

Lecture 7

* LDA, QDA

$\rightarrow \sum_i = \sigma^2 I$, $\sum_i = \Sigma$, Σ_i 's are not the same

* Bernoulli pmf

* Error bound

* MLE



Chernoff bound

$$\min(a, b) \leq a^\beta b^{1-\beta} \quad \text{for } a, b \geq 0 \text{ and } 0 \leq \beta \leq 1$$

Often it is difficult to find $P(\text{error})$ as it needs decision boundaries.
In high dimensions it is very hard to find decision boundaries.

Thus instead of exact $P(\text{error})$, we seek a bound on this error.

$$\begin{aligned} P(\text{error}) &= \int \min\{P(w_1|x), P(w_2|x)\} \delta(x) dx \\ &\leq \text{Upper bound} \\ \text{W.L.O.G. } a > b \mid b &\leq \left(\frac{a}{b}\right)^\beta \cdot b \end{aligned}$$

$$\begin{aligned}
 P(\text{error}) &\leq \int_{-\infty}^{\infty} P(w_1|x)^{\beta} P(w_2|x)^{1-\beta} \rho(x) dx \\
 &= P(w_1)^{\beta} P(w_2)^{1-\beta} \cdot \int_{-\infty}^{\infty} P(x|w_1)^{\beta} P(x|w_2)^{1-\beta} dx \\
 &= P(w_1)^{\beta} P(w_2)^{1-\beta} \underbrace{\int_{-\infty}^{\infty} P(x|w_1)^{\beta} P(x|w_2)^{1-\beta} dx}_{e^{-K(\beta)}}
 \end{aligned}$$

$$P(x|w_1) \sim N(\mu_1, \Sigma_1)$$

$$P(x|w_2) \sim N(\mu_2, \Sigma_2)$$

$$P(\text{error}) = P(w_1)^{\beta} P(w_2)^{1-\beta} e^{-K(\beta)}$$

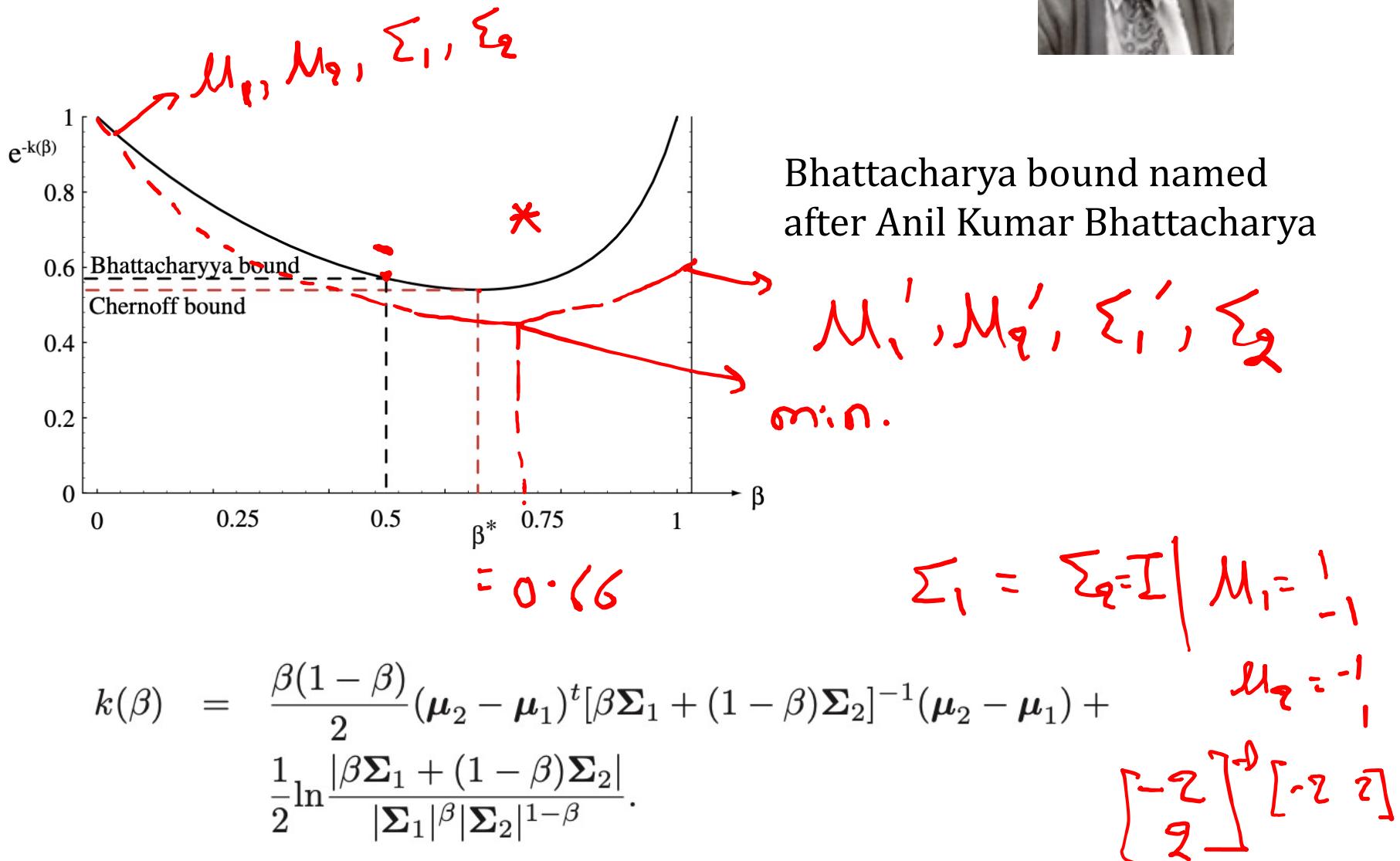
$$\int p(x) dx = f(x)$$

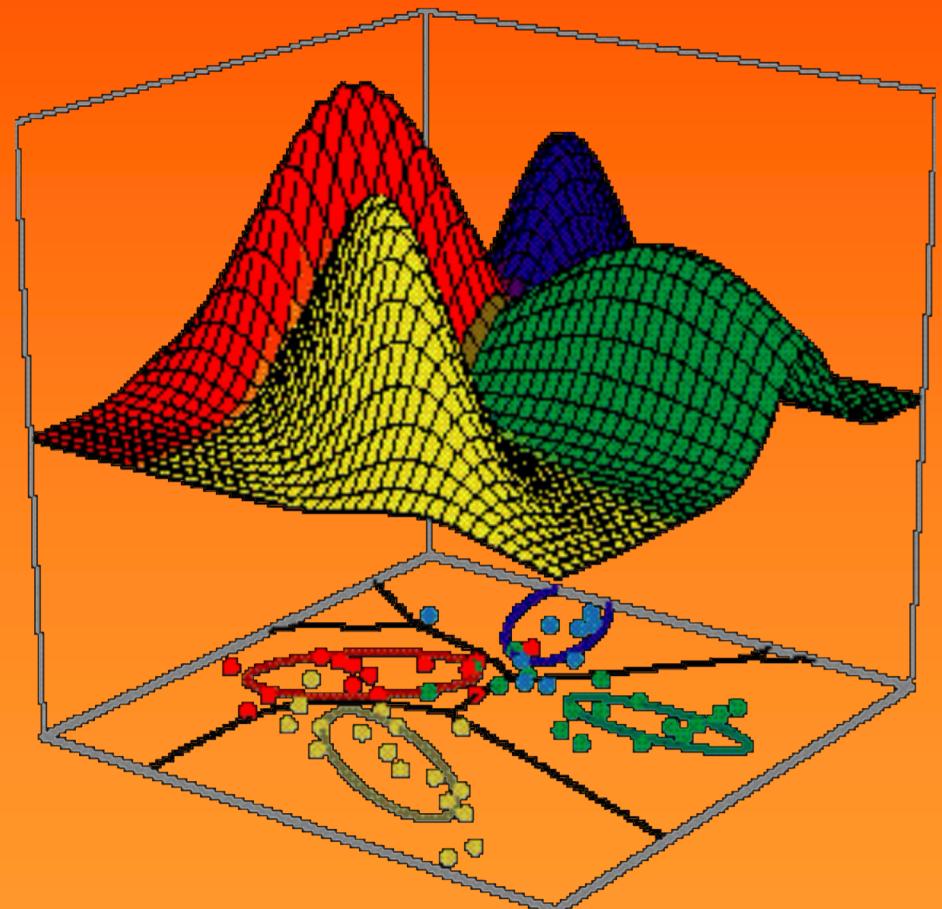
Bhattacharya bound, $\beta = 1/2$

$$P(\text{error}) = \underbrace{\sqrt{P(w_1) P(w_2)}}_{\sim} e^{-K(1/2)}$$

$$P(\text{error}) \lesssim \frac{1}{C} \sqrt{P(w_1) P(w_2)}$$

Chernoff bound





Pattern Classification

All materials in these slides were taken from
Pattern Classification (2nd ed) by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000
with the permission of the authors and the publisher

Chapter 3: Maximum-Likelihood & Bayesian Parameter Estimation (part 1)

- Introduction
- Maximum-Likelihood Estimation
 - Example of a Specific Case
 - The Gaussian Case: unknown μ and σ
 - Bias

Ronald Aylmer Fisher (1890~1962)



■ Known for :

- 1912 : Maximum likelihood
- 1922 : F-test
- 1925 : Analysis of variance
(Statistical Method for
Research Workers)

■ Notable Prizes :

- Royal Medal (1938)
- Copley Medal (1955)

• **Introduction**

- 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!)

- A priori information about the problem
- Normality of $P(x | \omega_i)$

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

- Characterized by mean and cov parameters
- Estimation techniques
 - Maximum-Likelihood (ML) and the Bayesian estimations
 - Results are nearly identical, but the approaches are different

- Parameters in ML estimation are fixed but unknown!
- Best parameters are obtained by maximizing the probability of obtaining the samples observed
- Bayesian methods view the parameters as random variables having some known distribution
- In either approach, we use $P(\omega_i \mid x)$ for our classification rule!

Maximum-Likelihood Estimation

$\Sigma \rightarrow \text{Symm.}$

$d > N$
no. of samples

- Has good convergence properties as the sample size increases
- Simpler than any other alternative techniques
- General principle
- Assume we have c classes and

$$\in \mathbb{R}^d$$

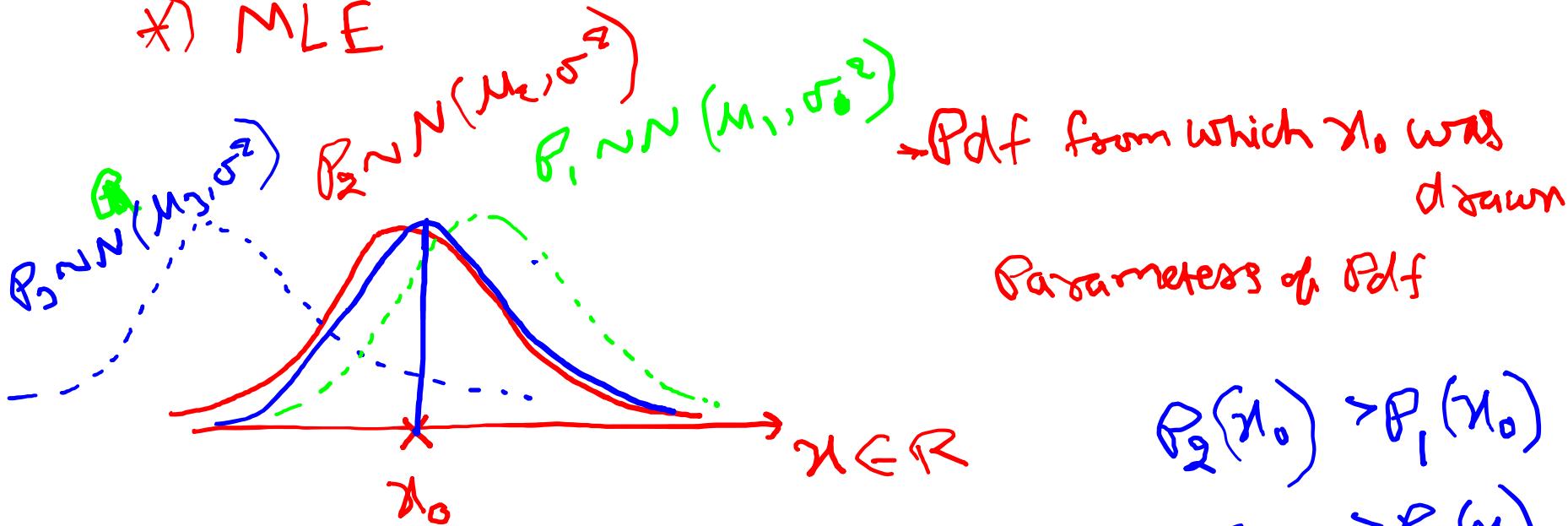
$c \cdot d$ param.
 $c \frac{d(d+1)}{2}^+$

$$\theta = (\mu_j, \Sigma_j) = (\mu_j^1, \mu_j^2, \dots, \sigma_j^{11}, \sigma_j^{22}, \text{cov}(x_j^m, x_j^n) \dots)$$

$\underbrace{\quad}_{c}$

Total no of parameters to be computed?

* MLE



Parameters of Pdf

$$P_2(x_0) > P_1(x_0)$$

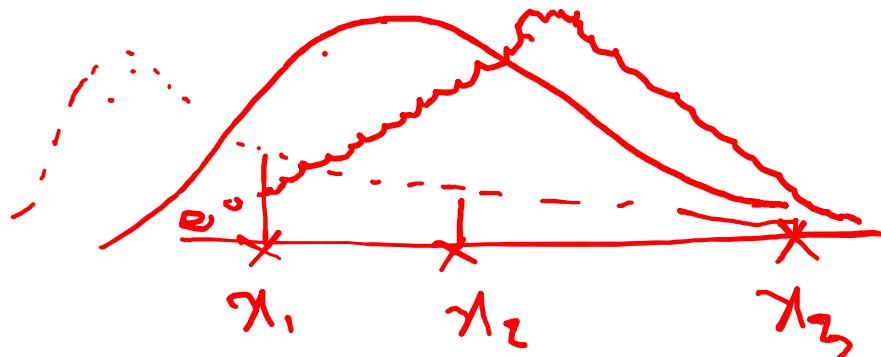
$$P_2(x_0) > P_3(x_0)$$

When you do not anything?

→ 1st step.

→ Assumption: Gaussian.

→ $\mu = x_0, \sigma^2 \rightarrow X$



σ is known.

$$P(x_1 | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{x_1 - \mu}{\sigma} \right)^2}$$

$$\mu = ?$$

$$P(x_2 | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{x_2 - \mu}{\sigma} \right)^2}$$

- ① Joint prob. of obs. x_1, x_2, x_3 from parameters μ, σ^2

MLE Demo

x_1, x_2, x_3 are iids. \rightarrow assumption.

$$P(x_1, x_2, x_3 | \mu, \sigma^2) = \prod_{i=1}^3 f(x_i | \mu, \sigma^2)$$

$$P(x_1, x_2, \dots, x_n | \mu, \sigma^2) = \prod_{i=1}^n f(x_i | \mu, \sigma^2)$$

likelihood w.r.t. Sampled.

Which Gaussian is most likely to have generated
 x_1, x_2, \dots, x_n

$$P(x|\theta) = \prod_{i=1}^n P(x_i|\mu, \sigma^2)$$

$$X = \{x_1, x_2, \dots, x_n\}$$

$\theta \rightarrow$ Parameters, $\theta = \mu,$

$$\ln P(x|\theta) = \sum_{i=1}^n \ln P(x_i|\mu, \sigma^2)$$

$$= \sum_{i=1}^n \ln \frac{1}{\sqrt{2\pi\sigma^2}} e^{\left\{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}\right\}}$$

$$= \sum_{i=1}^n \left[\ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} \right]$$

We want to maximize log-likelihood to obtain

$$\hat{\theta}, \text{ or here } \hat{\theta} = \bar{M}$$

$$\frac{\partial}{\partial \theta} \ln \theta(x|\theta) = -\frac{1}{2\sigma^2} \sum_i -2(x_i - \bar{M}) = 0$$

$$\bar{M} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$f(\theta) = \ln \theta(x|\theta)$$

$$\theta = \{\theta_1, \theta_2, \theta_3, \dots, \theta_m\}$$

~~$$\nabla_{\theta} f(\theta) = 0$$~~

We want

$$\Theta = \cancel{\mu}, \sigma^2$$

$$\frac{\partial}{\partial \mu} F(\theta) = 0 \Rightarrow \mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\frac{\partial}{\partial \sigma^2} F(\theta) = 0 \Rightarrow \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$\boxed{\begin{aligned}\mu_{ML} &= \frac{1}{n} \sum_{i=1}^n x_i \\ \sigma_{ML}^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{ML})^2\end{aligned}}$$

ML
Estimation

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

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

P(D | θ) is called the likelihood of θ w.r.t. the set of samples)

- ML estimate of θ is, by definition the value that maximizes $P(D | \theta)$
“It is the value of θ that best agrees with the actually observed training sample”



$$\mu, \sigma^2 \text{ or } \{\theta_1, \theta_2\}$$

- Optimal estimation

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

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

- We define $I(\theta)$ as the log-likelihood function

$$I(\theta) = \ln P(D | \theta)$$

- New problem statement:
determine θ that maximizes the log-likelihood

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

Set of necessary conditions for an optimum is:

$$\begin{aligned}(\nabla_{\theta} \mathcal{L} = \sum_{k=1}^{K=n} \nabla_{\theta} \ln P(x_k | \theta)) \\ \nabla_{\theta} \mathcal{L} = 0\end{aligned}$$

- ML Estimation:

- Gaussian Case: *unknown μ and σ*
 $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$

$$\boxed{\begin{aligned} I &= \ln P(x_k | \theta) = -\frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2} (x_k - \theta_1)^2 \\ \nabla_{\theta} I &= \begin{pmatrix} \frac{\sigma}{\sigma\theta_1} (\ln P(x_k | \theta)) \\ \frac{\sigma}{\sigma\theta_2} (\ln P(x_k | \theta)) \end{pmatrix} = \mathbf{0} \\ \begin{cases} \frac{1}{\theta_2} (x_k - \theta_1) = \mathbf{0} \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} = \mathbf{0} \end{cases} \end{aligned}}$$

Summation:

$$\left\{ \begin{array}{l} \sum_{k=1}^{n-1} \frac{1}{\hat{\theta}_2} (x_k - \theta_1) = 0 \\ \end{array} \right. \quad (1)$$

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

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

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

ML estimate for multivariate case

$$P(x|\theta) = \theta e^{-\theta x}, x > 0 \quad \text{Pdf.}$$

$\hat{\theta}_{ML}$

Likelihood $F(\theta) = \ln P(X|\theta)$

$$X = \{x_1, x_2, \dots, x_n\}$$

$$\hat{\theta}_{ML} = \frac{1}{\sum x_i}$$

~~$F(x) = \ln \theta$~~

$$\checkmark \hat{\theta}_{ML} = \frac{n}{\sum x_i}$$

$$F(\theta) = \ln \prod_{i=1}^n \theta e^{-\theta x_i}$$

$$= \sum_{i=1}^n \ln \theta - \theta x_i$$

$$\frac{\partial F(\theta)}{\partial \theta} = 0$$

When $n \rightarrow \infty$ $\hat{\mu}_{ML} = \mu$

Where ' μ ' is the true value.

* ~~Biased~~ Estimate 

x_i are i.i.d.
 $\sim N(\mu, \sigma^2)$

$$\hat{\mu}_{ML} = \frac{1}{n} \sum_{i=1}^n x_i$$

and est' is called unbiased,

$$\begin{aligned}\mathbb{E}[\hat{\mu}_{ML}] &= E\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \mu \\ &= \frac{1}{n} \sum_{i=1}^n E(x_i) = \frac{1}{n} \cdot n \mu\end{aligned}$$