

Youtube lecture 26/09/16

- likelihood function

- A means of generating estimators
- Bayesian inference
- sufficiency

define: $x^n = (x_1, \dots, x_n)$ joint density $p(x^n; \theta) = p(x_1, \dots, x_n; \theta)$
 $\theta \in \Theta$

(*) The likelihood function

$L: \Theta \rightarrow [0, \infty)$ (maps parameter space to some numbers)

$$L(\theta) \equiv L(\theta; x^n) = p(x^n; \theta)$$

(*) W: treat x^n as fixed, view likelihood as a function of parameters θ

(*) If data is IID then likelihood is:-

$$L(\theta) = \prod_{i=1}^n p(x_i; \theta) \quad \text{IID case only}$$

(*) If θ is fixed then $p(x^n; \theta)$ is a probability density for x

(*) If x^n is fixed and view $L(\theta)$, it is no longer a probability, but a function of θ .

(*) log-likelihood function

$$l(\theta) = \log L(\theta)$$

(*) - review example 1

W: Probability density fn and likelihood fns live on different space

likelihood functions are an equivalence class of functions
 (defined up to a constant of proportion.)

Example 2

$x_1, \dots, x_n \sim N(\mu, 1)$

$$L(\mu) = \prod_{i=1}^n \frac{1}{(2\pi)^{1/2}} \exp^{-\frac{1}{2}(x_i - \mu)^2} = \left(\frac{1}{2\pi}\right)^{n/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 \right\}$$

- Add and subtract \bar{x} ;
 throw away non μ terms

yielding $l(p) \propto \exp\left\{-\frac{n}{2}(\bar{x} - \mu)^2\right\}$

(A1): Review the proport. calculation ↑

(*) this is the likelihood function, a function of p , not x .

tw: Also looks like normal p.d. centered at \bar{x} ; in this case then they are the same; but in general, likelihood and p.d. are distinct

example 3

- let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$

Ⓢ

tw: Reversing roles; fix X at obs. values
 $X = x \dots$ and request prob. of data
as function of p .

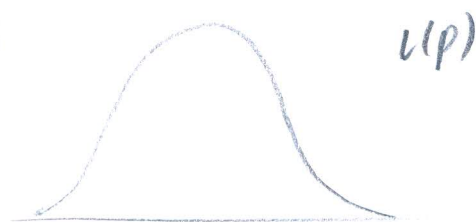
$$(*) l(p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} \propto p^X (1-p)^{n-X}$$

where $X = \sum_{i=1}^n x_i$ (sum of coin flips)

tw: sample space - sequences of 0s and 1s

parameter space - $[0, 1]$

- discrete probability function but
- continuous likelihood function



example 5

- let $X_1, \dots, X_n \sim \text{Unif}(0, \theta)$
 $\theta > 0$

PDF

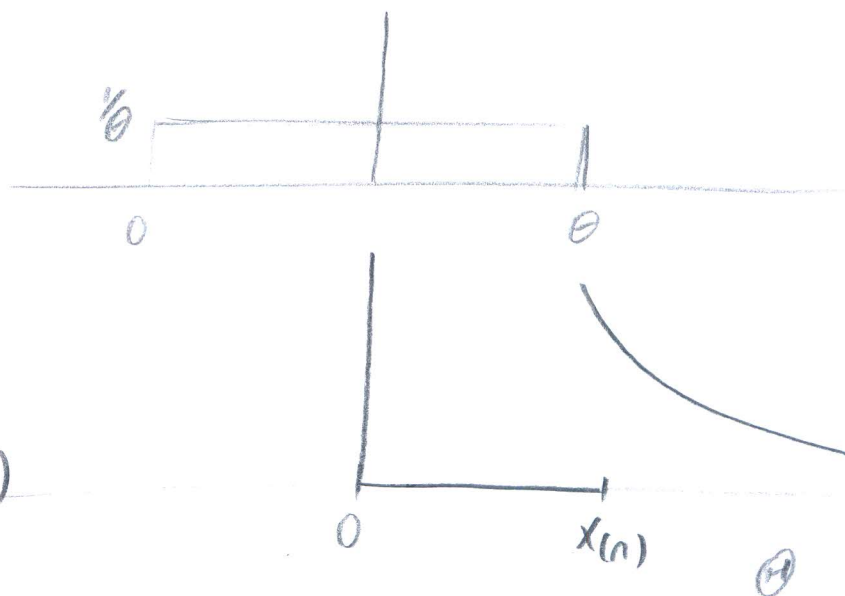
$p(x_i; \theta)$

tw: good way to think about likelihood function.

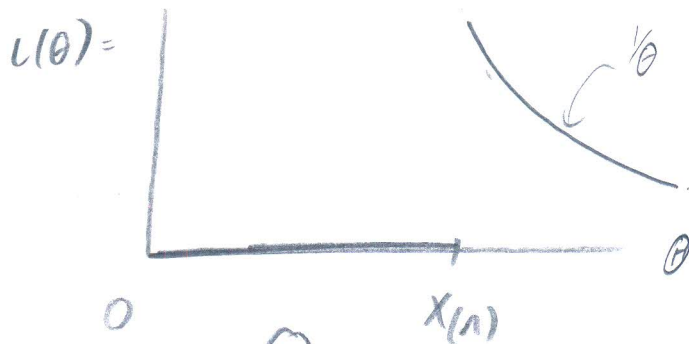
For every possible θ ,
what is probability
of getting a
particular dataset. (for
this) a particular θ

$$l(\theta) = \prod_{i=1}^n p(x_i; \theta)$$

$X_{(n)}$ - maximum observation.



$$L(\theta) = \begin{cases} 0 & \theta < x_{(n)} \\ \left(\frac{1}{\theta}\right)^n & \theta > x_{(n)} \end{cases}$$



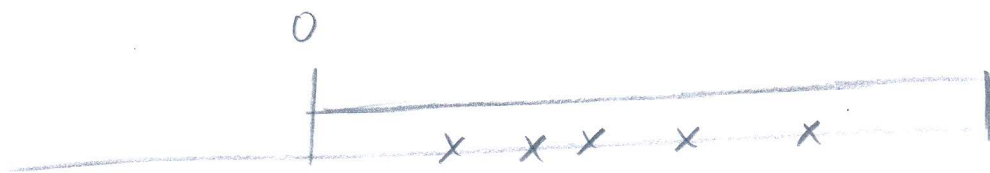
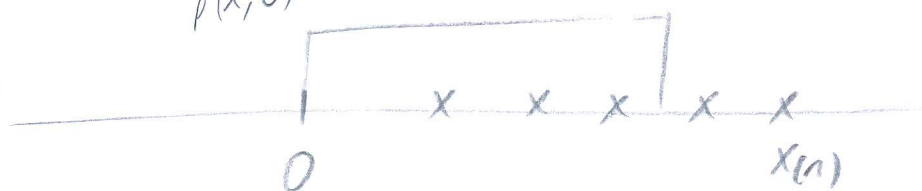
(?) Found this tricky; don't understand intuition (A2)

W: If $\theta \in [0, x_{(n)}] \rightarrow$ could not have gotten dataset

$\theta \in [x_{(n)}, \infty) \rightarrow$ non-zero probability of getting dataset (at that value of θ)

PDF
 $p(x, \theta)$

$\theta < x_{(n)}$



- okay; a little clearer intuitively

(*) Note distinction between PDF and likelihood.

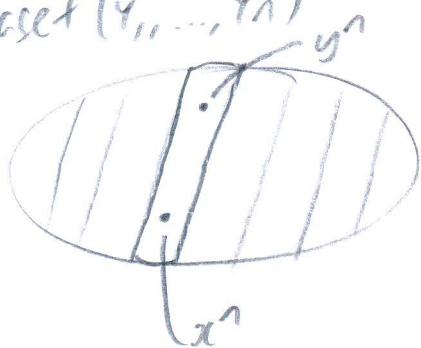
- notation: - likelihood also written $L_{x_1, \dots, x_n}(\theta)$ to remind reader of dataset

(*) Two likelihoods are the same if they are proportional \rightarrow gives partition.

$$(x_1, \dots, x_n) \sim (y_1, \dots, y_n) \text{ if } L_{x_1, \dots, x_n}(\theta) \propto L_{y_1, \dots, y_n}(\theta)$$

- 2 datasets are in the same partition of a partitioned sample space if the likelihood of dataset (x_1, \dots, x_n) is proportional to likelihood of dataset (y_1, \dots, y_n)

W:



• Take set of all possible datasets

• create partitions so

that datasets are in the same partition if likelihood of 2 datasets are the same shape.

- creates an equivalence class/partition.

(*) This is the minimal sufficient partition

- sufficiency and minimal sufficiency related via likelihood fn

(*) Asking if a statistic is sufficient

↳ can I compute the likelihood function from this statistic?

- YES \rightarrow it is a sufficient statistic

e.g. \bar{X} is a sufficient statistic as I can compute likelihood (for Normal)

entire dataset is — " ————— " —————

(*) minimal sufficiency \rightarrow that partition induced by statistic is in a se minimal reduction of data

AND

\rightarrow in a sense, smallest amount of info needed to construct likelihood

(*) Informally, likelihood function itself is minimal sufficient statistic

lw: likelihood induces a partition, telling you two datasets have the same likelihood; that partition is minimal sufficient

(*) (6). The minimal sufficient statistic has all information you need to compute the likelihood function ^{estimating the}

- whether likelihood is suitable for parameters to be estimated depends - different question

lw: many likelihoods need to be computed numerically
(impractical)

example 6

/ symm. P.D.

- $X_1, X_2, \dots, X_n \sim N(\mu, \Sigma)$

likelihood:

$$L(\mu, \Sigma) \propto |\Sigma|^{-n/2} \exp\left(-\frac{1}{2} \sum_i (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)\right)$$

- maximising over μ and Σ is a little nasty.

- consider $\ell(\mu, \Sigma)$

$$\ell(\mu, \Sigma) = -\frac{n}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

Alt - Review calc.

$$\sum_i (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \quad \text{- relaxing vector not.}$$

$$= \sum_i (x_i - \bar{x})^T \Sigma^{-1} (x_i - \bar{x}) + n(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

quad. form
i.e. scalar

$$= \sum_i \text{tr} \{ (x_i - \bar{x})^T \Sigma^{-1} (x_i - \bar{x}) \} + n(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

$$= \sum_i \text{tr} \{ \Sigma^{-1} (x_i - \bar{x}) (x_i - \bar{x})^T \} + n(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

$$= \text{tr} \sum_i \{ \Sigma^{-1} (x_i - \bar{x}) (x_i - \bar{x})^T \} + n(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

$$= \text{tr} \Sigma^{-1} \sum_i (x_i - \bar{x}) (x_i - \bar{x})^T + n(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

$$= n \text{tr} \Sigma^{-1} \frac{1}{n} \sum_i (x_i - \bar{x}) (x_i - \bar{x})^T + n(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

$$= n \text{tr} \Sigma^{-1} S_n + n(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

- Add/subtract \bar{x}

- cross prod. 0

- $\text{tr}(A) = \sum \text{diag. for sq.}$

- $\text{tr}(a) = a$

$$\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$$

- Σ^{-1} does not contain index i

$$S_n = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) (x_i - \bar{x})^T$$

(AS) - Find out correct typo.

Here;

$$\ell(\mu, \Sigma) = -\frac{n}{2} \log |\Sigma| - \frac{n}{2} \text{tr}(\Sigma^{-1} S) - \frac{n}{2} (\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)$$

- maximising this:-

$$\text{set } \mu = \bar{x} \text{ and } \Sigma = S$$

- LW: why focus on maximisation?

- LW: lower section 2. likelihood, sufficiency & likelihood Principle
after we review lecture notes 7 (later)

Lecture Notes 7 - Point Estimation

1. Introduction

WAG: clarify if r.v. is vector or scalar.

- $X_1, \dots, X_n \sim p(x; \theta)$. Want to estimate $\theta = (\theta_1, \dots, \theta_k)$ given data

(*) Take a fraction of data i.e. a statistic (but not quite)

$$\hat{\theta} = \hat{\theta}_n = w(X_1, \dots, X_n)$$

- $\hat{\theta}$ - indicates estimate of param

- $\hat{\theta}_n$ - remainder of sample size

EN

(*) Distinction between parameter and estimator.

θ - parameter - fixed unknown number (no randomness)

$\hat{\theta}$ - estimator - r.v. - has distrib, mean, variance.

(*) Most common estimator construction principles - parametric point estimation.

1. MOM - method of Moments

2. Maximum likelihood (MLE)

3. Bayesian estimators.

Context

1. - Invented in 1900s \rightarrow computationally v. simple \rightarrow resurgence in ML learning
2. - Most common point estimator for parametric models.
- under certain conditions; it is in a sense, optimal
3. - Bayesian estimators (not Bayesian inference, stats, rule, theorem).
- for some ML models; MLE difficult; MOM easier

WV: How do we know which to use?

- 2 questions \rightarrow i) method for constructing estimators.

\rightarrow ii) evaluating estimators according to some criteria

(*) Evaluating estimators:-

1. consistency \rightarrow if I gave you more data, will your estimator converge in probability to...

2. Bias-variance

3. Mean squared error

4. * minimax theory \rightarrow a way of formulating and evaluating estimators wrt optimality

5. Robustness \rightarrow estimator may be optimal under certain conditions; what if those conditions are relaxed?
- presence of outliers etc. \rightarrow fish in high-dim statistics

W: minimax theory

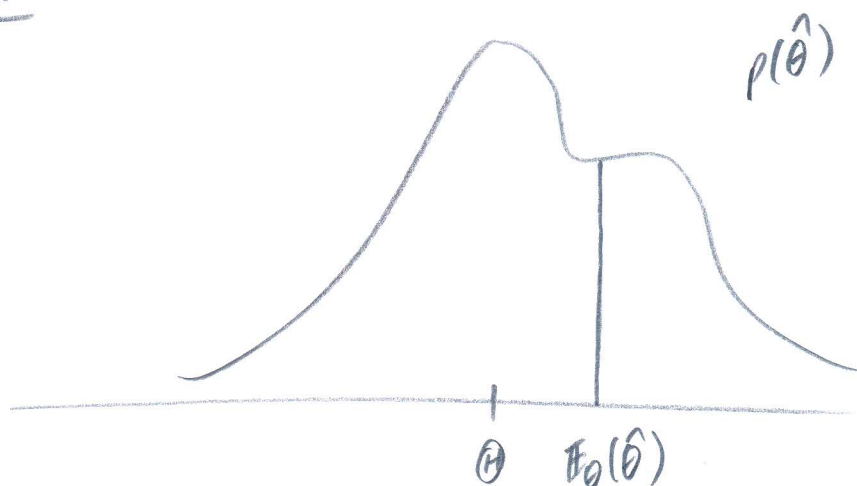
- covered in cursory detail \rightarrow see 36-702

(*) terminology

1. $\mathbb{E}_\theta(\hat{\theta}) = \int \dots \int \hat{\theta}(x_1, \dots, x_n) p(x_1; \theta) p(x_2; \theta) \dots p(x_n; \theta) dx_1 \dots dx_n$

- remarks: subscript $\theta \nrightarrow \theta$ is random - it is a fixed unknown no.
- informs us that we are putting in a particular value of θ when we compute $\hat{\theta}(x_1, \dots, x_n)$; when the value is θ .
- $\hat{\theta}$ is random

2. Bias: $-\mathbb{E}_\theta(\hat{\theta}) - \theta$



3. Sampling distribution
- distribution of $\hat{\theta}_n$

4. Standard error

- standard deviation in special context when r.v. is an estimator

$$\sqrt{\text{var}(\hat{\theta}_n)} = \text{se}(\hat{\theta}_n)$$

5. $\hat{\theta}_n$ is consistent if $\hat{\theta}_n \xrightarrow{p} \theta$ (large sample/asympt. property)
as $n \rightarrow \infty$