

## Lecture 2

### Parameter Estimation

How do we find the prob. dist. to a r.v.  $X$ ?

Three steps:

- i) choose a parametric model (e.g. Gaussian)  
call the parameters  $\theta$ .
- 2) assemble a collection of samples (observations)  
from  $X$ .

$$D = \{x_1, \dots, x_N\}$$

sample of  $X$

We assume  $x_i$ 's are independent.  $x_i$  are iid samples.  
independantly & identically distributed

- 3) Maximum likelihood (ML) principle.

"The optimal parameter  $\theta^*$  is that which maximizes  
the probability of the training data  $D$ ."  
(likelihood)

### ML estimate (MLE)

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \underset{\text{data}}{\underbrace{P(D|\theta)}} \quad \begin{array}{l} \text{likelihood of the } D \text{ wrt. } \theta. \\ \text{"likelihood function"} \end{array}$$

$$\begin{aligned} &= \underset{\theta}{\operatorname{argmax}} \log P(D|\theta) \quad \text{"log-likelihood" (LL)} \\ &= \underset{\theta}{\operatorname{argmin}} -\log P(D|\theta) \quad \begin{array}{l} \text{"negative log-likelihood"} \\ \text{(NLL)} \\ \text{"loss"} \end{array} \end{aligned}$$

Note:  $\underset{\theta}{\operatorname{argmax}} \log P(D|\theta)$   
 $D$  is known, so  $P(D|\theta)$  is a function of  $\theta$ .  
 This is not a prob dist in  $\theta$ !!!

Note:  $\log = \text{natural logarithm (log base e)}$

### Data log-likelihood term

$$\begin{aligned} l(\theta) &= \log P(D|\theta) \\ &= \log \prod_{i=1}^N p(x_i|\theta) \quad \downarrow \text{independence} \\ l(\theta) &= \sum_{i=1}^N \log p(x_i|\theta) \end{aligned}$$

$$\log(a \cdot b) = \log a + \log b$$

Assignment Project Exam Help  
<https://powcoder.com>  
 Add WeChat powcoder

To get the ML Solution:

if  $\theta$  is scalar, at local optimum of  $l(\theta)$ :

$$1) \frac{\partial}{\partial \theta} \log p(D|\theta) = 0 \text{ at } \theta^*$$

$$2) \frac{\partial^2}{\partial \theta^2} \log p(D|\theta) < 0 \text{ at } \theta^* \\ (\text{at a local maximum}) \\ (\text{concave})$$

3) check the boundary conditions  
on  $\theta$  (if necessary)

if  $\theta$  is a vector:  $\theta = [\theta_1, \dots, \theta_p]$

$$1) \nabla_{\theta} l(\theta) = \begin{bmatrix} \frac{\partial}{\partial \theta_1} l(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_p} l(\theta) \end{bmatrix} = 0$$

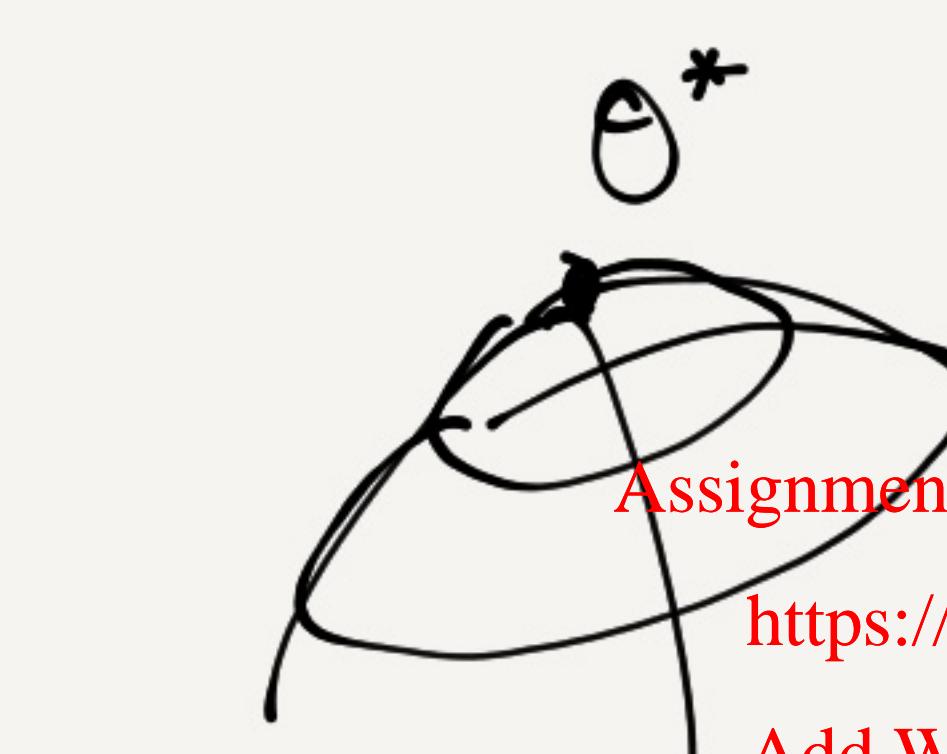
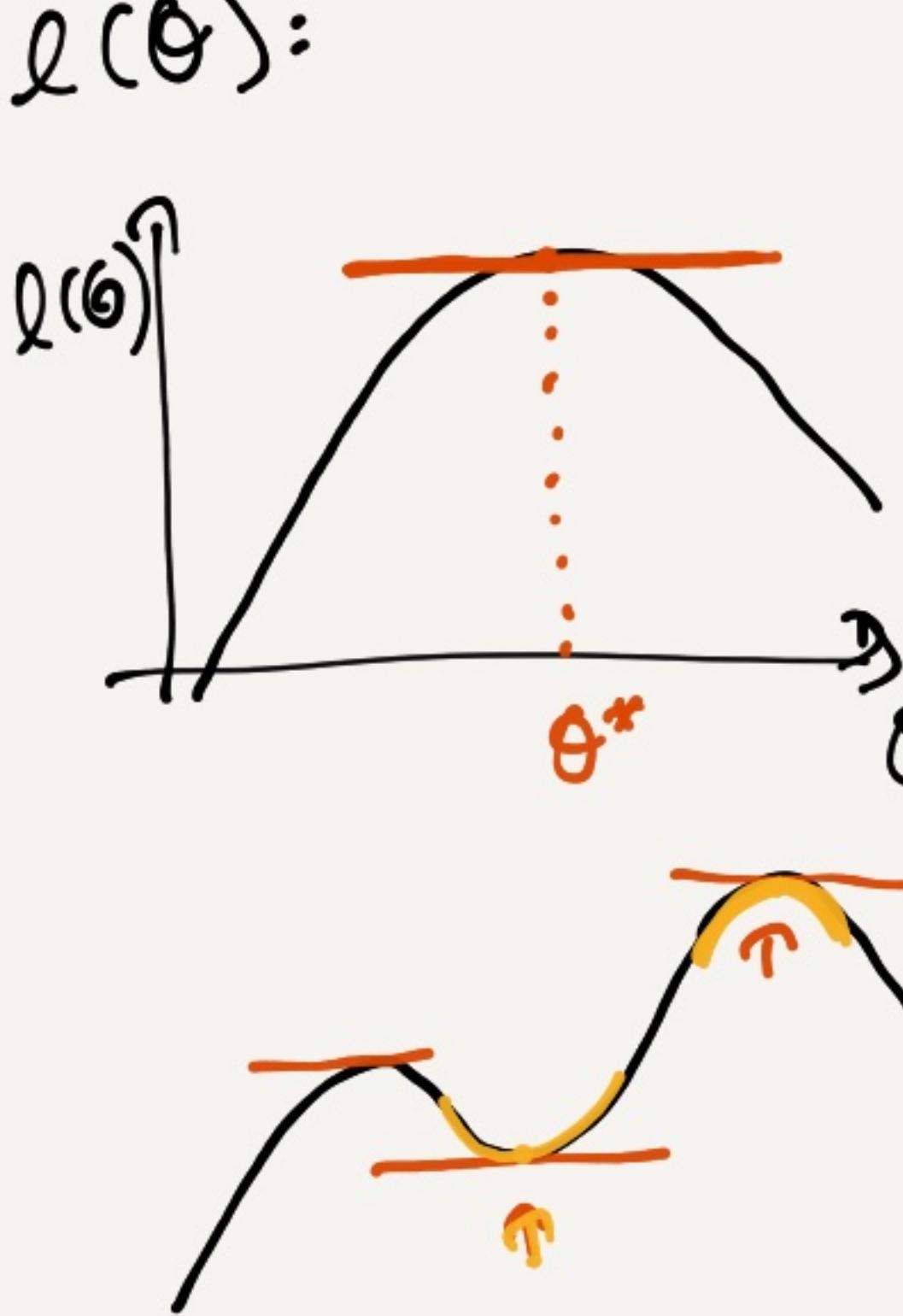
↑ gradient

$$2) \nabla_{\theta}^2 l(\theta) \preceq 0$$

↑ Hessian

negative definite.

$$\nabla_{\theta}^2 l(\theta) = \begin{bmatrix} \frac{\partial^2}{\partial \theta_1^2} l(\theta) & \cdots & \frac{\partial^2}{\partial \theta_1 \partial \theta_p} l(\theta) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial \theta_p \partial \theta_1} l(\theta) & \cdots & \frac{\partial^2}{\partial \theta_p^2} l(\theta) \end{bmatrix}$$



$H \preceq 0 \Rightarrow$   
 $\theta^T H \theta < 0 \quad \forall \theta$   
 (all directions are curving downwards)

"describes the curvature of  $l(\theta)$ "

Example: Bernoulli:

$$\theta = \pi, 0 \leq \pi \leq 1, x \in \{0, 1\}$$

$$\log a^b = b \log a$$

$$l(\theta) = \sum_{i=1}^N \log p(x_i | \theta)$$

$$= \sum_i \log \frac{\pi^{x_i} \cdot (1-\pi)^{1-x_i}}{1}$$

$$= \sum_i x_i \log \pi + (1-x_i) \log (1-\pi)$$

$$= \underbrace{\left( \sum_i x_i \right) \log \pi}_{\# \text{ of 1s}} + \underbrace{\left( \sum_i (1-x_i) \right) \log (1-\pi)}_{\# \text{ of 0s}}$$

$m = \sum_i x_i \leftarrow \text{"sufficient statistic"} - l(\theta) \text{ only depends on } N \text{ observations through this term.}$

$$= m \log \pi + (N-m) \log (1-\pi)$$

Solve for  $\pi$ :

$$1) \frac{\partial}{\partial \pi} l(\theta) = \left[ \frac{m}{\pi} + \frac{N-m}{1-\pi} (-1) \right] = 0 \quad \boxed{\times (\pi)(1-\pi)}$$

$$\frac{m}{\pi} (1-\pi) + (N-m)(-1)(\pi) = 0$$

$$m - N\pi = 0 \Rightarrow \boxed{\pi = \frac{m}{N} = \frac{1}{N} \sum_{i=1}^N x_i}$$

"fraction of 1's observed"  
(sample mean)

$$2) \frac{\partial^2}{\partial \pi^2} l(\theta) = \frac{\partial}{\partial \pi} \left( \frac{\partial}{\partial \pi} l(\theta) \right) = \frac{\partial}{\partial \pi} \left( \frac{m}{\pi} - \frac{(N-m)}{1-\pi} \right)$$

$$= \frac{-m}{\pi^2} - \frac{(N-m)(-1)(-1)}{(1-\pi)^2} < 0 \quad \checkmark$$

$$3) 0 \leq m \leq N \Rightarrow 0 \leq \pi \leq 1 \quad \checkmark$$

### Example: Gaussian

$\theta = (\mu, \sigma^2)$ ,  $\sigma^2$  is known.

$$l(\theta) = \sum_{i=1}^N \log p(x_i | \theta)$$

$$= \sum_i \left[ -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (x_i - \mu)^2 \right]$$

$$= -\frac{N}{2} \log 2\pi - \frac{N}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

what are the suff. stats?  
 $(x_i - \mu)^2 = x_i^2 - 2x_i\mu + \mu^2$

$$\sum_i x_i^2, \sum_i x_i$$

Solve for  $\mu$ :

$$\frac{\partial}{\partial \mu} l(\theta) = \cancel{\frac{1}{2\sigma^2}} \sum_{i=1}^N (x_i - \mu) \cdot \cancel{2} \cdot (-1) = 0$$

$$\begin{aligned} \sum_i x_i - \sum_i \mu &= 0 \\ \sum_i x_i - N\mu &= 0 \Rightarrow \boxed{\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i} \end{aligned}$$

Sample mean of  $x$ .

•  $\sigma^2$  is unknown,  $\mu$  is known.

$$1) \frac{\partial}{\partial \sigma^2} l(\theta) = -\frac{N}{2} \frac{1}{\sigma^2} - \frac{1}{2} \frac{1}{\sigma^4} (-1) \sum_i (x_i - \mu)^2 = 0 \quad \downarrow \times 6^4$$

$$= -\frac{N}{2} \sigma^2 + \frac{1}{2} \sum_i (x_i - \mu)^2 = 0$$

$$\Rightarrow \boxed{\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$

Sample variance  
of  $x$ .

### Multivariate Gaussian

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

## Estimators

- the estimate (e.g.  $\hat{\mu}$ ) is a number.
- the estimator is a r.v. over many datasets.
- the estimator  $\beta$  the value of the estimator for a given dataset D.

estimator:  $f(x_1, \dots, x_N) = \frac{1}{N} \sum_{i=1}^N x_i$       r.v. for each sample  
 $x_i \sim p(x_i | \theta)$   
true distribution

$$\text{ML estimate: } \hat{\mu} = f(x_1, \dots, x_N) \Big|_{x_1=x_1, \dots, x_N=x_N} \\ = \frac{1}{N} \sum_i x_i$$

- Since the estimator is a r.v., we can compute its mean & variance. Hence, we can quantify how good is the estimator.

Assignment Project Exam Help  
<https://powcoder.com>  
 Add WeChat powcoder

## Bias & Variance

- How do we measure "goodness"?
- $\hat{\theta} = f(x_1, \dots, x_N)$
- 1) Will it ever converge to the true value of  $\theta$ ?
  - $\text{Bias}(\hat{\theta}) = E_{x_1, \dots, x_N}[\hat{\theta} - \theta] = E[\hat{\theta}] - \theta$   
↑ true value
  - If the bias is non-zero, we can never get the true value of  $\theta$ , even with an infinite # of samples. "measures the expressiveness"

- 2) How long will it take to converge?  
 How many samples do we need?

$$\text{Var}(\hat{\theta}) = E_{x_1, \dots, x_N}[(\hat{\theta} - E[\hat{\theta}])^2]$$

"measures the uncertainty/variability"

Example: Gaussian (true value is  $\mu$ )

Estimator:  $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i$

Bias: mean  $E_{X_1 \dots X_N} [\hat{\mu}] = E_{X_1 \dots X_N} \left[ \frac{1}{N} \sum_i X_i \right]$

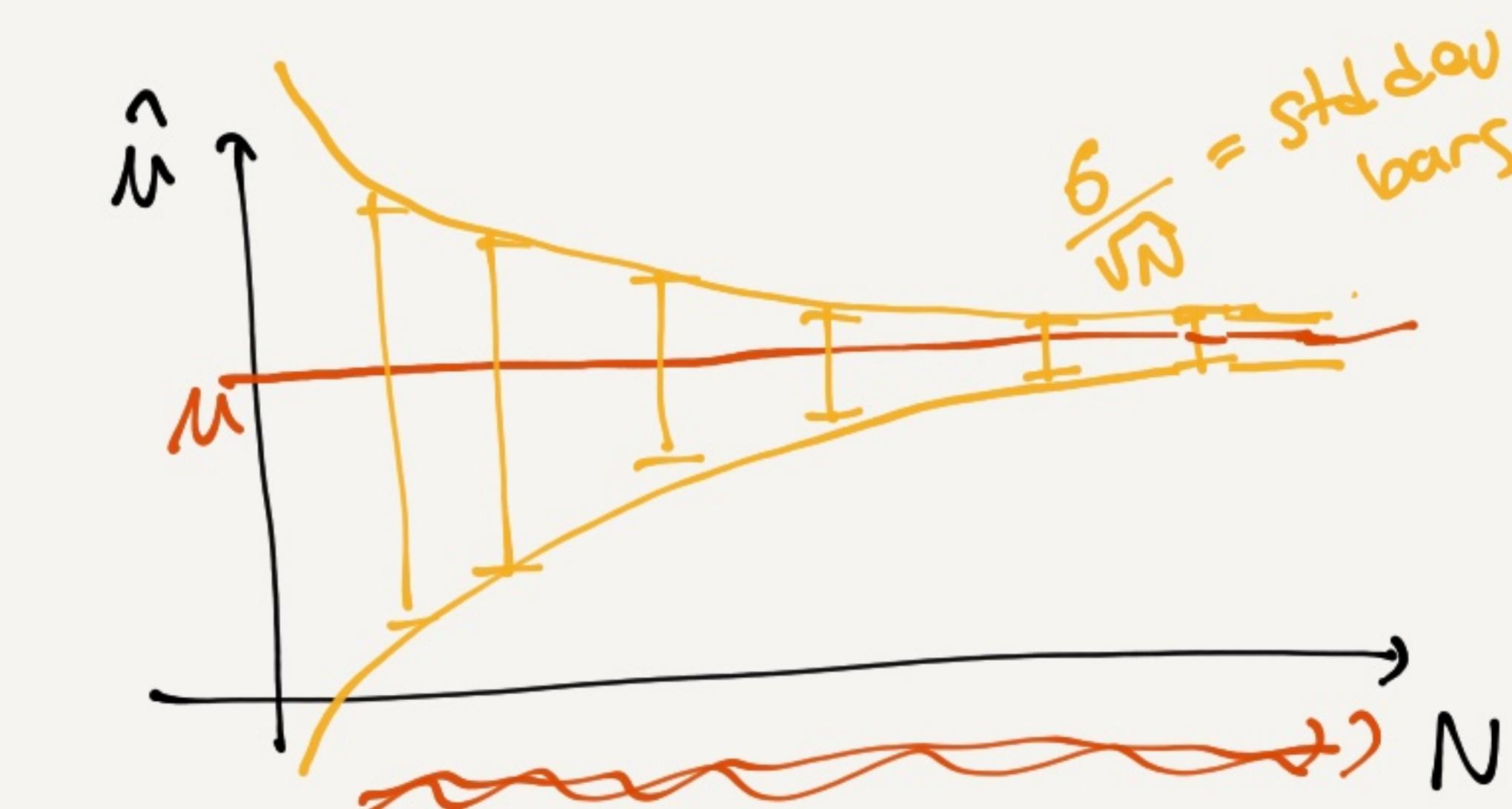
$$= \frac{1}{N} \sum_i E_{X_i} [X_i] = \frac{1}{N} \sum_i \mu = \mu$$

$\Rightarrow \boxed{\text{Bias}(\hat{\mu}) = \mu - \mu = 0}$  = "unbiased"

Variance:  $\text{Var}(\hat{\mu}) = E_{X_1 \dots X_N} \left[ (\hat{\mu} - E\hat{\mu})^2 \right]$

$$= E \left[ \left( \frac{1}{N} \sum_i X_i - \mu \right)^2 \right]$$
$$= \frac{1}{N^2} E \left[ \left( \sum_i (X_i - \mu) \right)^2 \right]$$
$$= \frac{1}{N^2} E \left[ \sum_i \sum_j (X_i - \mu)(X_j - \mu) \right]$$
$$= \underbrace{\frac{1}{N^2} \sum_i \sum_{i \neq j} (X_i - \mu)(X_j - \mu)}_{\text{if } i=j \Rightarrow E((X_i - \mu)^2) = \sigma^2} + \underbrace{\frac{1}{N^2} \sum_i \sum_{i \neq j} (X_i - \mu)(X_j - \mu)}_{\text{if } i \neq j \Rightarrow E((X_i - \mu)(X_j - \mu)) = \sigma_{ij}^2}$$
$$= \frac{1}{N^2} \left( \underbrace{\sum_{i=j} \sigma^2}_{N \sigma^2} + \underbrace{\sum_{i \neq j} \sigma_{ij}^2}_{=0} \right)$$

$\boxed{\text{Var}(\hat{\mu}) = \frac{\sigma^2}{N}}$  variance of the true Gaussian.



Variance converges to 0 as  $N \rightarrow \infty$   
i.e.  $p(\hat{\mu})$  concentrates around the true mean  $\mu$  as  $N$  increases.

Another example: variance of Gaussian. (PS 2-12)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$$E[\hat{\sigma}^2] = \frac{N-1}{N} \sigma^2$$

$$\boxed{\text{Bias}(\hat{\sigma}^2) = -\frac{1}{N} \sigma^2}$$

## Important Asymptotic Properties of MLE

1) consistent - as  $N \rightarrow \infty$ , the estimated value converges to the true value. i.e. asymptotically unbiased.

2) efficient - achieves the Cramér-Rao Lower Bound (CRLB) as  $N \rightarrow \infty$ .

CRLB is a theoretical lower bound on the variance of an unbiased estimator, for some  $p(x|\theta)$ .

(no unbiased estimator can have lower variance than the CRLB, & MLE achieves it)

Assignment Project Exam Help  
<https://powcoder.com>  
 Add WeChat powcoder

## MLE & Regression

• Supervised learning:  $D = \{(x_i, y_i)\}$

• input  $x \in \mathbb{R}$

• output  $y \in \mathbb{R}$

E.g. estimate a polynomial function

$$f(x, \theta) = \sum_{d=0}^k \theta_d x^d = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_d \end{bmatrix}^T \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_d \end{bmatrix} = \phi(x)^T \theta$$

$k=2: \theta_0 + \theta_1 x + \theta_2 x^2$

(linear function of parameters  $\theta$ )

Observe a noisy output  $y$  given an  $x$ :

$$y = f(x, \theta) + \epsilon, \quad \epsilon \sim N(0, \sigma^2) \text{ (iid)}$$

Equivalently:

$$p(y | x, \theta) = N(y | f(x, \theta), \sigma^2)$$

Given a dataset  $D = \{(x_i, y_i)\}$ , we can use MLE to find  $\theta$ :

$$\text{MLE: } \theta^* = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log p(y_i | x_i, \theta)$$

: (intuition)

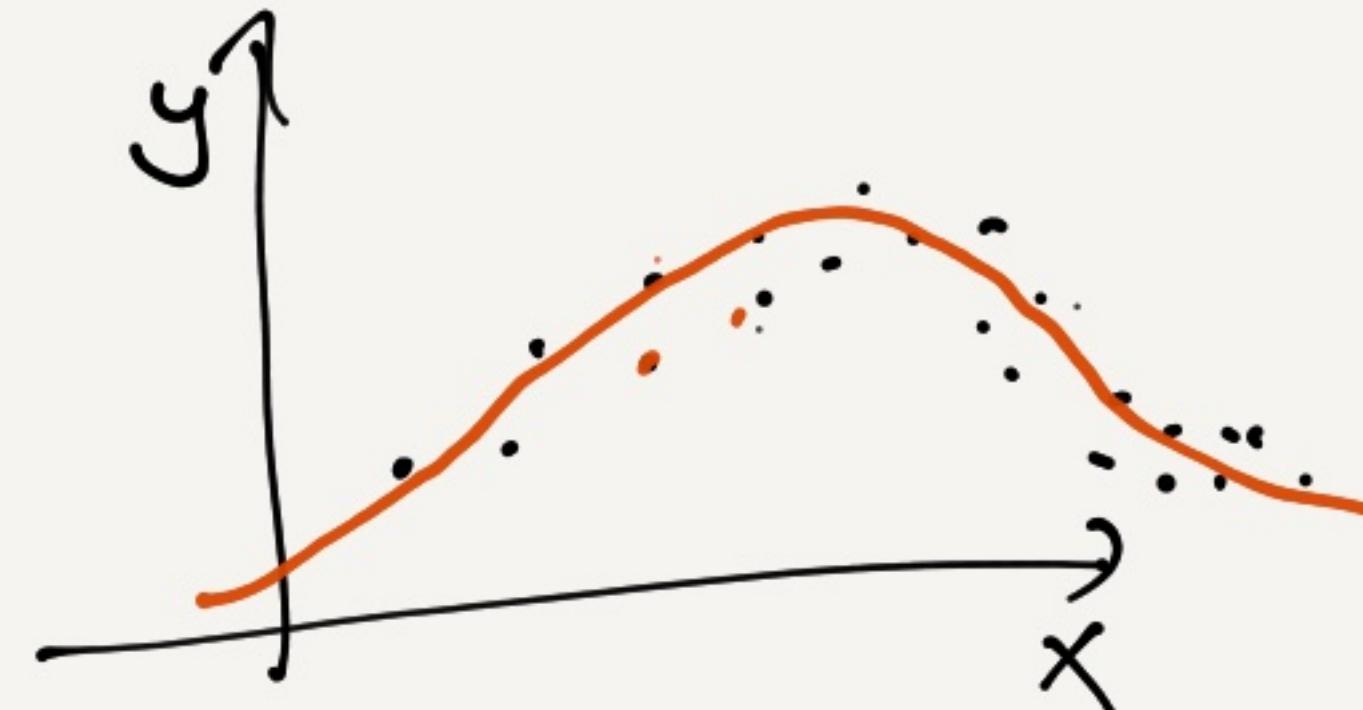
$$= \underset{\theta}{\operatorname{argmin}} \sum_i (y_i - f(x_i, \theta))^2$$

Squared error

$$\boxed{\theta^* = (\Phi \Phi^T)^{-1} \Phi y}$$

$$\Phi = [\phi(x_1) \dots \phi(x_N)]$$

$$y = [y_1 \dots y_N]$$



Notes:

- 1) ML is more general than LS
- 2) all the assumptions are explicit.
  - i) Gaussian Noise
  - ii) iid samples (iid noise)
  - iii)  $\mu = 0$ ,  $\sigma^2$  variance
- 3) We can change the assumptions to get other regression formulations (e.g. types of LS)
  - i) weighted LS (PS 2-8)
  - ii) regularized LS (lecture 3)
  - iii)  $L_p$  norms (PS 2-a)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

## Problem 2-6

MLE for m.v. Gaussian

$$D = \{x_1, \dots, x_N\}, \quad x_i \in \mathbb{R}^D$$

$$p(x) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} \|x - \mu\|_\Sigma^2\right)$$

Data LL

$$\log p(D) = \sum_{i=1}^N \log p(x_i) = \sum_{i=1}^N -\frac{1}{2} \|x_i - \mu\|_\Sigma^2 - \frac{1}{2} \log |\Sigma| - \frac{1}{2} \log 2\pi$$

a) Mean  $\mu$

$$l(\mu) = \sum_{i=1}^N -\frac{1}{2} \|x_i - \mu\|_\Sigma^2$$

$$= -\frac{1}{2} \sum_i (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$= -\frac{1}{2} \sum_i \underbrace{x_i^T \Sigma^{-1} x_i}_{x} - \underbrace{x_i^T \Sigma^{-1} \mu - \mu^T \Sigma^{-1} x_i}_{2x_i^T \Sigma^{-1} \mu} + \mu^T \Sigma^{-1} \mu$$

$$= -\frac{1}{2} \sum_i -2x_i^T \Sigma^{-1} \mu + \mu^T \Sigma^{-1} \mu$$

$$\frac{\partial}{\partial \mu} l(\mu) = \frac{1}{2} \sum_i -2\Sigma^{-1} x_i + 2\Sigma^{-1} \mu = 0$$

$$\Sigma * \left( \sum_i \Sigma^{-1} x_i + \Sigma^{-1} \mu \right) = 0$$

$$\sum_i [-x_i + \mu] = 0$$

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i$$

b) covariance

$$l(\Sigma) = \sum_{i=1}^N -\frac{1}{2} \|x_i - \mu\|_\Sigma^2 - \frac{1}{2} \log |\Sigma|$$

$$= \sum_{i=1}^N -\frac{1}{2} \underbrace{(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}_{a^T b} - \frac{1}{2} \log |\Sigma|$$

$$a^T b = \text{tr}(a^T b) = \text{tr}(ba^T)$$

$$= \sum_i -\frac{1}{2} \text{tr}(\Sigma^{-1} (x_i - \mu)(x_i - \mu)^T) - \frac{1}{2} \log |\Sigma|$$

$$\frac{\partial}{\partial \Sigma} \text{tr}(X^T A) = -(X^T A^T X^{-T}) \quad \frac{\partial}{\partial x} \frac{a}{x} = \frac{-a}{x^2}$$

$$\frac{\partial}{\partial \Sigma} \log |\Sigma| = X^{-T} \quad \Leftrightarrow \frac{\partial}{\partial x} \log x = \frac{1}{x}$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$$\frac{\partial l(\Sigma)}{\partial \Sigma} = \sum_i -\frac{1}{2} (-\Sigma^{-1} (x_i - \mu)(x_i - \mu)^T \Sigma^{-1}) - \frac{1}{2} \Sigma^{-1} = 0$$

premultiply & postmultiply by  $\Sigma$

$$\sum_i \left[ +\frac{1}{2} (x_i - \mu)(x_i - \mu)^T - \frac{1}{2} \Sigma \right] = 0$$

$$\Rightarrow \sum_i = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T$$

scalar version

$$\begin{cases} \frac{\partial}{\partial x} ax = a \\ \frac{\partial}{\partial x} ax^2 = 2ax \end{cases}$$

vector version

$$\begin{cases} \frac{\partial}{\partial x} a^T x = a \\ \frac{\partial}{\partial x} x^T A x = 2Ax \\ \frac{\partial}{\partial x} x^T A x = 2Ax \end{cases}$$

(if  $A$  symmetric)

$$\frac{\partial}{\partial x} x^T A x = Ax + A^T x$$

( $A$  not symmetric)

## Problem 2.8

function

$$f(x) = \underbrace{\phi(x)^T \theta}_{\text{features of } x} \quad \uparrow \text{parameters}$$

noisy observations

$$y = f(x) + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

$$p(y|x, \theta) = N(y | f(x), \sigma^2)$$

Dataset:  $D = \{(x_i, y_i)\}_{i=1}^N$

MLE

$$\ell(D) = \sum_{i=1}^N \log p(y_i|x_i, \theta)$$

$$= \sum_i -\frac{1}{2\sigma^2} (y_i - f(x_i))^2 - \underbrace{\frac{1}{2} \log \sigma^2}_{\times} - \underbrace{\frac{1}{2} \log 2\pi}_{\times}$$

$$= \sum_i -(y_i - f(x_i))^2$$

$$\underset{\theta}{\operatorname{argmax}} \ell(D) = \underset{\theta}{\operatorname{argmin}} -\ell(D)$$

$$= \underset{\theta}{\operatorname{argmin}} \sum_i (y_i - f(x_i))^2 \quad \text{in least squares objective.}$$

$$= \underset{\theta}{\operatorname{argmin}} \sum_i (y_i - \underbrace{\phi(x_i)^T \theta}_{\Phi_i})^2$$

$$= \underset{\theta}{\operatorname{argmin}} \|y - \Phi^T \theta\|^2$$

polynomial:  $\phi(x) = \begin{bmatrix} 1 \\ x \\ \vdots \\ x^p \end{bmatrix}$

$$\theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_p \end{bmatrix}$$

$$\underset{\theta}{\operatorname{argmin}} \|y - \Phi^T \theta\|^2 = \underset{\theta}{\operatorname{argmin}} (y - \Phi^T \theta)^T (y - \Phi^T \theta)$$

$$= \underset{\theta}{\operatorname{argmin}} y^T y - 2y^T \Phi^T \theta + \underbrace{\theta^T \Phi \Phi^T \theta}_{\alpha^T x}$$

$$\frac{\partial}{\partial \theta} = -2\Phi y + 2\Phi \Phi^T \theta = 0$$

$$(\Phi \Phi^T)^{-1} \times (\Phi \Phi^T \theta = \Phi y)$$

$$\boxed{\theta^* = (\Phi \Phi^T)^{-1} \Phi y}$$

(assume  $\Phi \Phi^T$  is invertible)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder