mathtt{statistical\ inference}$
- Statistical inference
  - Data generated in accordance with some unknown probability distribution must be analyzed
  - Some type of inference about the unknown distribution must be made like the characteristics (parameters) of the distribution generating the experimental data, the mean and variance etc.

The vector of random variables
$$X = \{X_1, X_2, ..., X_n\} \underbrace{\qquad \theta(X) \qquad}_{\text{estimator}} g(X|\boldsymbol{\Phi})$$

The vector of sample values
$$x = \{x_1, x_2, ..., x_n\} \underbrace{\qquad \theta(x) \qquad}_{\text{estimate}} g(x|\boldsymbol{\Phi})$$

$\boldsymbol{\Phi}$ :the parameters of the distribution

# Introduction

- Three common estimators (estimation methods)
  - Minimum mean square estimator
    - Estimate the random variable itself
    - Function approximation, curve fitting, …
  - Maximum likelihood estimator
    - Estimate the parameters of the distribution of the random variables
  - Bayes' estimator
    - Estimate the parameters of the distribution of the random variables

# Minimum Mean Square Error Estimation and Least Square Error Estimation

- There are two random variables $X$ and $Y$. When observing the value of $X$, we want to find a transform $\hat{Y} = g(X, \boldsymbol{\Phi})$ ( $\boldsymbol{\Phi}$ the parameter vectors of function $g$ ) to predict the value of $Y$

  - **Minimum Mean Square Error Estimation**

    $$\boldsymbol{\Phi}_{MMSE} = \arg \min_{\boldsymbol{\Phi}} \left[ E \left[ (Y - g(X, \boldsymbol{\Phi}))^2 \right] \right]$$

    If the joint distribution $f_{X,Y}(X,Y)$ Is known

  - **Least Square Error Estimation**

    $$\boldsymbol{\Phi}_{LSE} = \arg \min_{\boldsymbol{\Phi}} \sum_{i=1}^{n} \left[ y_i - g(x_i, \boldsymbol{\Phi}) \right]^2$$

    When samples of $(x_i, y_i)$ pairs are observed

- Base on the law of large numbers, when the joint probability $\hat{Y} = f_{X,Y}(X,Y)$ is uniform or the number of samples approaches to infinity, MMSE and LSE are equivalent

# Minimum Mean Square Error Estimation and Least Square Error Estimation

- Constant functions $g(X) = c$

  - **MMSE**

    $$\nabla_c E\left[(Y - c)^2\right] = 0$$
    $$\therefore c_{MMSE} = E[Y]$$

    <span style="color:blue">mean</span>

  - **-- LSE**

    $$\nabla_c \sum_{i=1}^{n} (y_i - c)^2 = 0$$
    $$\therefore c_{LSE} = \frac{1}{n} \sum_{i=1}^{n} y_i$$

    <span style="color:blue">sample mean</span>

- Linear functions $g(X) = aX + b$

  - **MMSE**

    $$\nabla_a E\left[(Y - (aX + b))^2\right] = 0 \quad aE\left[X^2\right] + bE[X] - E[XY] = 0$$
    $$\nabla_b E\left[(Y - (aX + b))^2\right] = 0 \quad aE[X] + b - E[Y] = 0$$

    $$a = \frac{\text{cov}(X,Y)}{Var(X)}$$

    $$b = E[Y] - \rho_{XY} \frac{\sigma_Y}{\sigma_X} E[X]$$

# Minimum Mean Square Error Estimation and Least Square Error Estimation

- Linear functions
  - **LSE**
    - Suppose that $x$ are d-dimensional vectors and y are scalars

$$\hat{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = XA = \begin{bmatrix} \overset{c_0}{1} & \overset{c_1}{x_1^1} & \cdots & \overset{c_d}{x_1^d} \\ 1 & x_2^1 & \cdots & x_2^d \\ \vdots & \vdots & & \vdots \\ 1 & x_n^1 & \cdots & x_n^d \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix}$$
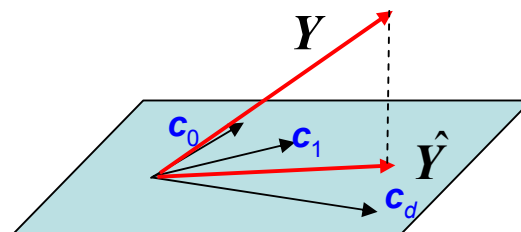
$$e(A) = \left\| \hat{Y} - Y \right\| = \sum_{i=1}^{n} \left( A^t x_i - y_1 \right)^2$$

$$\nabla e(A) = \sum_{i=1}^{n} 2\left( A^t x_i - y_i \right) x_i = 2 X^t (XA - Y) = 0$$

$$\Rightarrow X^t X A = X^t Y$$

$$\Rightarrow A = \left( X^t X \right)^{-1} X^t Y$$

.....



6

# Maximum Likelihood Estimation (MLE/ML)

- ML is the most widely used parametric estimation method

- A set of random samples $X = \{X_1, X_2, ..., X_n\}$ is to be drawn independently according to a distribution with the pdf $p(x|\boldsymbol{\Phi})$

  - Given a sequence of random samples $\boldsymbol{x} = (x_1, x_2, ..., x_n)$ the likelihood of it is defined as $p_n(\boldsymbol{x}|\boldsymbol{\Phi})$, a joint pdf of $(x_1, x_2, ..., x_n)$

    $$p_n(\boldsymbol{x}|\boldsymbol{\Phi}) = \prod_{k=1}^{n} p(x_k|\boldsymbol{\Phi}), \because X_1, X_2, ... X_n \text{ are iid}$$

  - Maximum likelihood estimator of $\boldsymbol{\Phi}$ is denoted as

    $$\boldsymbol{\Phi}_{ML} = \arg\max_{\boldsymbol{\Phi}} p_n(\boldsymbol{x}|\boldsymbol{\Phi}) = \arg\max_{\boldsymbol{\Phi}} \prod_{k=1}^{n} p(x_k|\boldsymbol{\Phi})$$

  - Since the logarithm function is *monotonically increasing function*, the parameter set $\boldsymbol{\Phi}_{ML}$ that maximizes the log-likelihood should also maximize the likelihood. The log-likelihood can be expressed as: $l(\boldsymbol{\Phi}) = log\ p_n(\boldsymbol{x}|\boldsymbol{\Phi}) = \sum_{k=1}^{n} \log p(x_k|\boldsymbol{\Phi})$

# Maximum Likelihood Estimation (MLE/ML)

- If $p_n(x|\boldsymbol{\Phi})$ is differentiable function of $\boldsymbol{\Phi}$, $\boldsymbol{\Phi}_{ML}$ can be attained by taking the partial derivative with respect to $\boldsymbol{\Phi}$ and setting it to zero

  – Let $\boldsymbol{\Phi}$ be a M-component parameter vector $\boldsymbol{\Phi} = (\Phi_1, \Phi_2, ..., \Phi_M)^t$

$$\nabla_{\boldsymbol{\Phi}} l(\boldsymbol{\Phi}) = \nabla_{\boldsymbol{\Phi}} \sum_{k=1}^{n} log \; p(x_k|\boldsymbol{\Phi}) = \begin{bmatrix} \dfrac{\partial l(\boldsymbol{\Phi})}{\partial \Phi_1} \\ . \\ . \\ . \\ \dfrac{\partial l(\boldsymbol{\Phi})}{\partial \Phi_M} \end{bmatrix} = 0$$

- Example: $p(x|\boldsymbol{\Phi})$ is a univariate Gaussian pdf with the parameter set $(\mu, \sigma^2)$

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

$$log \; p_n(x|\boldsymbol{\Phi}) = \sum_{k=1}^{n} log \; p(x_k|\boldsymbol{\Phi}) = \sum_{k=1}^{n} log\left( \frac{1}{\sqrt{2\pi}\,\sigma} exp\left[ -\frac{(x_k-\mu)^2}{2\sigma^2} \right] \right) = -\frac{n}{2} log\left(2\pi\sigma^2\right) - \frac{1}{2\sigma^2} \sum_{k=1}^{n} (x_k - \mu)^2$$

# Maximum Likelihood Estimation (MLE/ML)

- Example: univariate Gaussian pdf (cont.)
  - Take the partial derivatives of the above expression and set them to zero

$$\frac{\partial}{\partial \mu} \log p_n\left(x | \boldsymbol{\Phi}\right) = \frac{1}{\sigma^2} \sum_{k=1}^{n}\left(x_k - \mu\right) = 0$$

$$\frac{\partial}{\partial \sigma^2} \log p_n\left(x | \boldsymbol{\Phi}\right) = -\frac{n}{\sigma^2} + \sum_{k=1}^{n} \frac{\left(x_k - \mu\right)^2}{\sigma^4} = 0$$

  - The maximum likelihood estimates for $\mu$ and $\sigma^2$ are

$$\mu_{ML} = \sum_{k=1}^{n} x_k = E\left(x\right)$$

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

  - The maximum likelihood estimation for mean and variance is just **the sample mean and variance**

# Maximum Likelihood Estimation (MLE/ML)

- Example: multivariate Gaussian pdf (cont.)

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

  – The maximum likelihood estimates for $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are

$$\hat{\boldsymbol{\mu}}_{MLE}=\frac{1}{n}\sum_{k=1}^{n}x_{k}$$

$$\hat{\boldsymbol{\Sigma}}_{MLE}=\frac{1}{n}\sum_{k=1}^{n}\left(x_{k}-\hat{\boldsymbol{\mu}}_{MLE}\right)\left(x_{k}-\hat{\boldsymbol{\mu}}_{MLE}\right)^{t}$$

$$=E\left[\left(x_{k}-\hat{\boldsymbol{\mu}}_{MLE}\right)\left(x_{k}-\hat{\boldsymbol{\mu}}_{MLE}\right)^{t}\right]$$

  - The maximum likelihood estimation for mean vector and variance matrix is just **the sample mean vector and variance matrix**

- **In fact, $\boldsymbol{\Phi}_{MLE}$ itself is also a Gaussian distribution**

# Bayesian Estimation

- Bayesian estimation has a different philosophy than maximum likelihood (ML) estimation
  - ML assumes the parameter set $\boldsymbol{\Phi}$ is fixed but unknown (non-informative, uniform prior)
  - Bayesian estimation assumes the parameter set $\boldsymbol{\Phi}$ itself is a random variable with a prior distribution $p(\boldsymbol{\Phi})$

  - Given a sequence of random samples $\boldsymbol{x} = (x_1, x_2, ..., x_n)$, which are i.i.d. with a joint pdf $p(\boldsymbol{x}|\boldsymbol{\Phi})$, the posterior distribution of $\boldsymbol{\Phi}$ can be the following according to the Bayes' rule

$$p(\boldsymbol{\Phi}|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|\boldsymbol{\Phi})p(\boldsymbol{\Phi})}{p(\boldsymbol{x})} \propto p(\boldsymbol{x}|\boldsymbol{\Phi})p(\boldsymbol{\Phi})$$

# Bayesian Estimation

- $p(\boldsymbol{\Phi}|\boldsymbol{x})$ : the posterior probability, the distribution of $\boldsymbol{\Phi}$ after we observed the values of random variables
- $p(\boldsymbol{\Phi})$ : a conjugate prior of the random variables (or vector) is defined as the prior distribution for the parameters of the density function (e.g. $\boldsymbol{\Phi}$ ) of the random variables (or vectors)
  - Before we observed the values of random variables

- The joint pdf/likelihood function

$$p(\boldsymbol{x}|\boldsymbol{\Phi}) = \frac{1}{(2\pi)^{n/2}\sigma^n} \exp\left[-\frac{1}{2}\sum_{i=1}^{n}\left(\frac{x_i - \Phi}{\sigma}\right)^2\right] \propto \exp\left[-\frac{1}{2}\sum_{i=1}^{n}\left(\frac{x_i - \Phi}{\sigma}\right)^2\right]$$

- The prior is also a Gaussian distribution

$$p(\boldsymbol{\Phi}) = \frac{1}{(2\pi)^{1/2}\nu} \exp\left[-\frac{1}{2}\left(\frac{\Phi - \mu}{\nu}\right)^2\right] \propto \exp\left[-\frac{1}{2}\left(\frac{\Phi - \mu}{\nu}\right)^2\right]$$

# Maximum a Posterior Probability (MAP)

- The MAP chooses a estimate $\Phi_{MAP}$ that maximizes the posterior probability $p(\Phi|x)$ is the most common Bayesian estimator

$$\Phi_{MAP} = \arg \max_{\Phi} p(\Phi|x) = \arg \max_{\Phi} p(x|\Phi)p(\Phi)$$

$$\Phi_{MAP} = \arg \max_{\Phi} [log \ p(x|\Phi) + log \ p(\Phi)]$$

$$\frac{\partial \ log \ p(x|\Phi)}{\partial \Phi} + \frac{log \ p(\Phi)}{\partial \Phi} = 0$$

- For example, the conjugate prior for the mean of a Gaussian pdf is also a Gaussian pdf
  - Supposed in previous example, $X = \{X_1, X_2, ..., X_n\}$ is drawn from a Gaussian which mean $\Phi$ is unknown and variance $\sigma^2$ is known, while the conjugate prior (is a Gaussian) with mean $\mu$ and variance $v^2$
  - The MAP estimated $\Phi$ is:

$$\Phi_{MAP} = \frac{\sigma^2 \mu + n v^2 \bar{x}_n}{\sigma^2 + n v^2}, \ n \text{ is no. of training samples, } \bar{x}_n \text{ is the sample mean}$$

13

# Bayes' Decision Theory

- A decision-making based on both the posterior knowledge obtained from specific observation data and prior knowledge of the categories
  - Prior class probabilities $P(\omega_i), \;\; \forall \; \text{class } i$
  - Class-conditioned probabilities (likelihoods) $P(x|\omega_i), \;\; \forall \; \text{class } i$

$$k = \arg \max_i P(\omega_i|x) = \arg \max_i \frac{P(x|\omega_i)P(\omega_i)}{P(x)} = \arg \max_i \frac{P(x|\omega_i)P(\omega_i)}{\sum_{j=1} P(x|\omega_j)P(\omega_j)}$$

$$\therefore \; k = \arg \max_i P(x|\omega_i)P(\omega_i)$$

# Bayes' Decision Theory

- Bayes' decision rule designed to minimize the overall risk involved in making decision
  - The **expected loss** (conditional risk) when making decision $\delta_i$

$$R(\delta_i|x) = \sum_j l(\delta_i|\omega_j, x) P(\omega_j|x), \quad \text{where} \quad l(\delta_i|\omega_j, x) = \begin{cases} 0, & i = j \\ 1, & i \neq j \end{cases}$$

a decision    The class x might belong to    loss function

  - The overall risk  (Bayes' risk)

$$R = \int_{-\infty}^{\infty} R(\delta(x)|x) p(x) dx, \ \delta(x): \text{the selected decision for a sample } x$$

# Bayes' Decision Theory

- Minimize the overall risk (classification error) by computing the conditional risks and select the decision $\delta_i$ for which the conditional risk $R(\delta_i|x)$ is minimum, i.e., $P(\omega_i|x)$ is maximum

$$R(\delta_i|x) = \sum_j l(\delta_i|\omega_j, x) P(\omega_j|x) = \sum_{j \neq i} P(\omega_j|x)$$

$$= \sum_j P(\omega_j|x) - P(\omega_i|x)$$

$$= 1 - P(\omega_i|x)$$

the decision should be made

$$\delta(x) = \arg\max_i P(\omega_i|x) = \arg\max_i P(x|\omega_i) P(\omega_j)$$

  – Called the **minimum-error-rate decision rule** which minimizes the classification error rate

16

# Bayes' Decision Theory

- Two-class pattern classification

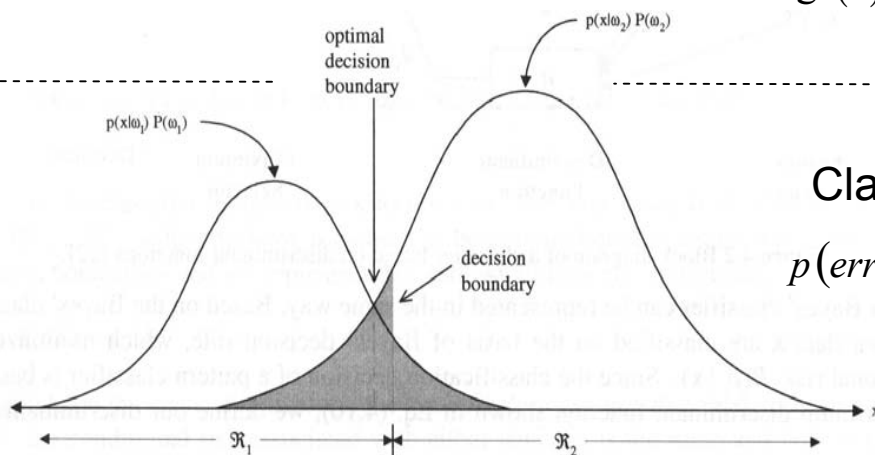$$d_1(x) = P(\omega_1|x) \cong P(x|\omega_1)P(\omega_1), \ \ d_2(x) = P(\omega_2|x) \cong P(x|\omega_2)P(\omega_2)$$

**Bayes' Classifier**

$$P(x|\omega_1)P(\omega_1) \overset{\omega_1}{\underset{\omega_2}{\gtrless}} P(x|\omega_2)P(\omega_2)$$

$P(\omega_2|x)$     $P(\omega_2|x)$

Likelihood ratio or log-likelihood ratio:

$$l(x) = \frac{P(x|\omega_1)}{P(x|\omega_2)} \overset{\omega_1}{\underset{\omega_2}{\gtrless}} \frac{P(\omega_2)}{P(\omega_1)}$$

$$\log l(x) = \log P(x|\omega_1) - \log P(x|\omega_2) \overset{\omega_1}{\underset{\omega_2}{\gtrless}} \log P(\omega_2) - \log P(\omega_1)$$



Figure 4.1 Calculation of the likelihood of classification error [22]. The shaded area represents the integral value in Eq. (4.9).

X falls in $R_2$, but the true class is $\omega_1$

Classification error:

$$p(error) = P(x \in R_1, \omega_2) + P(x \in R_2, \omega_1)$$
$$= P(x \in R_1|\omega_2)P(\omega_2) + P(x \in R_2|\omega_1)P(\omega_1)$$
$$= \int_{R_1} P(x|\omega_2)P(\omega_2)dx + \int_{R_2} P(x|\omega_1)P(\omega_1)dx$$