

lecture 2

• parameter Estimation

How do we find a probability distribution for a r.v. X

Three steps:

1) choose a parametric model (e.g. Gaussian)
 $\theta = \text{parameters}$

2) collect samples from r.v. X :

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

we assume x_i 's are independent; x_i are iid (independant identically distribution) samples

3) maximum likelihood principle:

the optimal parameter θ^* is that which maximizes the probability (likelihood) of the training data

$$\theta^* = \underset{\theta}{\operatorname{argmax}} p(\mathcal{D} | \theta)$$

$\underset{\theta}{\operatorname{argmax}}$ ↗ likelihood of data w.r.t. θ .

a.k.a.: "likelihood function"

$$= \underset{\theta}{\operatorname{argmax}} \underbrace{\log p(\mathcal{D} | \theta)}$$

$L(\theta) = \log \text{likelihood function}$

$$= \underset{\theta}{\operatorname{argmin}} \underbrace{-\log p(\mathcal{D} | \theta)}$$

negative log likelihood function (neg log)

Note: \mathcal{D} is known, so $p(\mathcal{D} | \theta)$ is a function of

θ . It is not a probability w.r.t θ .

and $\log = \text{natural log} = \ln$

- data LL (Log likelihood)

$$L(\theta) = \log P(D|\theta)$$

independent assumption

$$= \log \prod_{i=1}^n P(x_i|\theta)$$

$$L(\theta) = \sum_{i=1}^n \log P(x_i|\theta)$$

To get the MLE solution

- if θ is a scalar. at local optimum

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

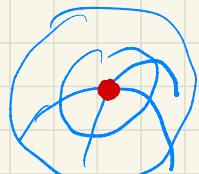
$$\Rightarrow \frac{\partial^2}{\partial \theta^2} \log P(D|\theta) < 0 \text{ at } \theta^*$$

(local maximum ; concave) : [2]

- 3) check the boundary conditions of θ (if necessary)

- If θ is a vector:

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



$$\Rightarrow \nabla_{\theta}^2 L(\theta) = \begin{bmatrix} \frac{\partial^2}{\partial \theta_1^2} & \cdots & \frac{\partial^2}{\partial \theta_1 \partial \theta_p} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial \theta_p \partial \theta_1} & \cdots & \frac{\partial^2}{\partial \theta_p^2} \end{bmatrix} L(\theta) \not\succeq 0$$

Hessian

This Hessian need to be negative definite

to make sure the local maximum

$H \leq 0$: negative definite: $\theta^T H \theta < 0, \forall \theta$

"mountain" concave in all directions

$H \geq 0$: positive definite: $\theta^T H \theta > 0, \forall \theta$

"bowl" convex in all directions

凸函数: 任意两点连成的线段皆位于图像的上方

Example: Bernoulli

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

$$L(\theta) = \sum_{i=1}^N \log P(x_i | \theta) = \sum_{i=1}^N \log [\pi^{x_i} (1-\pi)^{1-x_i}]$$

$$= \sum [x_i \log \pi + (1-x_i) \log (1-\pi)]$$

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

$$m = \sum x_i \leftarrow \text{'sufficient statistic' } = L(\theta) \text{ only}$$

depends on the N observations (dataset)

through this value

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

find the max:

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

$$(1-\pi)m - \pi(1-m) = 0$$

$$\Rightarrow \bar{x}_v = \frac{m}{N} = \frac{1}{N} \sum_{i=1}^N x_i \quad (=1)$$

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

$$2) \frac{\partial}{\partial \pi} \left(\frac{\partial}{\partial \pi} L(\theta) \right) = \frac{\partial}{\partial \pi} \left(\frac{m}{\pi} - \frac{N-m}{1-\pi} \right) \\ = -\frac{m}{\pi^2} - \frac{N-m}{(1-\pi)^2} \leftarrow \text{(check)} \\$$

$$3) \text{boundary condition: } 0 \leq m \leq N$$

$$\Rightarrow 0 \leq \frac{m}{N} \leq 1 \quad (\text{check})$$

Example: Gaussian

$$\theta = \mu \quad (6^2 \text{ known}) \quad (\text{case 1})$$

$$L(\theta) = \sum \log P(x_i | \theta)$$

$$= \sum \left[-\frac{1}{2} \log 2\pi v - \frac{1}{2} \log 6^2 - \frac{1}{26^2} (x_i - \mu)^2 \right]$$

$$= -\frac{N}{2} \log 2\pi v - \frac{N}{2} \log 6^2 - \frac{1}{26^2} \sum (x_i - \mu)^2$$

what are the sufficient statistics?

$$\left\{ \sum_i x_i, \sum_i x_i^2 \right\}$$

$$\hookrightarrow \frac{\partial L}{\partial \mu} = -\frac{1}{2\sigma^2} \sum 2(x_i - \mu) (-1) = 0$$

$$\Rightarrow \sum (x_i - \mu) = 0$$

$$\Rightarrow \hat{\mu} = \frac{1}{N} \sum x_i \quad (\text{sample mean})$$

$\theta = \sigma^2$ (μ is known) (case 2)

$$(x \neq \sigma^2) \cdot \underbrace{\frac{\partial L}{\partial \sigma^2}}_{\cdot} = -\frac{N}{2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4} \sum (x_i - \mu)^2 = 0$$

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

"sample variance"

M.v. Gaussian ..

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

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

• Estimators

the Estimator (e.g. $\hat{\mu}$) is a number

the Estimator is a r.v. (over possible datasets)

$$\text{estimator } f(x_1, \dots, x_N) = \frac{1}{N} \sum_{i=1}^N x_i$$

(x_i is r.v. for each sample, $x_i \sim p(x_i | \theta)$ true def.)

the estimate is the value of the estimator
for a given dataset D

$$\hat{\theta} = f(X_1, X_2, \dots, X_N) \Big|_{X_1 = x_1, \dots, \underset{\substack{\uparrow \\ \text{Sample}}}{X_N = x_N}} = \frac{1}{N} \sum_i^N x_i$$

Since the estimator is a r.v., we can derive the mean and the variance to qualify the "goodness"

- Bias and variance $\hat{\theta} = f(X_1, \dots, X_N)$

1) Will it converge to the true value of θ

$$\text{Bias}(\hat{\theta}) = E_{X_1, \dots, X_N} [\hat{\theta} - \theta] = \underline{E_X [\hat{\theta}]} - \theta$$

if the Bias is non-zero, then we can never get the true value even if infinite samples
means the estimator

2) How long will it take to converge?

(How many samples do we need?)

$$\text{Var}(\hat{\theta}) = E_{X_1, \dots, X_N} [(\hat{\theta} - E\hat{\theta})^2]$$

Example: Gaussian

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

$$\text{Mean of } \hat{\mu} = E_{X_1, \dots, X_N} \left[\frac{1}{N} \sum_i^N x_i \right]$$

$$= \frac{1}{N} E_{X_1, \dots, X_N} [\bar{x}]$$

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

Bias of $\hat{\mu} = 0$

$$\text{var of } \hat{\mu} = E_{X_1, \dots, X_N} \left[(\hat{\mu} - E[\hat{\mu}])^2 \right]$$

$$= E \left[\left(\frac{1}{N} \sum_i (x_i - \mu) \right)^2 \right]$$

$$= \frac{1}{N^2} E \left[\left(\sum_i (x_i - \mu) \right)^2 \right]$$

$$= \frac{1}{N^2} E \left[\sum_i \sum_j (x_i - \mu)(x_j - \mu) \right] = E \left[\frac{1}{N^2} \left(\sum_i (x_i - \mu) \right)^2 \right]$$

$$i=j \Rightarrow E[(x_i - \mu)^2] = 6^2$$

$$i \neq j \Rightarrow E(x_i - \mu)(x_j - \mu) = 0 \leftarrow \text{independent}$$

$$= \frac{1}{N^2} (N \cdot 6^2) = \frac{6^2}{N} = \text{var}(\hat{\mu})$$

variance converges to 0 as $N \rightarrow \infty$

Gaussian variance (PS 2-12)

$$E(\hat{\sigma}^2) = \frac{N}{N-1} 6^2 \Rightarrow \text{Bias}(\hat{\sigma}^2) = \frac{1}{N} 6^2 \neq 0$$

to make it unbiased

$$\hat{\sigma}^2 = \frac{N}{N-1} \hat{\sigma}^2 = \frac{N}{N-1} \cdot \frac{1}{N} \sum (x_i - \mu)^2 = \frac{1}{N-1} \sum (x_i - \mu)^2$$

Important Asymptotic Properties of MLE

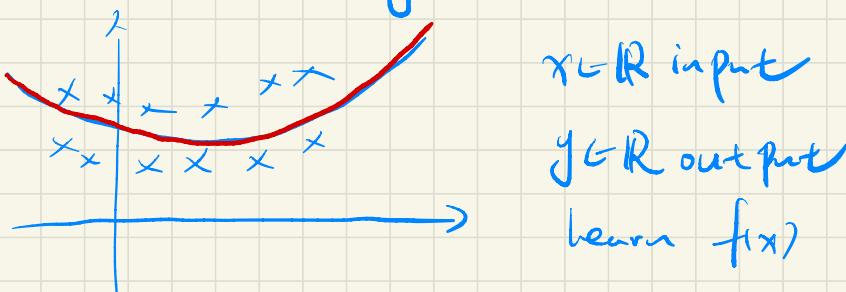
1) Consistent - As $N \rightarrow \infty$, the estimate converge to the true value. Asymptotic

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

(CRLB is a theoretical bound on the variance of any unbiased estimator for a given $P(x|\theta)$)

No unbiased estimator can get lower variance than MLE.

• MLE for Regression.



Consider a Polynomial function (kth order)

$$f(x, \theta) = \sum_{d=0}^N x^d \theta_d = \begin{bmatrix} 1 & x & x^2 & \dots & x^N \end{bmatrix}^T \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_N \end{bmatrix} = \phi(x)^T \theta$$

$$\overbrace{\phi(x)}^{\text{function}} \quad \overbrace{\theta}^{\text{parameters}}$$

observe a noisy output:

$$y = f(x, \theta) + \xi \quad \underbrace{\xi}_{\text{noise}} \sim N(0, \sigma^2) \text{ i.i.d}$$

equivalently (y is a r.v.)

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

Given dataset $\{(x_i, y_i)\}_{i=1}^N$, estimate θ

using MLE:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \sum_i \log P(y_i|x_i, \theta)$$

$$= \underset{\theta}{\vdots}$$

$$= \underset{\theta}{\operatorname{argmin}} \sum_i (y_i - f(x_i, \theta))^2 \quad \text{least square formulation}$$

$$= \underset{\theta}{\operatorname{argmin}} \|y - \Phi^T \theta\|^2, \quad \Phi = [\phi(x_1), \dots, \phi(x_N)]$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

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

Notes:

\Rightarrow MLE is more general than LS

2) Assumptions are explicit

i) Gaussian noise

ii) $\mu = 0$, σ^2 variance (fixed)

iii) noise is iid

3) MLE can describe other LS formulations

i) weighted LS (ps. 2.8)

ii) regularised LS (lecture 3)

iii) L_p-norm (ps 2.9)