

## Interestori

course page: [http://jasonmiao.com](#)

raazesh.sarmadi@math.ucla.edu : subject M305

Kolmogorov basis:  
Theoretical foundation for descriptive stats  
(HT19)

### Sets

A set is a collection of distinct elements.

e.g.  $\{0, 0\}$ . We ignore number of sets, e.g.  $A = \{0, 0\}$ .

$A = \{0, 0\}$  is not a set (confused)  
 $\emptyset = \{\}$  is the empty set.

An element belongs to (or doesn't belong to) a set and we write ' $\in$ ' or ' $\notin$ ' respectively. e.g.  $0 \in \{0, 0\}$  and  $0 \notin \{0, 0\}$ , etc. etc. definition below.

### Set operations

We can add elements to an existing set by union operation.

$$\text{e.g. } \{1\} \cup \{2\} = \{1, 2\}$$

$$\{0, 0\} \cup \{0\} = \{0, 0, 0\}$$

$$\{0, 0\} \cup \{0\} = \{0, 0\}$$

Def  $A \cup B := \{x \mid x \in A \text{ or } x \in B\}$

Intersection:  $A \cap B := \{x \mid x \in A \text{ and } x \in B\}$

Set difference:  $A \setminus B := \{x \mid x \in A \text{ and } x \notin B\}$

### complement

Given a universal set  $U$ ,  $A^c := \{x \mid x \notin A\} = U \setminus A$

## Maps

A map or a function associates each element in the set domain with exactly one element in the set range (codomain).

### Formal definition

A function is a specific kind of relation between elements in the domain and range.

### Large image interpretation

The large image of a function  $f: X \rightarrow Y$  is  $\{f(x) : x \in X\}$  where  $f^{-1}(y) = \{x \in X \mid f(x) = y\}$  for every  $y \in Y$ .

$$f^{-1}(\emptyset) = \{x \in X \mid f(x) \in \emptyset\} \text{ and } f^{-1} = f(Y) = \sigma_f(X)$$

is called the preimage of  $X$

### Note

In this case,  $N = \{1, 2, 3, \dots\}$  and  $Z = \mathbb{Z}_{\geq 0} = \{0, 1, 2, 3, \dots\}$

## Probability

### Language

An experiment is an activity that produces distinct observable outcomes.  
The set of such outcomes is called the sample space of the experiment, denoted by  $\Omega$ .

An event is a subset of the sample space.

Probability is a function

$$P: \{\text{events}\} \rightarrow [0, 1] \quad \text{where } P \text{ is defined}$$

$$\text{i)} \forall \text{ event } A, \quad 0 \leq P(A) \leq 1$$

$$\text{ii)} P(\Omega) = 1$$

$$\text{iii)} A \cap B = \emptyset \Rightarrow P(A \cup B) = P(A) + P(B)$$

$$\text{iv)} P(\bigcup_i A_i) = \sum_{i=1}^n P(A_i) \quad \text{where } A_i \text{ are pairwise disjoint}$$

### Motivation for definition:

Ideas of long-term relative frequency of independent experiments (experiments repeated many times). If we repeat an experiment a large number of times, the fraction of times the event  $A$  occurs will be close to  $P(A)$ .

Formally, let  $N(A)$  be the number of times  $A$  occurs in the first  $n$  trials.

$$P(A) = \lim_{n \rightarrow \infty} \frac{N(A, n)}{n}$$

$$\text{with} \rightarrow i) \quad 0 \leq \frac{N(A, n)}{n} \leq 1$$

$$ii) \quad \frac{N(A, n)}{n} \xrightarrow{n \rightarrow \infty} \text{Something between zero and one}$$

$$iii) \quad A \cap B = \emptyset \Rightarrow N(A \cup B, n) = N(A, n) + N(B, n)$$

iv) If this is ignored the mathematics would be much harder

1.1 Tossing a fair coin. Can each result be defined probability?

$$\Omega = \{H, T\}$$

Bernoulli random variable. Bernoulli( $\theta$ )

$$2^{\Omega} = \mathcal{F}_{\Omega}$$

$\frac{e^{\theta}}{N} \in \text{Lotto } (40 \text{ balls})$

Label of the 1st ball that come in lotto:

$$\Omega = \{1, 2, \dots, 40\}$$

Discrete random variable,

Discrete( $\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n}$ ) - Discrete( $\theta, \theta_1, \dots, \theta_n$ )

$$n=10$$

Here we have  $\theta_1 = \theta_2 = \dots = \theta_{10} = \frac{1}{10}$ , can Equi-probable Discrete variable.



$$2^{\Omega} = \mathcal{F}_{\Omega}$$

$$\begin{aligned} \text{So, what is } P(\text{"one number"}) &= P(\{1, 2, \dots, 40\}) = \\ &= P(\{1\}) + P(\{2\}) + P(\{3\}) + \dots + P(\{40\}) = \\ &= 20 \cdot \frac{1}{40} = \frac{1}{2} \end{aligned}$$

Properties

i.  $P(A) = 1 - P(A^c)$

ii.  $A, B$  events.  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$



The domain of  $P$  is called a sigma field or sigma algebra.  
It's denoted by  $\mathcal{F}(\Omega)$  or  $\mathcal{F}_{\Omega}$  or just  $\mathcal{F}$  if  $\Omega$  is clear from context.

We see that  
i)  $\Omega \in \mathcal{F}$ , ii)  $A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}$ , iii)  $A_1, A_2, \dots \in \mathcal{F} \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$  by Kolmogorov.

Bernoulli experiments (one trial)

The triple  $(\Omega, \mathcal{F}_\Omega, P)$  is called the probability space.  
 $\Omega$  is a set of probabilities,  $\mathcal{F}_\Omega = (\mathcal{F}_\Omega, \mathcal{F}_\Omega, \dots)$  is called a collection  
 of events.

Result

events A, B are independent  $\Rightarrow P(A \cap B) = P(A)P(B)$

The present experiment

$$X_1, X_2, \dots, X_n \in \Omega$$

Trial 1, Trial 2, ..., Trial n

Conditional probability

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

provides  $P(B) > 0$

conditional probability  $\Rightarrow$  prob. 1. Yes, basically B is given A

Graph:  
 Constructing random graph from uniform dist.

Toss coin with  $P(H) = \Theta(1/n^2)$  n times.

Graph:  
 Let  $V = \{v_1, \dots, v_n\}$  be set of vertices and let  $E \subseteq V^2$  and  $|E| = k$   
 (n vertices, k edges). Let

edges  $e_1, e_2, \dots, e_k$  where  $A_{ij} \in \{0, 1\}$  if  $e_j$  connects  $v_i$  and  $v_j$ .

$$\Pr[A] = \prod_{j=1}^k \Pr[e_j]$$

$$\text{so, } E(|E|) = |V|^2$$

Whichever is very close to central.

$$pp_{\text{min}} = \left[ \pi_{16} \cdot D(p_0(x)) + i \right] \text{ for } i$$

$$p(0) = \frac{1}{11}, \quad p(1) = \frac{14}{11}, \quad p(2) = \frac{15}{11}, \quad p(3) = \frac{14}{11}, \quad p(4) = \frac{1}{11}$$

$$(b) \quad \tau = \frac{V_0 L^2}{2 \cdot 100 \cdot \sigma} = \frac{100 \cdot \left(\frac{1}{2}\right)^2}{2 \cdot 100 \cdot 0.01} = 25 \text{ s}$$

卷之三

25.  $\frac{d}{dx} \int_{0}^{x^2} f(t) dt =$  \_\_\_\_\_

Exodus 15

## Principle to Decision Theory

### Estimation in parametric models



Data = Random variable ( $R.V.$ )

Probability distribution ( $P.D.F.$ )

### Problem

Let  $X \in R.V$  with c.d.f. in  $\mathbb{X}$  and  $L(x) = \text{Loss}(x)$  is loss of  $x$ .

The parameter space  $\Theta$  (bold italic):

$L(X) \in \{P_\theta | \theta \in \Theta\}$ , here we assume  $\Theta \subseteq \mathbb{R}$  for simplicity.

The decision problem is to estimate a function  $g(\theta)$  based on  $n$  realizations of  $X$ .

### Topically (Outline)

$X = (X_1, \dots, X_n), X_i \sim X_i$

$n$  is called sample size, initially  $X$  is countable or  $\subseteq \mathbb{R}^n$ .

### Def

A statistic  $T$  is an arbitrary function of the observed  $R.V. X$  (italic)

$$\{t_{\theta}(x) | \theta \in \Theta\}$$

### Def

An estimator  $T$  of  $g(\theta)$  we denote by  $T: \mathbb{X} \rightarrow g(\Theta)$

Concretely:

$\begin{array}{c} \text{I} \\ \downarrow \\ \mathbb{X} \\ \downarrow \\ \text{II} \end{array}$

$\begin{array}{c} \text{I} \\ \downarrow \\ P_\theta, \eta_\theta \\ \downarrow \\ \text{III} \end{array}$

$\begin{array}{c} \text{I} \\ \downarrow \\ T \\ \downarrow \\ \text{III} \end{array}$

gives us an estimate of  $g(\theta)$

$$T = g(\Theta)$$

## Indicator Function

$$\text{Law}(X) \sim \text{Bernoulli}(\theta) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{if } x \notin A \end{cases} \quad \text{is Bernoulli}$$

$$\text{Law}(X) \sim \text{Bernoulli}(\theta) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{if } x \notin A \end{cases} \quad \text{is Bernoulli}$$

Suppose the data vector  $x = (x_1, \dots, x_n)$  is a realization of  $X \sim \text{Binomial}(\theta)$ , i.e.  $x \in \mathbb{R} \cdot \{0, 1\}^n$ , for unknown  $\theta \in \Theta = [0, 1]$ .

Note that  $\mathbb{P}_\theta(x_{i,1}) = \prod_{j=1}^n \theta^{x_{j,1}} (1-\theta)^{1-x_{j,1}} = \theta^{x_{1,1}} (1-\theta)^{1-x_{1,1}}$ .

$$\text{Let } T(x) = \sum_{i=1}^n x_{i,1}$$

Thus, the prob. of data (i.e.  $\mathbb{P}_\theta(x_{1:n})$ ) only depends on the statistic  $T$ .

Now, consider another statistic: (sample mean)

$$T(\theta) = \frac{1}{n} \sum_{i=1}^n x_{i,1} = \bar{X}_{n,1}$$

Then  $\theta$  becomes  $\theta^{\bar{x}}(1-\theta)^{1-\bar{x}} = \theta^{\bar{x}_{1:n}} (1-\theta)^{1-\bar{x}_{1:n}}$

An estimator of  $\theta$ ,  $\hat{\theta}(\theta)$ , (say  $g(\theta)$ ) should converge towards to the "true" but unknown  $\theta$  to be estimated, as the sample size  $n \rightarrow \infty$

Def

A sequence  $\bar{T}_n := \bar{T}_n(x_1, \dots, x_n)$  of estimators (each based on a sample of size  $n$ ) for a parameter  $\theta$  is called (asymptotically) consistent if  $\forall \delta > 0 : P_{\theta,n}(|\bar{T}_n(x^{(n)}) - \theta| > \varepsilon) \rightarrow 0$  as  $n \rightarrow \infty$  or, in shorter notation,

$$\bar{T}_n \xrightarrow[n \rightarrow \infty]{P} \theta \quad \text{if } \text{Law}(X^{(n)}) = P_\theta$$

or (non product ex)

The estimator  $T_n(x_1, \dots, x_n) = \bar{x}_n$  is consistent for  $\theta$

Proof

Since  $x_1, \dots, x_n$  are iid  $B(1, \theta)$  R.V.s, we have  $E(x_i) = \theta$  and the result follows from the law of large numbers.

So, consistency can be seen as a minimal requirement for estimators. But this still leaves a lot of consistent estimators to choose from. A quantitative comparison of estimators is made possible by the approach of statistical decision theory.

We choose a loss function  $\text{Loss}(t, \theta)$  which measures the loss (inaccuracy) of the unknown parameter  $\theta$  is estimated by  $t$ . This can be done by different choices of  $\text{Loss}(T(\bar{x}), \theta)$ .

Absolute error:  $\text{Loss}(t, \theta) = |t - \theta|$

Quadratic error:  $\text{Loss}(t, \theta) = (t - \theta)^2$

or  $\text{Loss}(t, \theta) = 1_{\{t > \theta\}}(t - \theta)$  for some  $\delta > 0$  to emphasize that distance being less than  $\delta$

Note

$\text{Loss}$  is a Q.N.  $\Leftrightarrow \text{Loss}(T(x), \theta)$  needs to account for random errors.

Def

The risk of an estimator  $T$  of parameter  $\theta$  is

$$R(t, \theta) := E_\theta(\text{Loss}(T(x), \theta))$$

A risk function of  $T$

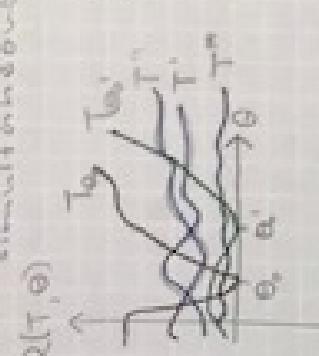
Note Risk might exist

Note Risk might not exist since the expectation might not exist

class of

$$Q(\tau, \theta) = \lim_{n \rightarrow \infty} Q(\tau, \theta) \quad \text{for any } \tau, \theta \in \Theta$$

We say  $\theta_0$  is an estimator of  $\theta$ . We minimize the unbiased function



a uniformly best estimator (in  $\mathcal{X}$ ).

Argument

$$\text{For each } \theta_0, \theta \in \Theta, \text{ consider } E_\theta T = Q_\theta$$

$$\text{so } P(T_{\theta_0}, \theta_0) = 0 \Rightarrow \text{bad if } \theta_0 \text{ is true}$$

If  $T$  is TRUE, then it would have to compete with  $T_{\theta_0}$ .

Unbiased estimators

Consider an estimator  $T$  s.t.  $E_\theta T$  exists.

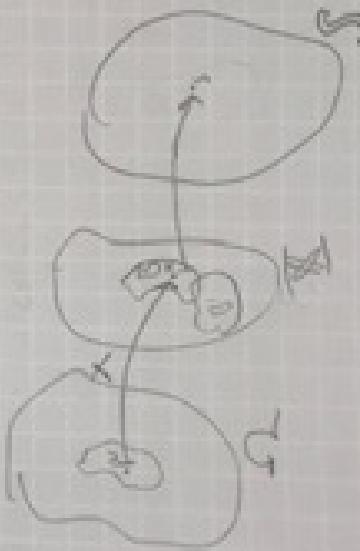
Def

Then  $E_\theta(T) - g(\theta)$  is called bias of the estimator.  
 $\text{If } E_\theta(T) = g(\theta) \text{ for any } \theta \in \Theta$ , then the estimator  $T$  is called unbiased

for  $g(\theta)$

Def

A statistic  $S$  is called sufficient for  $\theta$  if  $P_\theta(x, B|S(x)=s)$  is independent of  $\theta$ , for all values of  $s$  and  $x$  such that



In words, conditional distribution of  $X$  given  $S(x)=s$  does not depend on the parameter  $\theta$

For discrete experiments

$$P_\theta(X, B|S(x)=s) = \begin{cases} \frac{P_\theta(X=x \cap B|S(x)=s)}{P_\theta(S(x)=s)} & \text{if } P_\theta(S(x)=s) > 0 \\ 0 & \text{if } P_\theta(S(x)=s) = 0. \end{cases}$$

Let's clarify again:

Originally, the law of  $X$  depended on  $\theta$  ( $L_{\theta}(x) = P_\theta(x)$ ).

After the value of the sufficient statistic  $S(x)$  is known,  
then the condition now  $P_\theta(\cdot | S(x)=k)$  is no longer dependent  
on  $\theta$ .

Since we are interested in making inference about  $\theta$ ,  
the conditional law is interesting. For our purpose, it will be  
convenient:

After having  $S(x)$  into account, the remaining randomness  
does not depend on  $\theta$  anymore.

### Remark (Exercise to prove Q.E.D.)

The data itself is sufficient, i.e.

if  $S(x)=x$ , then for  $\theta \in \Theta$ ,  $P_\theta(x|S(X)=x) = P_\theta(x) = \begin{cases} 1, & x=k \\ 0, & x \neq k \end{cases}$

Proof

$$I_n(\bar{x}; \theta)(\theta) \text{ s.t. the sample mean } \bar{x}_n \text{ is a sufficient statistic.}$$

Proof

$$n\bar{x}_n = \sum_{i=1}^n x_i \sim \text{Bin}(n, \theta)$$

Suppose  $\bar{x}_n = k$  where  $k$  is one of the possible values. This means  
 $n\bar{x}_n = k$  for some  $k \in \{0, 1, \dots, n\}$ . Then, for any  $x = (x_1, \dots, x_n)$ ,

$$P_\theta(x_1=x_1, x_2=x_2, \dots, x_n=x_n | n\bar{x}_n=k) = \frac{P_\theta(x_1=x_1, x_2=x_2, \dots, x_n=x_n)}{P_\theta(n\bar{x}_n=k)} \quad (\textcircled{X})$$

If  $\theta \in \mathcal{U}$ ,  $\sum_{i=1}^n x_i = k$ , then  $\theta$  is  $\frac{(\theta^k)(1-\theta)^{n-k}}{(\theta^k)(1-\theta)^{n-k}} = \frac{1}{\binom{n}{k}}$  which clearly is  
indep. of  $\theta$ .

$\square$  If for this  $X_i \sim \sum_{k=1}^n k \pi_k$ , then the numerator is 0 which is indep of  $\theta$ , so  $\Theta$  is indep of  $\theta$ .  $\square$

Say we want to estimate  $\theta$  in this example and we limit ourselves to estimators that are functions of the sufficient statistic  $\bar{X}_n$ :

$$\tau(x) = h(\bar{X}_n)$$

Additionally, suppose we limit ourselves further to unbiased estimators:

$$E(\tau(X)) - \theta = 0 \Rightarrow E(\tau(X)) = \theta$$

Proof.

Under sufficiency and unbiasedness restrictions on allowed estimators of  $\theta$  in our  $\mathbb{R}[\Theta]$  exp., the only possible estimator is  $\bar{X}_n$ .

Proof.

$$\begin{aligned} 0 &= E_\theta h(\bar{X}) - \theta \stackrel{\downarrow}{=} E_\theta (h(\bar{X}_n) - \bar{X}_n) = \overbrace{\sum_{k=0}^n c_k (h(\frac{k}{n}) - \frac{k}{n})}^{C_k} \theta^{(1-\theta)^n} = \\ &= (-\theta)^n \sum_{k=0}^n C_k n^k \text{ for } n = \frac{\theta}{1-\theta} \end{aligned}$$

Now, if  $\theta \in [0, 1)$ , then  $n \in \mathbb{N}, \infty$ ), hence the above polynomial can only be zero if every  $C_k$  is also zero.

Thus,  $0 = \sum_{k=0}^n (h(\frac{k}{n}) - \frac{k}{n}) \Rightarrow h(\frac{k}{n}) = \frac{k}{n}$  for  $k \in \{1, \dots, n\}$  and  $h(\bar{X}_n) = \bar{X}_n$  for all possible values of  $\bar{X}_n$ .  $\square$

Def

A statistic  $T$  is called complete if, for all  $h: \mathbb{T} \rightarrow \mathbb{R}$ :

$$(E_\theta(h(T(X))) = 0 \text{ for all } \theta \in \Theta) \Rightarrow$$

$$\Rightarrow P(\underbrace{h(T(X))=0}_{\text{almost surely}}) = 1 \text{ for all } \theta \in \Theta$$

In statistics, completeness means there is "no superfluous information" or "no redundancy" - the complete statistic  $T$

Consider an event  $T(X) \in \mathcal{B}$  and suppose  $P_\theta(T(X) \in B) = c$  for  $c$  indep. of  $\theta$ .

By taking  $h(\cdot) = \mathbf{1}_B(\cdot)$ , completeness  $\Rightarrow P_\theta(\mathbf{1}_B(T(X)) = c) = 1$  for all  $\theta \in \Theta$  which means that  $c$  is either 0 or 1. Thus, for any such  $T(X) \in \mathcal{B}$  which has a non-trivial probability ( $\neq 0, 1$ ), this prob. must depend on  $\theta$ .

Thm  $\otimes_{\mathcal{B}_\theta} (\Theta)$

The statistic  $T(X) = (X, \bar{X})$  is sufficient but not complete.

$h(X, \bar{X}) = X, -\bar{X}$  has  $E(h) = 0$  but  $h(X, \bar{X})$  is not almost surely 0.

Prop

I.  $\otimes_{\mathcal{B}_\theta} (\Theta)$ ,  $T(X) = \bar{X}_n$  is sufficient and complete.

Proof

Suppose for some function  $h$ ,

$$E[h(\bar{X}_n)] = 0$$

This means  $\mathbb{E} + \sum_{k=0}^{\infty} h\left(\frac{k}{n}\right)\left(\frac{n}{k}\right)\theta^k(1-\theta)^{n-k} = 0 \quad \forall \theta \in [0, 1]$ .

Then,  $h\left(\frac{k}{n}\right) = 0$  for  $k \in \{0, 1, \dots, n\}$  as per the earlier argument used to show that  $\bar{X}_n$  is the only unbiased and sufficient estimator

L114 - Conv. of sentence of R.N. X:  
X, ..., X, conv. to C.R.N. X:

X, D.A.S.H.; b.m.s.  
X, m.s.

-Marlboro's man, Chayhleben's men.

X, p.d.k

T - Prod - B.W.B.

151-156

③ Suppose  $T$  is a sufficient statistic with values in a set  $\bar{T}$   
and  $S: \bar{T} \rightarrow S$  is a one-to-one mapping with values in  $S$  (ie,  
there exists an inverse mapping  $S^{-1}: S \rightarrow \bar{T}$  such that  $S^{-1}(S(t)) = t$  for  
each  $t \in \bar{T}$ )  
Show that the statistic  $S(T(x))$  is sufficient.

④ In class it was claimed that the statistic  $T(x) = (x, \bar{x}_n)$  is  
sufficient for the  $\bigoplus_{i=1}^n$  Bernoulli( $\theta$ ) experiment. Prove this  
claim.

### Problem Set Week 3

- ① Show that the estimator  $\bar{T}_n$  is consistent in  $\theta = \bigcup_{i=1}^n \text{Bernoulli}(p_i)$  expand
- ② Suppose the data  $X$  in a statistical experiment can take values in a countable set  $\tilde{X}$  (i.e.,  $\text{Law}(X)$  is discrete).  
In class it was claimed that the data itself are a sufficient statistic (*i.e.*,  $T(X) = X$  is sufficient). Write down the argument that proves this claim (it can be a short paragraph).

- ③ Let  $X_1, \dots, X_n$  be independent and identically distributed with  $\text{Inv}(h)\text{-Poisson}(\lambda)$ ,  $\lambda \in (0, \infty)$  is unknown. Show that the sample mean  $\bar{X}_n$  is a sufficient statistic.

## Limits of R.V.s : Chapter 8 in CSE Book.pdf 161-156

8.1 - conv. of sequence of R.V.

$X_1, \dots, X_n$  conv. to  $\bar{X}$  R.V.  $X$ :

- In distribution  
 $X_n \rightsquigarrow X (\Leftarrow)$   $\forall t$  where  $F_X$  cont:  $\lim_{n \rightarrow \infty} F_{X_n}(t) = F_X(t)$  (Pointwise convergence of distribution function)
- In probability

$$X_n \xrightarrow{\text{prob}} X \Leftrightarrow \forall \epsilon > 0: \lim_{n \rightarrow \infty} P(|X_n - X| > \epsilon) = 0 \quad (\text{uniform conv. of dist. function})$$

- Markov's ineq.:  $P(|X| > \epsilon) \leq \frac{E(|X|)}{\epsilon}$  (Check by Jensen's ineq.)

$$\forall \epsilon > 0: P(X > \epsilon) \leq \frac{E(X)}{\epsilon}$$

$$\Rightarrow P(|X| > \epsilon) = P(|X| > E^+ + \epsilon^-) \leq \frac{E(X^+)}{\epsilon^-}$$

Problem set week 3

$$P(|X - E(X)| > \epsilon) \leq \frac{V(X)}{\epsilon^2}$$

- ① Show that the estimator  $\bar{X}_n$  is consistent for  $\theta$  in  $\otimes \text{Be}(\theta)$  experiment  
 ② Suppose the data  $X$  in a stat. experiment can take values in a countable set  $\mathbb{X}$  (i.e.  $\text{Law}(X)$  is discrete).

In class it was claimed that the data itself is a sufficient statistic (i.e.  $T(X) = X$  is sufficient). Write down an argument for that claim that proves it.

- ③ Suppose  $T$  is a sufficient statistic with values in a set  $\mathbb{T}$  and  $S: \Pi \rightarrow \mathbb{S}$  is a mapping (bijective) with values in  $\mathbb{S}$  (i.e.  $\exists g^n: \mathbb{S} \rightarrow \Pi$  s.t.  $S^{-1}(S(t)) = t \quad \forall t \in \mathbb{T}$ ). Show that the statistic  $S(T(X))$  is sufficient.  
 ④ In class it was claimed that the statistic  $T(X) = (X_1, \bar{X}_n)$  is sufficient for the  $\otimes \text{Be}(\theta)$  experiment. Prove this claim.  
 ⑤ Let  $X_1, \dots, X_n$  be independent and identically distributed (iid) with  $P_{\text{true}}(X_i) \sim P_\theta(\lambda)$ ,  $\lambda \in (0, \infty)$  is unknown. Show that the sample mean  $\bar{X}_n$  is a sufficient statistic for this data.

### Weak Law of Large numbers

$X_1, X_2, \dots, X_n, \text{ IID } X_i, E(X_i) \text{ exists. Then, } \bar{X}_n \xrightarrow{P} E(X_i)$

Proof (In the case where  $V(X_i) < \infty$ )

$$\text{Let } \Omega, P(|\bar{X}_n - E(\bar{X}_n)| \geq \varepsilon) = \frac{V(\bar{X}_n)}{\varepsilon^2} = \frac{1}{n} \frac{V(X_i)}{\varepsilon^2}.$$

Chебышев.  $X_1, X_2, \dots, X_n, \text{ IID } X_i$ .

We also know that  $E(\bar{X}_n) = E(X_i)$  since  $X_1, X_2, \dots, X_n, \text{ IID } X_i$ ,

$$\therefore E(\bar{X}_n) = E\left(\frac{\sum X_i}{n}\right) = \frac{\sum E(X_i)}{n} = \frac{nE(X_i)}{n} = E(X_i).$$

$$\text{So } P(|\bar{X}_n - E(\bar{X}_n)| \geq \varepsilon) = P(|\bar{X}_n - E(\bar{X}_n)| \geq \varepsilon) = \frac{1}{n} \frac{V(X_i)}{\varepsilon^2} \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Thus,  $\bar{X}_n \rightarrow E(X_i)$ .  $\square$

### Central Limit theorem (CLT)

$\exists \epsilon, X_1, X_2, \dots, X_n, \text{ and } E(X_i), V(X_i) \text{ exist, then}$

$$\bar{X}_n = \frac{\sum X_i}{n} \rightsquigarrow X \sim \text{Normal}\left(E(X_i), \frac{V(X_i)}{n}\right).$$

Cor.

$$i) \bar{X}_n - E(X_i) \rightsquigarrow X - E(X_i) \sim \text{Normal}(0, \frac{V(X_i)}{n})$$

$$ii) \sqrt{n}(\bar{X}_n - E(X_i)) \rightsquigarrow \sqrt{n}(X - E(X_i)) \sim \text{Normal}(0, V(X_i))$$

$$iii) Z_n := \frac{\bar{X}_n - E(\bar{X}_n)}{\sqrt{V(\bar{X}_n)}} = \frac{\sqrt{n}(\bar{X}_n - E(X_i))}{\sqrt{V(X_i)}} \rightsquigarrow Z \sim \text{Normal}(0, 1)$$

$$iv) \lim_{n \rightarrow \infty} P\left(\frac{\bar{X}_n - E(\bar{X}_n)}{\sqrt{V(\bar{X}_n)}} \leq z\right) = \lim_{n \rightarrow \infty} P(z_n \leq z) = P(Z \leq z) = \phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{x^2}{2}} dx$$

### Approximation

For "large"  $n$  ( $\approx n \gg 30$ ), we can use

$$P\left(\frac{\bar{X}_n - E(\bar{X}_n)}{V(\bar{X}_n)} \leq z\right) \approx P(Z \leq z) = \phi(z)$$

## Exercises

1. Prove that the statistic  $T(X_1, \dots, X_n) = \frac{\sum_{i=1}^n X_i}{n} = \bar{X}_n$  is consistent for  $\theta$  in  $\hat{\Theta}$  Bernoulli( $\theta$ ) exp.

Sol.

We have to show that  $\bar{X}_n$  is constant, i.e.  $\bar{X}_n \xrightarrow{P} \theta$ .  
 Let  $\varepsilon > 0$ .  $S_n = \sum_{i=1}^n X_i$ ; so  $\bar{X}_n = \frac{S_n}{n}$ .  $P(|\bar{X}_n - \theta| > \varepsilon) = P(|S_n - n\theta| > n\varepsilon) =$

$$= P(|S_n - n\theta| > n\varepsilon). \text{ Since } V(S_n) = n\theta(1-\theta), \text{ by Chebyshev:} \\ P(|\bar{X}_n - \theta| > \varepsilon) = P(|S_n - n\theta| > n\varepsilon) \leq \frac{V(S_n)}{\varepsilon^2 n} = \frac{\theta(1-\theta)}{\varepsilon^2 n} < \frac{\frac{1}{4}}{\varepsilon^2 n} = \frac{1}{4\varepsilon^2 n} \xrightarrow{n \rightarrow \infty} 0 \quad \square$$

$$\left. \begin{aligned} P(|X - E(X)| > \alpha) &= P((X - E(X))^2 > \alpha^2) \stackrel{\text{Markov}}{\leq} \frac{1}{\alpha^2} E((X - E(X))^2) = \frac{V(X)}{\alpha^2} \end{aligned} \right\} \text{Chebyshev}$$

2. Suppose  $X$  takes values in  $\mathbb{X}$  where  $|\mathbb{X}| \leq |\mathbb{N}|$ .

Let  $S(X) = s$ , if  $X = s$ . We want to show that  $P_\theta(X \in B | S(X) = s)$  is independent of  $\theta$ .  $B \subseteq \mathbb{X}$

$$P_\theta(X \in B | S(X) = s) = P_\theta(X \in B | X = s) = \begin{cases} \frac{P_\theta(\{X \in B\} \cap \{X = s\})}{P_\theta(X = s)}, & P_\theta(X = s) > 0 \\ 0 & , P_\theta(X = s) = 0 \end{cases}$$

$$P_\theta(\{X \in B\} \cap \{X = s\}) = \begin{cases} P_\theta(X = s), & s \in B \\ 0, & s \notin B, \text{ so} \end{cases}$$

$$P_\theta(X \in B | S(X) = s) = \begin{cases} \frac{P_\theta(X = s)}{P_\theta(X = s)}, & P_\theta(X = s) > 0, s \in B \\ 0, & \text{otherwise} \end{cases} = \begin{cases} 1 & \text{if } s \in B \\ 0 & \text{if } s \notin B \end{cases}$$

Alt:

Assume  $T(X) = X$ , let  $t = x$ . For any event  $B \in \mathcal{F}_X$   
 $P_\theta(X \in B | T(X) = t) = P_\theta(X \in B | X = x) = \mathbf{1}_B(x) = \begin{cases} 1, & x \in B \\ 0, & x \notin B \end{cases}$

3. Suppose  $T$  sufficient  $T: \mathbb{X} \rightarrow \bar{\mathbb{T}}$  and  $S: \bar{\mathbb{T}} \rightarrow \mathcal{S}$  inf.  
Show that  $S$  is sufficient.

sol.  $\forall s \in \mathcal{S} := \{s : S(t) = s\}$  and  $\forall B \in \mathcal{F}_{\bar{\mathbb{T}}}$  the cond. prob.

$$\frac{P_{\theta}(\{x \in B \mid S(T(X)) = s\})}{P_{\theta}(S(T(X)) = s)} = \frac{P_{\theta}(\{x \in B \mid T(X) = t\})}{P_{\theta}(T(X) = t)}$$

$= P_{\theta}(X \in B \mid T(X) = t)$  which is indep. of  $\theta$  because of  $T$  suff.  $\square$

Prop.

In the  $\bigotimes_{i=1}^n \text{Be}(\theta)$  experiment, the estimator  $\bar{X}_n$  is uniformly best among unbiased estimators for quadratic loss (i.e., for any unbiased estimator  $T$ :  $R(\bar{X}_n, \theta) = E_{\theta}((\bar{X}_n - \theta)^2) \leq E_{\theta}((T(X) - \theta)^2) = R(T, \theta)$  for all  $\theta \in [0, 1]$ ).

Before we prove this, let's revisit the notion of cond. expectation:

Let  $Y$  be a R.V. with finite support in  $\mathbb{Y}$  i.e.  $|\mathbb{Y}| < \infty$ .

Let  $U$  be a stochastic variable with values in  $\mathbb{U}$  and  $u \in \mathbb{U}$ ,  $h$  a real-valued fct. of  $Y$ .

Def

The cond. expectation of  $h(Y)$  given  $U=u$ , written

$E(h(Y) \mid U(Y)=u)$  is defined as the expectation of  $h(Y)$  under the cond. distribution of  $Y$  given  $U(Y)=u$ :  $P(Y=y \mid U(Y)=u) = \frac{P(Y=y, U(Y)=u)}{P(U(Y)=u)}$   
if  $P(U(Y)=u) \neq 0$  (0 otherwise)

$$E(h(Y) \mid U(Y)=u) = \sum_{y \in \mathbb{Y}} h(y) P(Y=y \mid U(Y)=u).$$

In special case  $h(y) = 1_B(y)$  for some  $B \subseteq \mathbb{Y}$ , then  
 $E(1_B(y) \mid U(Y)=u) = \sum_{y \in B} P(Y=y \mid U(Y)=u) = P(Y \in B \mid U(Y)=u)$ .

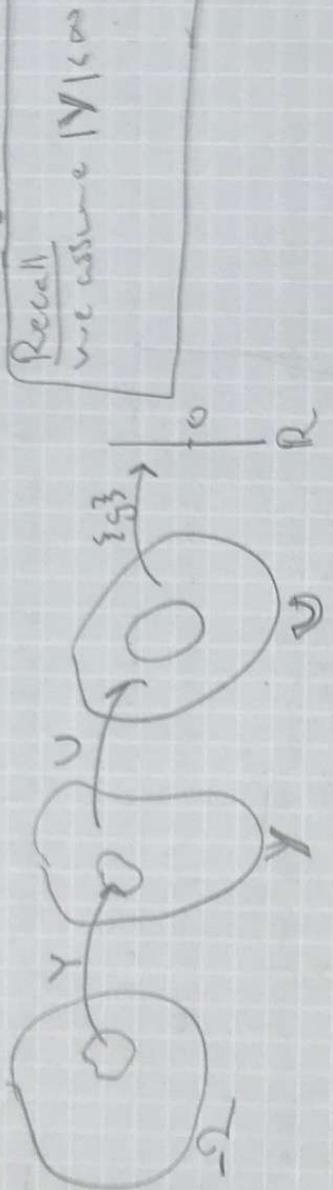
① Def

The conditional expectation can be a R.V.

$$E(h(Y)|U)$$

Note the following properties of cond. exp. as a R.V.:

Let  $\mathcal{M}_U$  be the set of all real-valued R.V.s  $Z$  that are functions of  $U$ , i.e.  $Z = g(U)$  for some function  $g: U \rightarrow \mathbb{R}$ .



$$\text{i)} E(E(h(Y)|U)) = E(h(Y))$$

proof

exercise 0

ii) For any  $z \in \mathcal{M}_U$ ,  $E(z h(Y)|U) = z E(h(Y)|U)$  ("almost surely")  
Intuition

R.V.s  $X, Y$  equal w.p. 1 means  $1 = P(X=Y)$

iii) For any function  $h: Y \rightarrow \mathbb{R}$  s.t.  $h = h_1 + h_2$ , then  $E(h(Y)|U) = E(h_1(Y)|U) + E(h_2(Y)|U)$  w.p. 1

iv) For any  $z \in \mathcal{M}_U$ ,  $E(z \cdot h(Y)) = E(z \cdot E(h(Y)|U))$

proof

exercise 1

v) For any  $z \in \mathcal{M}_U$ ,  $E((h(Y)-z)^2) = E((E(h(Y)|U)-z)^2) + E(h(Y)-E(h(Y)|U))^2$   
proof  
exercise 2 (assuming iii) and iv)  
hint: add and subtract  $E(h(Y)|U)$  in  
and  $E(E(h(Y)|U)) = E(h(Y)|U)E(1|U) = E_h(Y)|U)$

### Remark

The conditional expectation can be seen as an operator.  
Let  $\mathcal{M}_Y$  be all R-valued R.V.s which are functions of  $Y$ .

Then, since  $U$  is  $U(Y)$  we have  $\mathcal{M}_U \subseteq \mathcal{M}_Y$

and both  $\mathcal{M}_U$  and  $\mathcal{M}_Y$  are vector spaces.

$$\begin{array}{ccc} \Omega & \xrightarrow{Y} & \mathbb{R} \\ \mathcal{M}_Y = \{f(Y)\} & \xrightarrow{U(Y)} & \mathcal{M}_U \\ & \Omega & \mathbb{R} \end{array}$$

Let  $H \in \mathcal{M}_Y$ , then the cond. expect.  $E(H|U) \in \mathcal{M}_U$  and we define the operator  $\Pi(H) := E(H|U)$ , i.e.  $\Pi: \mathcal{M}_Y \rightarrow \mathcal{M}_U$

The " $\leq$ " in prop. V tells us that  $\underbrace{E(H - \Pi(H))^2}_{E[(h(X) - E(h(Y)|U))^2]} = \min_{z \in \mathcal{M}_U} E(H - z)^2$

In English: the cond. expect. of  $H$  is the element of  $\mathcal{M}_U$  which is closest to  $H$ . You can see this as the "projection"  $\Pi$  of  $H$  onto the space  $\mathcal{M}_U$ .

The " $=$ " in prop. V is the orthogonal decomposition

$$E(H - z)^2 = E((\Pi(H) - z)^2) + E((H - \Pi(H))^2)$$

We are now ready to prove that  $\bar{X}_n$  is uniformly best among all unbiased estimators of  $\Theta$  in  $\bigcap_{\Omega} \text{Be}(\Theta)$  exp.:

### Proof

N.T.S. (need to show) that the risk of  $\bar{X}_n$  is the lowest.

i.e.  $R(\bar{X}_n, \theta) := E_\theta((\bar{X}_n - \theta)^2) \leq E_\theta((T(X_1, \dots, X_n) - \theta)^2) =: R(T, \theta)$   
for all  $\theta \in \Theta$ , i.e. and all unbiased estimators  $T$ .

Define a.R.V.:  $g(\bar{X}_n) = E_\theta(T | \bar{X}_n) = E(T | \bar{X}_n)$  and regard it as  
an estimator of  $\theta$   
 $\bar{X}_n$  is suff. for  $\theta$

by prop. (i),  $E(g(\bar{X}_n)) = E(E(T | \bar{X}_n)) = E(T) = \hat{E}_T$  unbiased  
so,  $g$  is unbiased.

Before, we showed that  $g(\bar{X}_n) = \bar{X}_n$  is the only unbiased  
estimator of  $\theta$  ( $\bar{X}_n$  is complete)

Hence,  $g(\bar{X}_n) = \bar{X}_n$

In prop. (v), set  $Z = \Theta$ , (point-mass R.V.  $P(Z=z) = \begin{cases} 1, & z=0 \\ 0, & z \neq 0 \end{cases}$ ), so  $Z \in \mathcal{Y}$ ,  
 $h(Y) = T$ ,  $U = \bar{X}_n$ ,  $E(h(Y) | U) = g(\bar{X}_n)$   
to get  $E_\theta(T - \theta)^2 \geq E_\theta(g(\bar{X}_n) - \theta)^2$

so, the estimator  $g(\bar{X}_n) = \bar{X}_n$  is at least as good as any  $T$ .  $\square$

For unbiased estimators, the quadratic risk is also the estimators variance:

$$E_\theta((\bar{T} - \theta)^2) = V(\bar{T}) := E_\theta((\bar{T} - E_\theta(\bar{T}))^2$$

$\therefore \bar{X}_n$  for  $\sum_{i=1}^n Be_i(\theta)$  exp is also called a

uniform minimum variance unbiased estimator (UMVUE)

### Bayes estimators

In the Bayesian approach to estimation of the param.  $\theta$  underpinning the law of data  $X = (X_1, \dots, X_n)$  in an exp., one assumes that prior distribution is given over the parameter space  $\Theta \ni \theta$

ex) In the  $\bigotimes_{i=1}^n \text{Be}(\theta)$  exp., assume that the prior distribution is given in the form of density on  $\Theta = [0, 1]$ .

### Def (Integrated risk)

For an estimator  $T$  of  $\theta$  the prior  $g(\theta)$  can be used to reduce the risk function  $R(T, \theta)$  for each  $\theta \in [0, 1]$  to a single number,

$$\text{by integration: } B(R(T)) = B_r(T) := \int_0^1 R(T, \theta) g(\theta) d\theta = \int_0^1 E_\theta(T - \theta)^2 g(\theta) d\theta$$

$B_r(T)$  is called integrated risk or mixed risk under risk  $R$ .

### Def

A Bayes estimator  $T_B$  of  $\theta$  is the estimator that minimizes the integrated risk, i.e.  $T_B := \arg \min_T B_r(T)$  and  $B_r(T_B)$ , the minimal integrated risk, is called Bayes Risk

The more "Bayesian" comes from Bayes formula:  
 $\{D_i\}_{i=1}^n$  partition  $\Omega$ , then  $P(A) = \sum_{i=1}^n P(A|D_i)P(D_i)$

Motivation: Consider the case when  $P$  is the joint distribution of  $(X, U)$  where  $X$  is data and  $U$  is a R.V. that takes k possible values in  $\Theta := \{\theta_1, \dots, \theta_k\}$ . Then  $A = "X \in A"$  and  $D := "U = \theta_i"$ . Posterior prob of  $A$  given data  $\Theta$ :  $P(U=\theta_i | X \in A) = \frac{P(X \in A | U=\theta_i)P(U=\theta_i)}{\sum_{j=1}^k P(X \in A | U=\theta_j)P(U=\theta_j)}$

The Bayesian approach views  $\{P_\theta : \theta \in \Theta\}$  as a family of conditional distributions given  $\theta$ , without a prior dist.

Ex In  $\mathbb{R}^+$   $\text{Be}(\theta)$  exp, consider the family of prior densities for  $\theta \in [0, 1]$ :

$$\text{for each } (\alpha, \beta) \in (0, \infty)^2 =: \mathbb{R}_{>0}^2 \quad (\alpha > 0, \beta > 0)$$

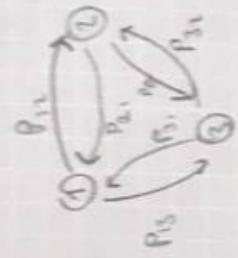
$$g_{\alpha, \beta}(\theta) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)}, \quad \theta \in [0, 1] \quad \text{where } B \text{ is beta function:}$$

$$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}, \quad \text{and} \quad \Gamma(\alpha) := \int_0^\infty x^{\alpha-1} e^{-x} dx.$$

Consider the whole family  $g_{\alpha, \beta}$  for specification of prior density.

## Marginal Classes

If they "line", then it applies



$$P = \begin{pmatrix} 0 & p_{12} & p_{13} \\ p_{11} & 0 & p_{13} \\ 0 & p_{12} & 0 \end{pmatrix} \quad Q = \begin{pmatrix} -q_{11} & q_{12} & q_{13} \\ q_{21} & -q_{22} & q_{23} \\ q_{31} & q_{32} & -q_{33} \end{pmatrix}$$

$$q_{ij} = q_{ik} p_{kj} \text{ if } i+j$$

$$P(t) = e^{Qt}$$

Note Bayes cont.

We consider the whole family  $\mathcal{G}_\theta$ , induced by  $(\alpha, \beta) \in \mathbb{R}_{>0}^2$  to allow a wide range of prior beliefs about  $\theta$ .

It will become clear that the Bayesian approaches are very useful even to prove non-Bayesian prop. of estimators.

Prop. In  $\bigcup_{\theta \in \Theta} \mathcal{G}_\theta$  case, let  $\mathcal{J}$  be an arbitrary prior density for  $\theta \in [\theta_0, 1] = \mathbb{B}$ , then the Bayes estimator is

$$\bar{\theta}_B(x) : \mathbb{X} \rightarrow \mathbb{B} = \frac{\int_0^1 s(x)^{(1-\theta)} (1-\theta)^{\theta-1} g(\theta) d\theta}{\int_0^1 s(x)^{(1-\theta)} (1-\theta)^{\theta-1} g(\theta) d\theta} \quad \text{where } s(x) = \sum_{i=1}^n x_i$$

Note

$T_B$  is a function of the sufficient statistic  $X_n = \frac{s(x)}{n}$

Proof exercise

### Def

An estimator  $T$  of a parameter  $\theta$  is called admissible if, for every estimator  $S$  of  $\theta$ , the relation given by

$$\Theta : R(s, \theta) \leq R(T, \theta) \quad \forall \theta \in \Theta \text{ implies}$$

$$\Theta : R(s, \theta) = R(T, \theta)$$

Thus, admissibility of  $T$  means that there can be no estimator  $S$  which is uniformly at least as good [ $\theta$  holds] and strictly better for one  $\theta_0 \in \Theta$ , because then  $R(s, \theta_0) < R(T, \theta_0)$  and thus ( $\star\star$ ) is contradicted.

So, non-admissibility of  $T$  means that  $T$  can be improved by another estimator  $S$ .

### Prop

Suppose that in  $\Theta \subset \text{Be}(\theta)$  exp, the prior density  $g$  is such that  $g(\theta) > 0$  for each  $\theta \in [0, 1]$  with the exception of a finite number of points. Then the Bayes estimator  $T_g$  for this prior density is admissible for quadratic loss.

Proof exercise

When trying these 2 exercises, show them to Rueze.

## Dg

### Minimax Estimators

Let the maximal risk of an estimator  $T$  be  $M(T) := \max_{\theta \in \Theta} R(T, \theta)$

An estimator  $T_\mu$  is called Minimax if

$$M(T_\mu) = \min_T \max_{\theta \in \Theta} R(T, \theta)$$

some collection  
of estimators

### Prop

In  $\text{Ber}(\theta)$  exp., the Bayes Estimator  $T_{\alpha, \beta}$  for  $\alpha = \beta = \frac{\sqrt{2}}{2}$  is a minimax estimator

### Proof

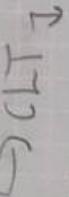
### Outline of course:

#### ① Point Estimators

- Moments Estimators { solve example moments = theoretical moments likelihood

- Maximum Likelihood estimators {  $L(\theta) \propto P(x_1, \dots, x_n | \theta)$

$$\hat{\theta} = \arg \max_{\theta \in \Theta} L(\theta)$$



#### ② Set Estimators

point est.

$\hat{\theta}_n$  set, VERY high  $\hat{\theta}_n$  prob. to contain  $\theta_0$  (true law)

#### ③ Hypothesis testing

#### ④ Regression: $f_{Y|X}$

## Parametric Inference

Recall the D.T. setting!

Given a param. exp.  $(\Omega, \mathcal{F}_n, P)$ ,  $\mathcal{P} = \{\mathbb{P}_{\theta}; \theta \in \Theta, |\Theta| < \infty\}$   
we are interested in some function  $g(\theta)$  based on data  $x$ .  
Typically,  $x = (x_1, \dots, x_n)$ ,  $x_1, \dots, x_n \stackrel{i.i.d.}{\sim} X$  has law  $\mathbb{P}_{\theta}$  is from a product  
exp. element,  $x_1, \dots, x_n \sim F_{\theta} = F_{\theta}(x_i; \theta)$  (or  $\mathbb{P}_{\theta} = f_{\theta}(x_i; \theta)$ )  
e.g.  $x_1, \dots, x_n \sim \underset{i.i.d.}{\stackrel{n}{\sim}} \text{Be}(\theta)$

Ex1

$x_1, \dots, x_n \sim N(\mu, \sigma^2)$ . Then  $\Theta = \{\theta = (\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0\}$

Suppose  $X_i$  is the outcome of a blood-alcohol breath-test and we

are interested in the fraction of testers with score  $> 1$ .

$$\text{What is } g(\theta) ? : g(\theta) = P(X_i > 1) = 1 - P\left(\frac{x_i - \mu}{\sigma} < \frac{1 - \mu}{\sigma}\right) = \\ := \bar{Z} \sim N(0, 1) \\ = 1 - \Phi\left(\frac{1 - \mu}{\sigma}\right)$$

Ex2  $X \sim \Gamma(\alpha, \beta)$ ,  $\Theta = \{\theta = (\alpha, \beta) = (\alpha_1, \alpha_2) \in \mathbb{R}_{>0}^2$ , recall  $f(x; \alpha, \beta) = \frac{1}{\beta^{\alpha_1} \Gamma(\alpha)} x^{\alpha_1 - 1} e^{-x/\beta}$ .

Often  $X$  is used to model lifetime of items.

Then we are interested in an estimate of the mean lifetime:  $\hat{\theta} = \frac{\sum x_i}{n}$

$$g(\theta) = g(\alpha, \beta) = E(X; \alpha, \beta) = \alpha \beta$$

## Method of Moments Estimator (MOME)

Suppose the parameter  $\theta$  have components  $(\theta_1, \theta_2, \dots, \theta_k)$ , then

For i.e. jth, define the j<sup>th</sup> moment:

$$\alpha_j := \alpha_j(\theta) = E(X^j; \theta) = \int x^j dF(x)$$

and the j<sup>th</sup> sample moment:

$$\hat{\alpha}_j := \frac{1}{n} \sum_{i=1}^n X_i^j$$

The MOME  $\hat{\theta}_m(X) = \hat{\theta}_m$  is defined to be the value of  $\theta$  s.t.

$$\begin{cases} \alpha_1(\hat{\theta}_m) = \hat{\alpha}_1 \\ \alpha_2(\hat{\theta}_m) = \hat{\alpha}_2 \\ \vdots \\ \alpha_k(\hat{\theta}_m) = \hat{\alpha}_k \end{cases} \quad \text{i.e. } \hat{\theta}_m \text{ is the MOM estimate and is the solution to} \\ \text{to this system of k equations in k unknowns.}$$

Note

MOME, if they exist, are typically suboptimal but it is often easy to compute and can be used to initialize the iterative strategies used to obtain more optimal estimators (e.g. MLE)

ex 3  $X_1, \dots, X_n \sim \text{Be}(\theta)$ . Find MOME for  $\theta$ .

$$\begin{aligned} \text{sol. } \alpha_1(\hat{\theta}) &= E(X; \theta) = \hat{\alpha}_1 \\ \hat{\alpha}_1 &:= \frac{1}{n} \sum_{i=1}^n x_i = \bar{X}_n \end{aligned} \quad \left. \begin{array}{l} \text{By equating we solve for } \hat{\theta}_m \\ \hat{\theta}_m = \bar{X}_n \end{array} \right\}$$

ex 4  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ . Find mome.

$$\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$$

$$\begin{aligned} \text{Sol. } \alpha_1 &= E(X_1; \mu, \sigma^2) = \mu \\ \alpha_2 &= E(\hat{X}^2; \mu, \sigma^2) = V(X_1; \mu, \sigma^2) + [E(X_1; \mu, \sigma^2)]^2 = \sigma^2 + \mu^2 \end{aligned}$$

$$\begin{aligned} \text{Solve } \left\{ \begin{array}{l} \hat{\mu}_m = \bar{X}_n \\ \hat{\sigma}_m^2 + \hat{\mu}_m^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 \Rightarrow \sigma_m^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \end{array} \right. \end{aligned}$$

ex5

Show that MOMC for parameter  $\lambda$  in a product  $P_\theta(x)$  exp

$$\text{is } \bar{X}_n = \frac{1}{n} \sum_{i=1}^n x_i. \quad \text{Recall } f(x; \lambda) = \prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} \frac{e^{-\lambda}}{\lambda^n}$$

Exercise: homogeneous densities are possible (e.g.  $\lambda$ )

### Properties of MOMC

Under "appropriate conditions" on the experiment, the following hold:

1)  $\hat{\theta}_n$  exists with prob  $\rightarrow 1$  (equations are solvable for large  $n$ )

2)  $\hat{\theta}_n \xrightarrow{P} \theta$  (asymptotic) consistency.

3)  $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D} N(0, \Sigma)$ ,  $\Sigma = g E_\theta(Y Y^\top) g^\top$  where  $Y = [x, x^\top, \dots, x^\top]^\top$ ,  
 $g = (g_1, \dots, g_n)$ ,  $g_i = \frac{\partial}{\partial \theta} \alpha_i(\theta)$

### Maximum likelihood estimators (MLE)

$X_1, \dots, X_n \stackrel{iid}{\sim} X$ , and  $X_i \sim F(x_i; \theta)$  with PMF or PDF  $f(x_i; \theta) = \underbrace{f(x_i; \theta)}_{f_\theta(x_i) = \partial F_\theta(x_i)}$

Def

The likelihood function is defined by: if product

$$L_n(\theta) := L_n(\theta; x_1, \dots, x_n) = f(x_1, \dots, x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$$

Note

$$\frac{L_n(\theta)}{L_n(\theta_0)}: \Theta \rightarrow [0, \infty) \equiv R_{>0} \neq \emptyset \text{ in large }$$

Def

$$\log\text{-likelihood function is } \lambda_n(\theta) := \lambda_n(\theta; x_1, \dots, x_n) = \ln(L_n(\theta))$$

Thus, the likelihood function is just the joint density of the data, except that we view this as a function of the parameter  $\theta$  random (because  $x = (x_1, \dots, x_n)$  is random)

Qof

The MLE  $\hat{\theta}_n(x_1, \dots, x_n) = \hat{\theta}_n = \underset{\theta \in \Theta}{\operatorname{argmax}} L_n(\theta) = \underset{\theta \in \Theta}{\operatorname{argmax}} \lambda_n(\theta)$

ex1

Let  $x_1, \dots, x_n \sim \text{Beta}(\theta)$ . Find MLE of  $\theta$ .

$$\text{sol } L_n(\theta) = \prod_{i=1}^n \frac{1}{\Gamma(n)} \theta^{x_i} (1-\theta)^{n-x_i} = \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} = \theta^{\sum x_i} (1-\theta)^{n-\sum x_i}$$

$$\Rightarrow \lambda_n(\theta) = \lambda_n\left(\theta^{\sum x_i} (1-\theta)^{n-\sum x_i}\right) = s \lambda_n(\theta) + (n-s) \lambda_n(1-\theta)$$

$$\frac{\partial}{\partial \theta} \lambda_n(\theta) = \frac{1}{\theta} s - (n-s) \frac{1}{1-\theta} = 0 \Rightarrow \frac{1}{\theta} s = (n-s) \frac{1}{1-\theta} \Rightarrow (1-\theta)s = (n-s)\theta \Rightarrow$$

$$\Rightarrow s - \theta s = n\theta - s\theta \Rightarrow \theta = \frac{s}{n} = \bar{X}_n$$

$$\frac{\partial^2}{\partial \theta^2} \lambda_n(\theta) = -\frac{s}{\theta^2} + \frac{n-s}{(1-\theta)^2} \Rightarrow \frac{\partial^2 \lambda_n}{\partial \theta^2}(\frac{s}{n}) = -\frac{s}{(\bar{X}_n)^2} + \frac{n-s}{(1-\bar{X}_n)^2} = -<0> \Rightarrow \text{max!}$$

ex2 Show that MLE  $\hat{\lambda}_n$  of parameter  $\lambda$  in a product  $\text{Exp}(\lambda)^n$  up is  $\frac{1}{\bar{x}_n}$ .  
Recall  $X_1, \dots, X_n \sim \text{Exp}(\lambda)$  means  $f(x_i; \lambda) = \lambda e^{-\lambda x_i} \mathbf{1}_{(0, \infty)}(x_i)$  for  $\lambda \in (0, \infty)$

ex3  $X_1, \dots, X_n$  iid  $N(\mu, \sigma^2)$ . Find the MLE  $\hat{\theta}_n$  for the unknown parameter  $\theta$ :  $\theta = (\mu, \sigma^2)$

Sol Likelihood function (ignoring constants  
as they are irrel.)

(for convenience)

$$\lambda_n(\theta) = L_n(\mu, \sigma) = \prod_{i=1}^n \frac{1}{\sigma} e^{-\frac{(X_i - \mu)^2}{2\sigma^2}} = \sigma^{-n} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2} = \sigma^{-n} e^{-\frac{n S_n^2}{2\sigma^2} - \frac{n(X_n - \mu)^2}{2\sigma^2}}$$

$$\text{where } \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i, S_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

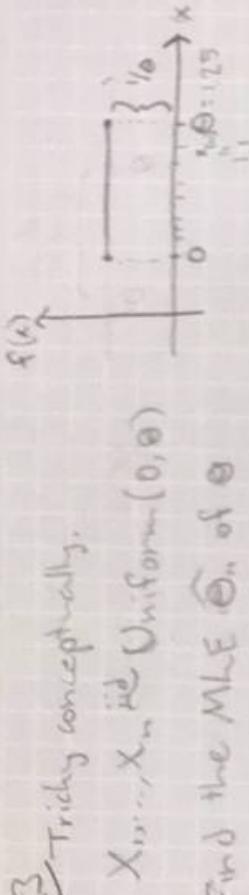
$$\text{The log-likelihood } \lambda_n(\mu, \sigma) = -n \ln(\sigma) - \frac{n S_n^2}{2\sigma^2} - \frac{n(X_n - \mu)^2}{2\sigma^2}$$

Now, solve the equations:  $\left( \frac{\partial}{\partial \mu} \lambda_n(\mu, \sigma) = 0 \right)$

$$\left( \frac{\partial}{\partial \sigma} \lambda_n(\mu, \sigma) = 0 \right)$$

exercise: show that  $\hat{\theta}_n = \left( \hat{\theta}_{n,1}, \hat{\theta}_{n,2} \right) = \left( \bar{X}_n, S_n \right)$

Ex3 Tricky conceptually.



$$\Omega(x_i; \theta) = \frac{1}{\theta} \mathbb{1}_{[0, \theta]}(x_i) = \begin{cases} \frac{1}{\theta}, & x_i \in [0, \theta] \\ 0, & x_i \notin [0, \theta] \end{cases}$$

Likelihood: consider  $\theta$  fixed, suppose  $0 < x_i$  for some  $i \in \{1, \dots, n\}$ .

The  $f(x_i; \theta) = 0$  and thus  $L_n(\theta) = 0$  if any  $x_i > \theta$   
 so,  $L_n(\theta) = 0$  if  $\max(\{x_1, \dots, x_n\}) = x_{(n)} > \theta$ .

Now, consider any  $0 < x_{(n)} \leq 0$ . Then for each  $x_i$  we have

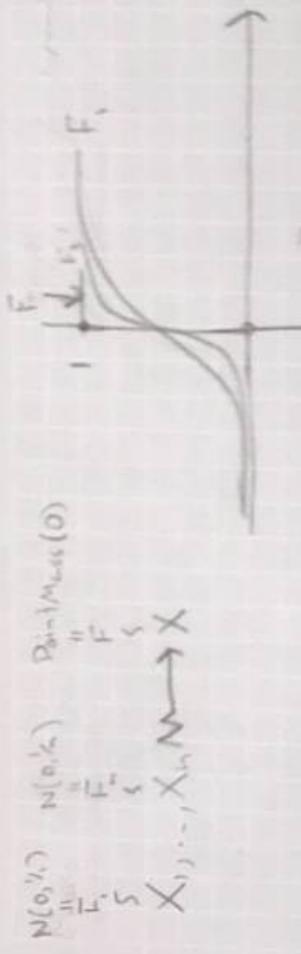
$$\begin{aligned} f(x_i; \theta) &= \frac{1}{\theta} \Rightarrow L_n(\theta) = \prod_{i=1}^n f(x_i; \theta) = \frac{1}{\theta^n} = \theta^{-n} \\ \Rightarrow L_n(\theta) &= \begin{cases} 0, & \text{if } \theta < x_{(n)} \\ 0, & \text{if } \theta \geq x_{(n)} \end{cases} \end{aligned}$$

So, MLE is  $x_{(n)}$  (ith order statistic)

### Properties of MLEs

- 1) The MLE is consistent:  $\hat{\theta}_n \xrightarrow{P} \theta_0$ , where  $\theta_0$  denote the true value of parameter.  
 If the model is "good" or "reasonable".
- 2) The MLE is equivariant: if  $\hat{\theta}_n$  is the MLE of  $\theta$ , then  $g(\hat{\theta}_n)$  is the MLE of  $g(\theta)$  (the function of interest)
- 3) The MLE is asymptotically normal:  $(\hat{\theta}_n - \theta_0) \xrightarrow{\sqrt{n}} N(0, 1)$   
 $\text{se}_n(\hat{\theta}_n) = \sqrt{\text{Var}(\hat{\theta}_n)}$  is the standard error
- 4) The MLE is asymptotically optimal (efficient):  
 Amongst all well-behaved estimators, the MLE has the smallest variance  
 (at least for large samples)
- 5) The MLE is approximately the Bayes Estimators

\* These things hold when certain regularity conditions are satisfied:  
 Mainly smoothness conditions on the density  $f(x; \theta)$   
 WHITENESS



Distr. conv. in distribution (P. 152, 154 in (SE))

Let  $X_1, \dots, X_n$  have DFs  $F_1, \dots, F_n$ . Let  $X$  be another RV with DF  $F$ .

Then we say  $X_n \xrightarrow{D} X$  if  $\forall \epsilon \in \mathbb{R}$  at which  $F$  is continuous, we have

$$\lim_{n \rightarrow \infty} F_n(t) = F(t) \quad (\text{Pointwise convergence of } DF_s)$$

$$\text{Result: } \lim_{n \rightarrow \infty} P(\{\omega : X_n(\omega) \leq t\}) = P(\{\omega : X(\omega) \leq t\})$$

Prop. (LT)

Let  $X_1, \dots, X_n$  &  $X$ , and  $E(X_i), V(X_i)$  exist.

$$\text{Then: } \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{D} X \sim N(E(X_i), \frac{V(X_i)}{n})$$

$$\therefore \bar{X}_n - E(X_i) \xrightarrow{D} X - E(X) \sim N(0, \frac{V(X_i)}{n})$$

$$\therefore \sqrt{n} (\bar{X}_n - E(X_i)) \xrightarrow{D} N(X - E(X), V(X_i))$$

$$\therefore Z_n := \frac{\sqrt{n}}{\sqrt{V(X_i)}} (\bar{X}_n - E(X)) = \frac{\bar{X}_n - E(X)}{\sqrt{V(X_i)}} \xrightarrow{D} Z \sim N(0, 1)$$

$$\therefore \lim_{n \rightarrow \infty} P\left(\frac{\bar{X}_n - E(\bar{X}_n)}{\sqrt{V(\bar{X}_n)}} \leq z\right) = \lim_{n \rightarrow \infty} P(Z_n \leq z) = \Phi(z) := \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$$

ex Suppose collection  $X_1, \dots, X_n$  model # errors in program runed

i.e. Suppose  $X_i \sim \text{Po}(\lambda=5)$  and  $X_i \text{ iid}$ ,  $E(X_i) = 5$

Suppose n=125 and we want to make prob. statement of  $\bar{X}_n$   
 $P(\bar{X}_n \leq 5.5) = P\left(\frac{\sum(X_i - E(X_i))}{\sqrt{n\lambda}} \leq \frac{\sum(5.5 - 5)}{\sqrt{n\lambda}}\right) = P(Z \leq \frac{\sqrt{n}(5.5 - 5)}{\sqrt{\lambda}}) =$   
 $= \Phi(2.5)$

Prob. back-fall

Def  
Conv. in prob

$X_n \xrightarrow{P} X$  if  $\forall \varepsilon > 0$ ,  $P(|X_n - X| > \varepsilon) \xrightarrow{n \rightarrow \infty} 0$  :  $(X_n \xrightarrow{P} c \text{ if } c \text{ is PointMass}(c))$

Def  
Conv. in dist.

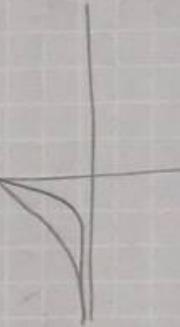
$X_n \rightsquigarrow X$  if  $F(t)$  continuous :  $\lim_{n \rightarrow \infty} F_n(t) = F(t)$

Def

conv. in quadratic mean (q.m.) or  $L_2$  conv. :

$X_n \xrightarrow{q.m.} X$  if  $E((X_n - X)^2) \xrightarrow{n \rightarrow \infty} 0$

$\xrightarrow{F_n \rightarrow F}$



$\Rightarrow X_n \sim N(0, 1/n)$ ,  $X \sim \text{PointMass}(0)$ .

How to formalize intuition " $X_n \rightarrow 0$ "?

conv. in dist.:

$\sqrt{n} X_n \sim N(0, 1) =: Z$

for  $t < 0$ :  $F_n(t) = P(X_n < t) = P(\sqrt{n} X_n < \sqrt{n}t) = P(Z < \sqrt{n}t) \xrightarrow{n \rightarrow \infty} 0$

for  $t > 0$ :  $F_n(t) = P(X_n < t) = P(\sqrt{n} X_n < \sqrt{n}t) = P(Z < \sqrt{n}t) \xrightarrow{n \rightarrow \infty} 1$

Hence,  $F_n(t) \rightarrow F(t) \quad \forall t \neq 0$  But this is ok since we only needs to be for  $t$  where  $F(t)$  cont.

Conv. in prob.:  $\forall \varepsilon > 0$ ,  $P(|X_n - 0| > \varepsilon) = P(|X_n|^2 > \varepsilon^2) \leq \frac{E(X_n^2)}{\varepsilon^2} = \frac{1}{n\varepsilon^2} \xrightarrow{n \rightarrow \infty} 0$

## Prop. (PNN)

a)  $X_n \xrightarrow{q} X \Rightarrow X_n \xrightarrow{p} X$

b)  $X_n \xrightarrow{p} X \Rightarrow X_n \rightsquigarrow X$

c) If  $X_n \rightsquigarrow X$  and  $\exists c \in \mathbb{R}: P(X_n = c) = 1$ , then  $X_n \xrightarrow{p} X$

So, conv. in prob.  $\Rightarrow$  conv. in dist.

ex conv. in prob  $\neq$  conv. in dist.  
 $\downarrow$   
 if  $X \sim \text{PointMass}(c)$

Let  $U \sim U(0, 1)$ ,  $X_n = \sqrt{n} \mathbf{1}_{(0, \frac{1}{n})}(U)$ ,  $X \sim \text{PointMass}(0) = 0$

Then  $P(|X_n| > \varepsilon) = P\left(\sqrt{n} \mathbf{1}_{(0, \frac{1}{n})}(U) > \varepsilon\right) = P\left(0 < U < \frac{\varepsilon}{\sqrt{n}}\right) = \frac{\varepsilon}{\sqrt{n}} \rightarrow 0$  so  $X_n \xrightarrow{p} X$

But  $E(X_n^2) = \int_0^{1/n} (n u)^2 \cdot n du = n \int_0^{1/n} u^2 du = n \left(\frac{1}{n} - 0\right) = 1$  for all  $n$  so  $X_n \not\rightsquigarrow X$

ex conv. in dist.  $\neq$  conv. in prob.

Let  $X \sim N(0, 1)$ ,  $X_n = -X$  for  $n = 1, 2, 3, \dots$

Hence  $X_n \sim N(0, 1)$  (goes "backwards"  $\rightsquigarrow \rightsquigarrow$  but since  $N(0, 1)$  symmetric, still the same)

$X_n$  has same DF as  $X$  for all  $n$ , so trivially  $\lim_{n \rightarrow \infty} F_n(x) = F(x) \forall x \Rightarrow X_n \rightsquigarrow X$

But  $P(|X_n - X| > \varepsilon) = P(|2X| > \varepsilon) = P(|X| > \frac{\varepsilon}{2}) \neq 0$

we may conjecture that if  $X_n \xrightarrow{p} c$ , then  $E(X_n) \rightarrow c$ . Not true in general.

(unless  $X_n$  is uniformly integrable)

Let  $X_n$  be s.t.  $P(X_n = n) = \frac{1}{n}$  and  $P(X_n = 0) = 1 - \frac{1}{n}$

Now,  $P(|X_n - 0| < \varepsilon) = P(X_n = 0) = 1 - \frac{1}{n} \xrightarrow{n \rightarrow \infty} 1 \Rightarrow X_n \xrightarrow{p} 0$

but  $E(X_n) = n \cdot \frac{1}{n} + 0 \cdot \left(1 - \frac{1}{n}\right) = n \Rightarrow E(X_n) \rightarrow \infty$

## Back to CLT

$$Z_n := \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{D} N(0, 1) \quad \text{for } X_1, \dots, X_n \stackrel{iid}{\sim} X_1,$$

$\nu := E(X_1) < \infty, \nu(X_1) := \sigma^2 < \infty$

We typically don't know  $\sigma$ .  
If we replace  $\sigma$  by its estimate  $s_n := \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2}$  it still works.

The Berry-Esseen inequality (tells us about accuracy of Normal approx.)

Suppose further that  $E(|X_1|^3) < \infty$ .

$$\text{Then } \sup_z |P(Z_n < z) - \Phi(z)| \leq \frac{33}{4} \cdot \frac{E(|X_1 - \mu|^3)}{\sqrt{n} \sigma^2}$$

## Back to properties of MLE

### Prop (consistency)

Let  $\theta_0$  denote the true (and possibly unknown) value of  $\theta$ .

$$M_n(\theta) := \frac{1}{n} \sum_{i=1}^n \ln \left( \frac{f(x_i; \theta)}{f(x_i; \theta_0)} \right), \quad M(\theta) = E_{\theta_0} \left( \ln \left( \frac{f(x_i; \theta)}{f(x_i; \theta_0)} \right) \right)$$

(maximizing (=) maximizing  $\lambda_{\theta_0}(\theta)$ )

Suppose that  $\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{P} 0$

and that  $\forall \varepsilon > 0 : \sup_{\theta: |\theta - \theta_0| \geq \varepsilon} M(\theta) < M(\theta_0)$

Let  $\hat{\theta}_n$  denote the MLE. Then  $\hat{\theta}_n \xrightarrow{P} \theta_0$ .

### Prop (equivariance)

Let  $\tau = g(\theta)$  be a function of  $\theta$  and  $\hat{\theta}_n$  be MLE of  $\theta$

Then  $\hat{\tau}_n = g(\hat{\theta}_n)$  is the MLE of  $\tau$

### Proof

Let  $h = g^{-1}$  denote the inverse of  $g$ . Then  $\hat{\theta}_n = h(\hat{\tau}_n)$ .

For any  $\tau$ ,  $L_n(\tau) = \prod_{i=1}^n f(x_i; h(\tau)) = \prod_{i=1}^n f(x_i; h(\tau); \theta) = L_n(\theta)$  where  $\theta = h(\tau)$ .

$$\text{Hence } \forall \tau, L_n(\tau) = L_n(\theta) \leq L_n(\hat{\theta}_n) = L_n(\hat{\tau}_n)$$

ex  
Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\theta, 1)$ . MLE of  $\theta$  is  $\hat{\theta}_n = \bar{X}_n$ .  
Let  $\tau = e^\theta$ . Then MLE for  $\tau$  is  $\tilde{\tau}_n = e^{\hat{\theta}_n} = e^{\bar{X}_n}$ .

Def  
The Score function is  $s(X; \theta) = \frac{\partial}{\partial \theta} \ln(f(x; \theta))$

Def  
Fisher information is  $I_n(\theta) = \nabla_\theta \left( \sum_{i=1}^n \ln(f(X_i; \theta)) \right) = \sum_{i=1}^n I_\theta(X_i; \theta)$   
where  $I_\theta$ ,  $I(\theta) := I_1(\theta)$

Prop  
 $I_n(\theta) = nI(\theta)$ . Also,  $I(\theta) = -E_\theta \left( \frac{\partial^2}{\partial \theta^2} \ln(f(x; \theta)) \right) = - \int_X \left( \frac{\partial^2}{\partial \theta^2} \ln(f(x; \theta)) \right)^2 f(x; \theta) dx$

Prop (Asymptotic Normality)

Let  $se = \sqrt{I(\hat{\theta}_n)}$ . Under appropriate regularity conditions the following hold:

1.  $se \approx \sqrt{\frac{1}{I(\theta)}}$ , then  $\frac{(\hat{\theta}_n - \theta)}{se} \xrightarrow{D} N(0, 1)$

2. Let  $\hat{se}_n = \sqrt{\frac{1}{I_n(\hat{\theta}_n)}}$  then  $\frac{(\hat{\theta}_n - \theta)}{\hat{se}_n} \xrightarrow{D} N(0, 1)$

So,  $\hat{\theta}_n \xrightarrow{D} N(\theta, \hat{se}_n^2)$    
we can then construct normal-based asymptotic confidence intervals for  $\theta$ .

$$\frac{P_{cap}}{L_c + C_n} = [\zeta_n, \bar{\zeta}_n] = (\hat{\Theta}_n - z_{\Theta_n^*} \hat{\zeta}_n, \hat{\Theta}_n + z_{\Theta_n^*} \hat{\zeta}_n)$$

$\text{Tr}_{\text{Lc}\mu} \tilde{\rho}_0(\theta \epsilon c_n) \rightarrow -\infty$  as  $n \rightarrow \infty$

Proof

Let  $\varepsilon$  be standard normal RV. Then  
 $Q_\theta(\theta \in C_n) = P_0(\widehat{\Theta}_n - z_{\alpha/2} \widehat{S}_n \leq \theta \leq \widehat{\Theta}_n + z_{\alpha/2} \widehat{S}_n) = P_0\left(-z_{\alpha/2} \leq \frac{\widehat{\Theta}_n - \theta}{\widehat{S}_n} \leq z_{\alpha/2}\right) \rightarrow P(x_{m_k} \leq Z \leq x_{m_k}) = 1 - \alpha$

For  $\alpha = 0.05$ ,  $z_{\alpha/2} = 1.96 \approx 2$ , so  $\hat{\theta}_n \pm 2\hat{\sigma}_n$  is an approx. 95% confidence interval.

*ext* When you read an opinion poll in a newspaper,

You see statements like "The poll is accurate to within one point 95% of the time." By this they are giving 95% confidence interval of the form  $\hat{\theta} \pm 2\sigma_{\hat{\theta}}$ .

By this, they were giving 95% confidence interval of the form  $\hat{\theta}_n \pm 2\hat{\sigma}_n$ .

Ex Let  $X_1, \dots, X_n$  be  $\text{Be}(\theta)$ . The MLE is  $\hat{\theta}_n = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ .

$f(x; \theta) = \theta^x (1-\theta)^{1-x}$ ,  $\ln(f(x_i; \theta)) = x_i \ln \theta + (1-x_i) \ln(1-\theta)$

$s(X; \theta) = \frac{x}{\theta} - \frac{1-x}{1-\theta}$  and  $s'(x; \theta) = \frac{x}{\theta^2} + \frac{1-x}{(1-\theta)^2} = \frac{1}{\theta(1-\theta)}$

Thus  $I(\theta) = E_\theta(s'(x; \theta)) = \frac{\theta}{\theta^2} + \frac{(1-\theta)}{(1-\theta)^2} = \frac{1}{\theta(1-\theta)} = \frac{1}{V(X)}$

Hence,  $\hat{s}\hat{\sigma}_n = \frac{1}{\sqrt{I(\hat{\theta}_n)}} = \frac{1}{\sqrt{n I(\hat{\theta}_n)}} = \left( \frac{\hat{\theta}(1-\hat{\theta})}{n} \right)^{1/2}$

An approx. 95% conf. interval is:  $\hat{\theta}_n \pm 2 \left( \frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n} \right)^{1/2}$

Q) Let  $X_1, \dots, X_n$  be  $N(\theta, \sigma^2)$  where  $\sigma^2$  is known.  
 $s(x; \theta) = (x - \theta)/\sigma$ , and  $s'(x; \theta) = \frac{1}{\sigma}$ , so  $I_1(\theta) = \frac{1}{\sigma^2}$ .

The MLE of  $\theta$  is  $\hat{\theta} = \bar{x}$ ,  $\bar{x} \sim N(\theta, \sigma^2)$

In this case the Normal approx. is actually exact.

### Exercise

Find 95% conf. int. for  $\mu$ .  $X_1, \dots, X_n$  iid  $P_o(\lambda)$

### Optimality of MLE:

Suppose  $X_1, \dots, X_n \sim N(\theta, \sigma^2)$ . The MLE of  $\theta$  is  $\bar{X}$ .

Another reasonable estimator is sample median  $\tilde{\theta}_n$ .

We know that MLE satisfies  $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D} N(0, \sigma^2)$

One can show that  $\tilde{\theta}_n$  satisfies  $\sqrt{n}(\tilde{\theta}_n - \theta) \xrightarrow{D} N(0, \sigma^2 \frac{1}{2})$

So, the median converges to the right value  $\theta$  but has a larger variance than MLE.

More generally, consider two estimators  $T_n$  and  $U_n$ .

Suppose  $\sqrt{n}(T_n - \theta) \xrightarrow{D} N(0, \sigma_T^2)$  and  $\sqrt{n}(U_n - \theta) \xrightarrow{D} N(0, \sigma_U^2)$

Def

The asymptotic relative efficiency (ARE) is given by

$$ARE(T, U) = \frac{\sigma_U^2}{\sigma_T^2}$$

So,  $ARE(\tilde{\theta}_n, \bar{X}) = \frac{2}{1} = 0.63$ , in words: if you use the median-based estimate  $\tilde{\theta}_n$ , then you are effectively using only a fraction of the samples (0.63%).

Prop.

If  $\hat{\theta}_n$  is the MLE and  $\tilde{\theta}_n$  is any other estimator, then  
 $ARE(\hat{\theta}_n, \tilde{\theta}_n) \leq 1$ , so MLE is asymptotically optimal

All of this MLE stuff is only as good as the  $\Theta$ -parametric family of models!

## Delta Method

Let  $\tau = g(\theta)$ ,  $g$  smooth  $g \in C^1$ .

Due to equivariance of MLE we see that  $\hat{\tau} = g(\hat{\theta})$  but what about the distribution of  $\hat{\tau}_n$ ?

Proof (The Delta Method)

If  $\tau = g(\theta)$  where  $g$  is differentiable and  $g'(\theta) \neq 0$ , then

$$\frac{(\hat{\tau}_n - \tau)}{\hat{s}\hat{\sigma}(\hat{\tau})} \xrightarrow{D} N(0, 1) \text{ where } \hat{\tau}_n = g(\hat{\theta}_n) \text{ and } \hat{s}\hat{\sigma}(\hat{\tau}_n) = |g'(\hat{\theta}_n)| \hat{s}\hat{\sigma}(\hat{\theta}_n).$$

Hence, if  $c_n = (\hat{\tau}_n - \tau) \cdot \hat{s}\hat{\sigma}(\hat{\tau}_n), \hat{\tau}_n - \tau \approx \hat{s}\hat{\sigma}(\hat{\tau}_n)$ , then

$$P(c_n \in C_n) \xrightarrow{n \rightarrow \infty} 1 - \alpha$$

Ex Let  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Beta}(\theta)$  and let  $\psi = g(\theta) = \lambda_n \left( \frac{\theta}{1-\theta} \right)$

The Fisher information function is  $I(\theta) = \frac{1}{\theta(1-\theta)}$  so the estimated se of the MLE  $\hat{\theta}_n$  is  $\hat{s}\hat{\sigma} = \sqrt{\frac{\lambda_n(1-\hat{\theta}_n)}{\hat{\theta}_n(1-\hat{\theta}_n)}}$ . The MLE of  $\psi$  is  $\hat{\psi} = \lambda_n \left( \frac{\hat{\theta}_n}{1-\hat{\theta}_n} \right)$ .

Since  $g'(\theta) = \frac{1}{\theta(1-\theta)}$ , according to the delta method,

$$\hat{s}\hat{\sigma}(\hat{\psi}_n) = |g'(\hat{\theta}_n)| \hat{s}\hat{\sigma}(\hat{\theta}_n) = \frac{1}{\sqrt{n \hat{\theta}_n(1-\hat{\theta}_n)}} \text{ and an approx. 95% ci. is } \hat{\psi}_n \pm \frac{2}{\sqrt{n \hat{\theta}_n(1-\hat{\theta}_n)}}.$$

## Exercise

$X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ . Suppose  $\mu$  is known but not  $\sigma^2 > 0$ . We want to estimate  $\psi = \lambda_n(\sigma)$ .

## hints

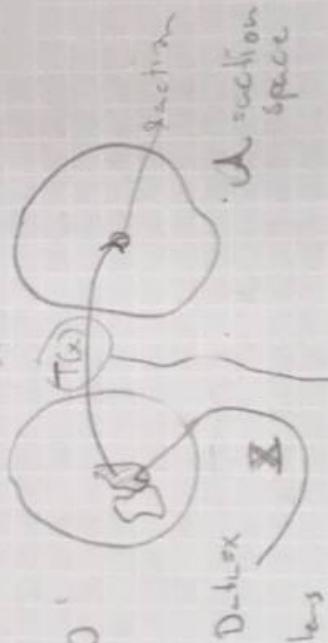
- Find MLE (log-likelihood etc.)
- $\hat{\psi} = \lambda_n(\hat{\sigma}_n)$
- find  $\hat{s}\hat{\sigma}(\hat{\sigma})$
- find  $s\sigma(\hat{\sigma})$

$$\frac{\partial \lambda_n}{\partial \sigma} = \sqrt{\frac{\lambda_n(\sigma)}{n}}$$

$$\hat{s}\hat{\sigma}(\hat{\sigma}_n) = \frac{1}{\hat{\sigma}_n} \cdot \frac{\hat{\sigma}_n}{\sqrt{2\pi}} = \frac{1}{\sqrt{2\pi}}$$

$$95\% \text{ ci.: } \hat{\psi}_n \pm \frac{2}{\sqrt{2\pi}}$$

Generally in decision theory:



Decision problems

estimation problem

$$\Lambda = \{0, T \text{ is } \hat{\Theta}_n \text{ (estimator)}$$

$\sqrt{n}(\theta - \theta_0)$

parametric

$$X_1, \dots, X_n \sim L_0$$

$$X_1, \dots, X_n \sim F \in \{\text{All DF}\}$$

$$\begin{aligned} \text{Point estimators!} & \cdot \text{MOME} \\ & \cdot \text{MLE } \hat{\theta}_n \end{aligned}$$

Set

$$\text{estimators: } [\underline{C}_n, \bar{C}_n], C_n \text{ (confidence interval)} \quad \cdot \quad [\underline{F}_n, \bar{F}_n] \text{ (interval in function space)}$$

(Dvoretzky-Kiefer-Wolfowitz) ineq.

Another Decision Problem: (hypotheses testing)

$$\Lambda = \{0, \beta, T \text{ is Test statistic}$$

space of hypotheses

scientific  
hypothesis  
Null hypothesis



Hyp. testing Problem

Should  $H_0$  be rejected? don't know

Data  $(X_1, \dots, X_n) \sim \text{Law}(H_0)$  or  $\text{Law}(H_A)$

So, the summing of outcome of hyp. testing:

We want to minimize  $\alpha$  at some "significance level", typically  $\alpha = 0.05$   
we also want to maximize Power for given size  $\alpha$  test

(fail to reject)

Retain Null | Reject Null

$H_0$  true |  $\checkmark$

Type II error |  $\beta$

$1 - \beta$  | Power

false pos.

$H_0$  false |  $\checkmark$

Type I error |  $\alpha$

$\checkmark$  |  $\alpha$

The Null hypothesis: "The disease rate is the same in two groups"

Alt. hypothesis: "The disease rate is not the same".

Formally

We partition index sets into two disjoint sets  $\Theta_0, \Theta_1$ ,  
and we want to test

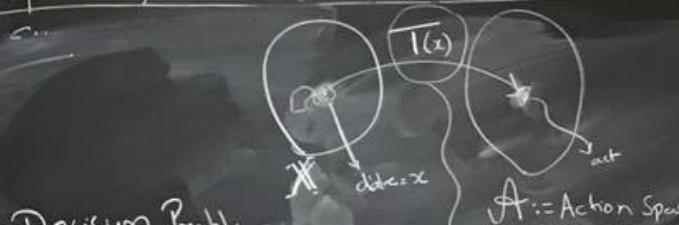
$$H_0: \theta \in \Theta_0 \text{ versus } H_1: \theta \in \Theta_1$$

↑  
Null Hyp.  
e.g.  $\{\bar{\theta}\}$   
 $\subseteq \Theta_0$

Def. Rejection region

$Q \subset \mathbb{R}^n$  s.t. if  $\hat{\theta}$  in  $Q$ , we reject  $H_0$ .

optimality of MLE ; Delta Method ; Hypothesis Testing ; Computer Exs?



Decision Problems

Estimation Problem

$A = \{\Theta\}$ ,  $\hat{T} = \hat{\Theta}_n$  is called an estimator

parametric / non-parametric

Point Estimator  
• MOME  
• MLE  $\hat{\Theta}_n$

Set Estimators

$[\underline{C}_n, \bar{C}_n], C_n \quad [\underline{\hat{F}}_n, \hat{F}]$  Cilivensko - Cantelli Lemma

Another Decision Problem is Hypothesis Testing Problem.

$A = \{0, 1\}$

$T$  is called Test Statistic

Falsifiability

Karl Popper

Demarcation Problem

Science?  
vs. Nonsense?

A scientific hypothesis is one that is falsifiable

Rejection Region

"Acceptance Region"

X

T

0

1

H<sub>0</sub>

H<sub>1</sub>

Nail Hypothesis

A

Hyp Testing Problem is

Should H<sub>0</sub> be rejected?

e.g. The Null hypothesis: "The disease rate is the same in two groups".

Alt. hypothesis: "The disease rate is not the same".

Formally

We partition index set  $\Omega$  into two disjoint sets  $\Omega_0, \Omega_1$ , and we want to tell

$$\begin{array}{c} H_0: \theta \in \Omega_0 \text{ versus } H_1: \theta \in \Omega_1 \\ \uparrow \\ \text{Null hyp.} \\ \text{e.g. } \{\bar{\theta}\} \end{array}$$

Let  $X$  be data,  
 $X: \Omega \rightarrow \mathbb{X}$   
Alt. hyp.

Def. Rejection region

$R \subseteq \mathbb{X}$  s.t. if data in  $R$ , we reject  $H_0$ .  
Then  $R$  is the rejection region

Outcomes of test:	Fail to reject $H_0$		Reject $H_0$
	$H_0$ true	Type II error	Type I error
$H_0$ false $\Leftrightarrow H_1$ true			✓

So, if  $x \in R$ , we reject  $H_0$ , otherwise we fail to reject  $H_0$ .

Typically, the rejection region is of the form

$$R = \{x: T(x) > c\}$$

where  $T$  is a test statistic and  $c$  is a critical value.

The problem of hyp. testing is to find appropriate  $T$  and  $c$ .

Q&

The power function of a test with rejection region  $R$  is

$$\beta(\theta) := P_{\theta}(X \in R).$$

The size of a test is  $\alpha := \inf_{\theta \in \Theta_0} \beta(\theta)$

A test is said to have level  $\alpha$  if its size is  $\leq \alpha$ .

(Types) of hypothesis:

- Simple hypothesis:  $H_0: \theta = \theta_0, H_1: \theta \neq \theta_0$
- Composite hypothesis:  $H_0: \theta < \theta_0, H_1: \theta > \theta_0$  or  $H_0: \theta \geq \theta_0, H_1: \theta < \theta_0$

Types of tests:

- Two-sided test:  $H_0: \theta = \theta_0, H_1: \theta \neq \theta_0$
- One-sided test:  $H_0: \theta \leq \theta_0, H_1: \theta > \theta_0$  or  $H_0: \theta \geq \theta_0, H_1: \theta < \theta_0$

Typically, tests are two-sided under a simple hypothesis.

Ex Let  $X_1, \dots, X_n$  iid  $N(\mu, \sigma^2)$ , where  $\sigma$  is known.

We want to test  $H_0: \mu \leq 0$  vs  $H_1: \mu > 0$ . ( $\Theta_0 = (-\infty, 0], \Theta_1 = (0, \infty)$ )

consider the test: reject  $H_0$  if  $T > c$  where  $T = \bar{X}_n$ ,

so the rejection region is  $R = \{(X_1, \dots, X_n) : T(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n X_i > c\}$

Let  $Z \sim N(0, 1)$ . Then the power function is  $\beta(\mu) = P_{\mu}(T > c) = P_{\mu}\left(\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} > \frac{\sqrt{n}(c - \mu)}{\sigma}\right) = P\left(Z > \frac{\sqrt{n}(c - \mu)}{\sigma}\right) = 1 - \Phi\left(\frac{\sqrt{n}(c - \mu)}{\sigma}\right)$ . Thus  $\beta(\mu)$  is increasing in  $\mu$ :

$$\text{Hence, size is } \sup_{\mu \in \Theta_0} \beta(\mu) = \beta(0) = 1 - \Phi\left(\frac{\sqrt{n}c}{\sigma}\right) \quad (\mu \leq 0)$$

For a size  $\alpha$  test, we set size =  $\alpha$  and solve for  $\mu$ :  $\alpha = 1 - \Phi\left(\frac{\sqrt{n}c}{\sigma}\right) \Rightarrow \mu = c = \Phi^{-1}(1-\alpha)\frac{\sigma}{\sqrt{n}}$

So, we reject the null hypothesis if  $\bar{X}_n > \Phi^{-1}(1-\alpha)\frac{\sigma}{\sqrt{n}}$  for a size  $\alpha$  test or equivalently, when  $\frac{\sqrt{n}(\bar{X}_n - 0)}{\sigma} > \Phi^{-1}(1-\alpha) = Z_{\alpha}$

It is desirable to find the test with the highest power under  $H_1$  ( $\theta \neq \theta_0$ ) among all size  $\alpha$  tests. Such a test (if it exists) is called the most powerful size  $\alpha$  test. (Generally very hard to find)

A common test is the so-called Wald Test:  
 let  $\hat{\theta}_n \in \mathbb{R}$ ,  $\hat{\theta}_n$  be an estimator of  $\theta$  with realization (estimate)  $\hat{\theta}_{n,0}$ ,  
 and let  $\hat{s}_{\theta,n}$  be the estimated standard error of  $\hat{\theta}_{n,0}$ .

Def Consider testing  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ .

Assume  $\hat{\theta}_n$  is asymptotically normal,  $\frac{\hat{\theta}_n - \theta_0}{\hat{s}_{\theta,n}} \xrightarrow{D} N(0, 1)$ .

Then the size  $\alpha$  Wald test is:

$$\text{reject } H_0 \text{ when } |W| > z_{\alpha/2}, \quad W = \frac{\hat{\theta}_n - \theta_0}{\hat{s}_{\theta,n}}$$

Prop.

Asymptotically, the Wald test has size  $\alpha$ , i.e.

$$P_{\theta_0}(|W| > z_{\alpha/2}) \rightarrow \alpha \text{ as } n \rightarrow \infty$$

Proof

$$\begin{aligned} &\text{Under } H_0: \theta = \theta_0, \quad \frac{(\hat{\theta}_n - \theta_0)}{\hat{s}_{\theta,n}} \xrightarrow{D} N(0, 1) \quad (\text{by assumption}), \text{ hence} \\ &\text{the probability of rejecting } H_0 \text{ when } H_0 \text{ is true is } P_{\theta_0}(|W| > z_{\alpha/2}) = P_{\theta_0}\left(\left|\frac{(\hat{\theta}_n - \theta_0)}{\hat{s}_{\theta,n}}\right| > z_{\alpha/2}\right) \rightarrow \\ &\rightarrow P_{\theta_0}(|Z| > z_{\alpha/2}) = \alpha \end{aligned}$$

Remark

Alternatively,  $W = \frac{(\hat{\theta}_n - \theta_0)}{\hat{s}_{\theta,n}}$  where  $\hat{s}_{\theta,n}$  is the standard error computed with  $\theta = \theta_0$ .  
 directly. So, one can use  $\hat{s}_{\theta,n}$  or  $s_{\theta,n}$  for a valid Wald test.

Let's consider the power of the Wald test when  $H_0$  is false.

Prop.

Suppose the true value of  $\theta$  is  $\theta^* \neq \theta_0$ . The power  $\beta(\theta^*)$ , i.e. prob. of rejecting  $H_0$ , is given approximately by  $\beta(\theta^*) = 1 - \Phi\left(\frac{\theta_0 - \theta^*}{\hat{s}_{\theta,n}} + z_{\alpha/2}\right) + \phi\left(\frac{\theta_0 - \theta^*}{\hat{s}_{\theta,n}} - z_{\alpha/2}\right)$

Note  
 $\hat{s}_{\theta,n} \rightarrow 0$  as  $n \rightarrow \infty$  so power is large if  
 i)  $|\theta_0 - \theta^*|$  is large  
 ii)  $n$  is large

Prop Wald test &  $1-\alpha$  asympt conf interval  $\hat{\theta}_n \pm z_{\alpha/2} s_{\theta_n}$

The size  $\alpha$  Wald test rejects  $H_0: \theta = \theta_0$  vs.  $H_1: \theta > \theta_0$ ,  
iff  $\theta_0 + \hat{C}_n = [\hat{\theta}_n, \bar{\theta}_n] \neq [\theta_0, \hat{\theta}_n, \bar{\theta}_n, \theta_0]$ .

Thus, testing the hypothesis is equivalent to checking if the null value  $\theta_0$  is in the conf. interval.

Ex Compare two predictive algorithms (A/B testing).

Suppose the first algorithm gives  $X$  many incorrect predictions out of  $n$  trials.  
Tand the second  $Y$   $\frac{1}{n}$   $\frac{1}{n}$   $\dots$   $\frac{1}{n}$   $n$  trials.

Assuming  $X, Y$  independent binomially distributed R.V.s, test the null hypothesis that their probabilities for incorrect predictions are the same.

sol Let  $X \sim \text{Bin}(n, \theta)$ ,  $Y \sim \text{Bin}(n, \theta)$ . We want to test  $H_0: \theta = \theta_0$  or  $\theta > \theta_0$  where  $\delta_0 = \theta - \theta_0$ . So,  $H_0: \theta = \theta_0$ ,  $H_1: \theta > \theta_0$  we can do Wald test.

The MLE of  $\theta$  is  $\hat{\theta} = \hat{\theta}_{\min} = \bar{\theta}_n - \hat{\theta}_n$  by equivalence prop of MLE.

The estimated std. error of  $\hat{\theta} - \hat{\theta}_0$  is

$$\hat{s}_{\theta_n} = \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} + \frac{\hat{\theta}_0(1-\hat{\theta}_0)}{n}$$

$$\begin{aligned} X &\sim \text{Bin}(n, \theta) \Rightarrow V(X) = n\theta(1-\theta) \\ \hat{\theta}_n &= \bar{X}_n \Rightarrow V(\hat{\theta}_n) = \frac{1}{n}V(X) = \frac{n\theta(1-\theta)}{n} \\ &= \theta(1-\theta) \end{aligned}$$

The size  $\alpha$  Wald test is to reject  $H_0$  when  
 $|W| > z_{\alpha/2}$ , where  $W = \frac{\hat{\theta} - \theta_0}{\hat{s}_{\theta_n}} = \frac{\bar{\theta}_n - \hat{\theta}_n}{\hat{s}_{\theta_n}}$ .

Suppose you observe 18 bad prediction out of 117 trials for alg. 1  
and  $21 - \frac{1}{117} - \frac{1}{117} - \dots - \frac{1}{117} = 111 - 11 = 102$

Does the size  $\alpha = 0.05$  Wald test reject the null hypothesis that  $\theta = \theta_0$

Exercise

Suppose  $X_1, \dots, X_n \sim \exp(\lambda)$  model waiting times at bus stop and you have values 7.6, 9.2, 11.8, 6.3, 13.2, 10.6, 6.3, 7.8, 8.9, 10.2.

First, obtain MLE  $\hat{\lambda}_n$  of  $\lambda$ ,  $s_{\lambda_n}$  of  $\hat{\lambda}_n$ ,  $(1-\alpha)$  c.i. for  $\lambda$ .

Then give a size  $\alpha = 0.05$  Wald Test to reject  $H_0: \lambda = \lambda_0$

## P-values

Instead of just reporting "reject  $H_0$ " or "fail to reject  $H_0$ " at a given size  $\alpha$ , we could make a more informed report of the test.  
 If a test rejects at level  $\alpha$ , it will also reject at any  $\alpha' < \alpha$ .  
 But there will be a smallest  $\alpha$  at which we fail to reject  $H_0$ .

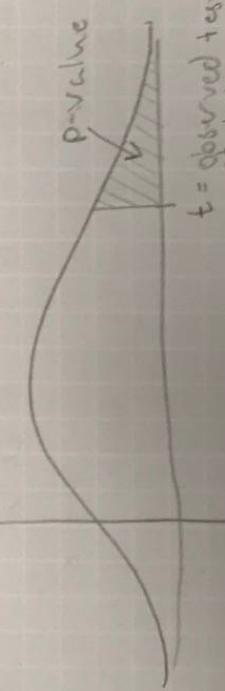
Def Suppose for every  $\alpha(\theta_0)$  we have a test of size  $\alpha$  with rejection region  $R_\alpha$ . Then, p-value is inf  $\{\alpha : T(x_0, \theta_0) \in R_\alpha\}$ .  
 i.e. p-value is the smallest  $\alpha$  at which we reject  $H_0$ .

### Remarks

p-value is a measure of evidence against  $H_0$ : the smaller the p-value,  
 the stronger the evidence against  $H_0$ .

p-value	evidence against $H_0$	
	very strong	strong
< 0.01		
0.01 - 0.05		
0.05 - 0.10		weak
> 0.1	kittle or no evidence	

↑ PDF of test statistic under  $H_0$



To find p-value for a particular test statistic  $T$  we find  $\alpha$  s.t.  
 the observed statistic  $t$  is just at the boundary of rejection region  $R$   
 $\Rightarrow$  p-value is tail area  $P(T > t)$

(if  $T$  is of the form  $T > c$ )

## Pearson's $\chi^2$ test for multinomial data

Recall, if  $X = (X_1, \dots, X_k)$  has Multinomial( $n, \theta$ ) distribution.

$$\begin{cases} \sum_{k=1}^K = \sum_{j=1}^k x_j \\ K=2 \Rightarrow \text{binomial} \end{cases}$$

$$P((X_1, X_2) = (x_1, x_2); n, \theta_1, \theta_2) = \binom{n}{x_1, x_2} \theta_1^{x_1} \theta_2^{x_2} = \frac{n!}{(n-x_1)! x_1!} \theta_1^{x_1} \theta_2^{x_2}$$

Multinomial:

$$2 \leq k < \infty$$

$$P((X_1, \dots, X_n) = (x_1, \dots, x_n); n, \theta_1, \dots, \theta_n) = \binom{n}{x_1, x_2, \dots, x_n} \theta_1^{x_1} \theta_2^{x_2} \dots \theta_n^{x_n} = \frac{n!}{x_1! x_2! \dots x_n!} \theta_1^{x_1} \theta_2^{x_2} \dots \theta_n^{x_n}$$

$$\text{Let } \Theta \in \Delta^{k-1} := \{(\theta_1, \dots, \theta_k) \in \mathbb{R}^k : \sum_{i=1}^k \theta_i = 1, \theta_i \geq 0 \forall i\}$$

we can check that the MLE of  $\theta = (\theta_1, \dots, \theta_n)$ ,  $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_n) = \left(\frac{x_1}{n}, \frac{x_2}{n}, \dots, \frac{x_n}{n}\right)$

Let  $\theta_0 = (\theta_{01}, \dots, \theta_{0n})$  be some fixed vector and suppose we want to test  $H_0: \theta = \theta_0$  vs.  $H_1: \theta \neq \theta_0$ .

Def Pearson  $\chi^2$ -statistic is

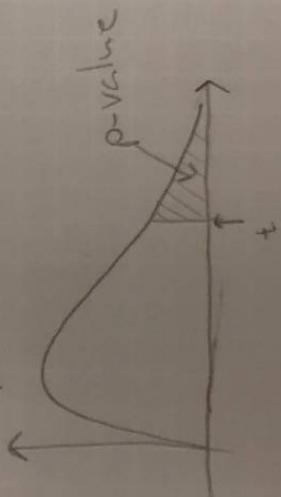
$$T = \sum_{j=1}^k \frac{(x_j - n\theta_{0j})^2}{n\theta_{0j}} = \sum_{j=1}^k \frac{(x_j - E_{\theta_0}(x_j))^2}{E_{\theta_0}(x_j)}$$

$x_j$  under  $H_0$ .

Prop.

Under  $H_0$ ,  $T \rightsquigarrow \chi^2_{k-1}$ . Hence the test rejects  $H_0$  if  $T \rightarrow \chi^2_{k-1, \alpha}$  has asymptotic level  $\alpha$ .

The p-value is  $P(\chi^2_{k-1} > t)$ , where  $t$  is observed test statistic.



### ex (Mendel's peas)

Mendel bred peas with round yellow seeds and green wrinkled seeds. There are 4 types of progeny:  $r_y$ ,  $r_g$ ,  $w_y$ ,  $w_g$ . The number of each type is modeled as a multinomial RV with probs  $\theta = (\theta_1, \theta_2, \theta_3, \theta_4)$ .

Mendel's theory of inheritance predicts that  $\theta = \theta_0 = (\frac{9}{16}, \frac{3}{16}, \frac{3}{16}, \frac{1}{16})$ . In  $n=556$  trials, Mendel observed  $X = (315, 101, 108, 32)$ .

We will test  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ .

Since  $n\theta_1 = 312.35$ ,  $n\theta_2 = n\theta_3 = 104.25$ ,  $n\theta_4 = 34.75$ ,

$$\text{the } \chi^2 \text{ test statistic is } \frac{(315 - 312.35)^2}{312.35} + \frac{(101 - 104.25)^2}{104.25} + \frac{(108 - 104.25)^2}{104.25} + \frac{(32 - 34.75)^2}{34.75} = 0.47.$$

The  $\alpha=0.05$  value for a  $\chi^2_3$  is 7.815. Since  $0.47 < 7.815$  we do not reject the null. The p-value is  $P(\chi^2 > 0.47) = 0.93$  ↑ table lookup which is not evidence against  $H_0$ .

### The $\chi^2$ -distribution (Prelude to Pearson's $\chi^2$ -test)

Def Let  $z_1, \dots, z_k \stackrel{\text{iid}}{\sim} N(0, 1)$  and  $Y = \sum_{i=1}^k z_i^2$ . Then  $Y$  has  $\chi^2$  distr.

with  $k$  degrees of freedom (d.f.) and we write  $Y \sim \chi^2_k$ .

The PDF of  $Y$  is  $f(y; k) = \frac{y^{\frac{k}{2}-1} e^{-y/2}}{2^{k/2} \Gamma(k/2)}$  for  $y > 0$ .

Facts  
 $E(Y) = k$ ,  $V(Y) = 2k$

Def Upper  $\alpha$  quantile is  $\chi^2_{k,\alpha} := F_{\chi^2,k}^{-1}(1-\alpha; k)$  where  $F$  is CDF of  $Y$ .  
i.e.  $P(Y > \chi^2_{k,\alpha}) = \alpha$

### Permutation test

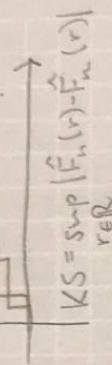
This is a non-parametric test (in A/B testing context).

So we want to test if two sets of samples come from the same distribution or not.  $X_1, \dots, X_m \sim F_X$  (all DFs) and  $Y_1, \dots, Y_n \sim F_Y$  (all DFs)

$H_0: F_X = F_Y$  vs  $H_1: F_X \neq F_Y$

Let  $T(X_1, \dots, X_m, Y_1, \dots, Y_n)$  be a given test statistic

examples of  $T$ :  $T = |\bar{X}_m - \bar{Y}_n|$  or  $T = KS$



Let  $N = m+n$  and consider all  $N!$  permutations of the data  $X_1, \dots, X_m, Y_1, \dots, Y_n$ .

Under  $H_0$  (distributions equal), all permutations should have same prob!

For each permutation, compute  $\bar{T}$  and denote them  $T_1, \dots, T_{N!}$

Under  $H_0$ , all the  $T_i$ 's have equal prob. So the distribution

$P_0$ , under  $H_0$ , puts mass  $\frac{1}{N!}$  on each  $T_i$  where  $1 \leq i \leq N!$

This  $P_0$  is called the permutation distribution of  $T$ .

Let  $t_{obs}$  be the observed statistic. Assuming  $T$  is uniform we reject for large values of  $T$ , the p-value is:

$$p\text{-value} = P_0(T > t_{obs}) = \frac{1}{N!} \sum_{j=1}^{N!} \mathbf{1}_{(T_j > t_{obs})}$$

son's

$t_1 / t_2$

does  
assume

### for (Toy example)

Suppose data are  $(X_1, X_2, Y) = (1, 9, 3)$ .

$$\text{Let } T(X_1, X_2, Y) = |X_1 - \bar{Y}| + |X_2 - \bar{Y}| = 2$$

Permutation	T	prob under $H_0$	p-value = $P(T > 2) =$
(1, 9, 3)	2 = $t_{\text{obs}}$	$\frac{1}{6} = k_1$	$\frac{1}{6} = 0.1667$
(4, 1, 3)	2	$\frac{1}{6}$	$\frac{1}{6} = 0.1667$
(1, 3, 9)	7	$\frac{1}{6}$	$\frac{1}{6} = 0.1667$
(3, 1, 9)	7	$\frac{1}{6}$	$\frac{1}{6} = 0.1667$
(3, 9, 1)	5	$\frac{1}{6}$	$\frac{1}{6} = 0.1667$
(4, 3, 1)	5	$\frac{1}{6}$	$\frac{1}{6} = 0.1667$

It's not practical to do  $N!$  permutations. However, we can approximate the p-value by sampling uniformly from the set of all permutations of  $\{1, 2, \dots, N\}$  due LHN.

### Alg. for Permutation test

1. Compute observed value of test statistic.  
 $t_{\text{obs}} = T(X_1, \dots, X_n, Y_1, \dots, Y_n)$
2. Randomly permute the data. Compute statistic again using permuted data.

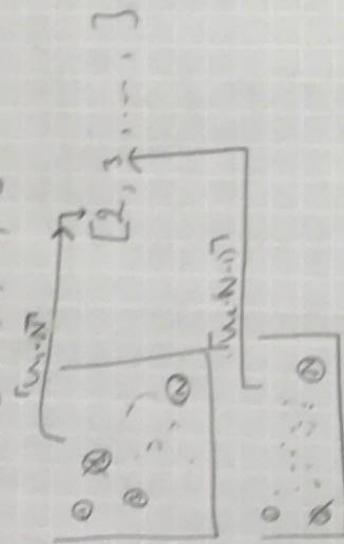
3. Repeat 2. B times and let  $T_1, \dots, T_B$  denote the results

4. The approx. p-value is  $\frac{1}{B} \sum_{j=1}^B \mathbb{1}(T_j > t_{\text{obs}})$

Uniform sampling:

$$U_{1:N} \sim \text{Uniform}(0,1)$$

$$\text{index}(x) \in \{1, 2, \dots, N\}$$



### Likelihood Ratio Test Statistic (LRTS)

with smaller than

A more general test when  $|\Theta| = r < n$

Def

Consider testing the following:

$$H_0: \theta \in \Theta_0 \text{ vs. } H_1: \theta \notin \Theta_0 \text{ (or } \theta \in \Theta \setminus \Theta_0 := \Theta_1)$$

The LRTS is:  
important  
$$\lambda_n = 2 \ln \left( \frac{\sup_{\Theta \in \Theta_0} L_n(\theta)}{\sup_{\Theta \in \Theta_1} L_n(\theta)} \right) = 2 \ln \left( \frac{L_n(\hat{\theta}_n)}{L_n(\hat{\theta}_{0,n})} \right)$$
 where  $\hat{\theta}_n$  is MLE and  $\hat{\theta}_{0,n}$  is MLE for  $\theta$  restricted to  $\Theta_0$ .

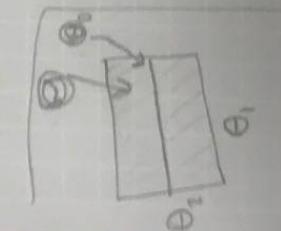
Prop.

Suppose  $\theta = (\theta_1, \dots, \theta_r)$  and let  $\Theta_0 = \{\theta : (\theta_{q+1}, \dots, \theta_r) = (\theta_{0,q+1}, \dots, \theta_{0,r})\}$   
 $|\Theta| = r$ ,  $|\Theta_0| = q$

Let  $\Lambda_n$  be the LRTS under  $H_0: \theta \in \Theta_0$ .

$$\Lambda_n(X_1, \dots, X_n) \xrightarrow[n \rightarrow \infty]{\text{distr}} \chi^2_{r-q}$$

The p-value of the test is  $P(\chi^2_{r-q} > \underline{\lambda_{n,obs}})$   
 $\uparrow$   
observed LRTS



For instance, if  $\theta = (\theta_1, \dots, \theta_d)$  and we want to test  $\theta_1 = \theta_2 = \dots = \theta_d = 0$ , then  $\Lambda_\theta$  has limiting distn. given by  $X_{d+3} = X_d^T Q \Lambda_\theta$ .

Consider a model  $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$  where  $\varepsilon_i \sim N(0, 1)$

for data of the form  $((X_i, Y_i))_{i=1}^n$

We are interested in testing  $H_0: \beta_1 = 0$  vs  $H_1: \beta_1 \neq 0$

Recall

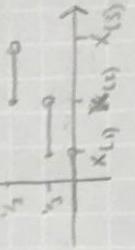
SageMath notebook from day 1:

Cart. on simple linear regression

### Non-parametric Estimation

Recall:  $E[DF] = \sum_{x \in \text{DF}} \mathbb{P}[X_i \in \text{DF}]$

$$E[\hat{F}_n(x)] = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i = x)$$



Proof

$$\text{At any fixed } x, E(\hat{F}_n(x)) = F(x)$$

$$V(\hat{F}_n(x)) = \frac{F(x)(1-F(x))}{n}$$

$$\text{So, MSE (Mean Squared Error)} = \frac{F(x)(1-F(x))}{n} \xrightarrow{n \rightarrow \infty} 0$$

Prop. (Givensko-Cantelli lemma)

Let  $X_1, \dots, X_n \stackrel{iid}{\sim} F$ . Then  $\sup_x |\hat{F}_n(x) - F(x)| \xrightarrow{P} 0$  as  $n \rightarrow \infty$

Prop. (The Dvoretzky-Kiefer-Wolfowitz (DKW) inequality.)

Let  $X_1, \dots, X_n \stackrel{iid}{\sim} F$ . Then, for any  $\varepsilon > 0$ ,

$$P(\sup_x |\hat{F}_n(x) - F(x)| > \varepsilon) \leq 2e^{-2n\varepsilon^2}$$

A non-parametric  $1-\alpha$  confidence band for  $F$ :

$$L(x) = \max \left\{ \hat{F}_n(x) - \varepsilon_k, 0 \right\}, \quad U(x) = \min \left\{ \hat{F}_n(x) + \varepsilon_k, 1 \right\} \quad \text{where } \varepsilon_k = \sqrt{\frac{1}{2n} \ln \left( \frac{2}{\alpha} \right)}$$

Due to DKW  $\#$ , for any fixed but unknown  $F \in \{ \text{DF} \}$ ,

$$P(\{L(x) \leq F(x) \leq U(x)\} \geq 1 - \alpha)$$

## Simple linear regression

Regression is a method of studying the relationship between a response RV Y and a covariate (predictor/feature) RV X.  
 [Galton (1822-1911)].

We summarize using regression function:  $r(x) := E(Y|X=x) = \int x f(y|x) dy$ .

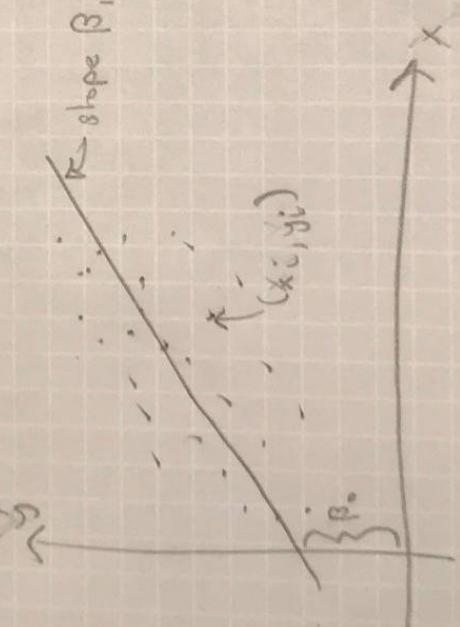
Our goal is to estimate  $r(x)$  from the data of the form

$$(Y_1, X_1), (Y_2, X_2), \dots, (Y_n, X_n) \sim F_{X,Y}(x, y; \theta)$$

Simplest version is called simple linear regression:

$X_i, Y_i$ : real-values.

$r(x) = \beta_0 + \beta_1 x$ , assume  $V(Y|X=x) = \sigma^2$  does not depend on  $x$ .



Def of SLR Model:

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad \text{where } E(\varepsilon_i | X_i) = 0 \text{ and } V(\varepsilon_i | X_i) = \sigma^2$$

We want estimates  $\hat{\beta}_0$  and  $\hat{\beta}_1$  for  $\beta_0, \beta_1$  from the data.

This gives us the fitted line:  $\hat{r}(x) = \hat{\beta}_0 + \hat{\beta}_1 x$

and the residuals:  $\hat{\varepsilon}_i = \hat{y}_i - \hat{r}(x_i) = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)$

$\hat{y}_i = \hat{r}(x_i)$  is the predicted value.

Residual sum of squares (RSS) is given by  $\sum_{i=1}^n \hat{\varepsilon}_i^2$

Def

Least squares estimate (LSE) is the value of  $\hat{\beta}_0$  and  $\hat{\beta}_1$  that minimizes RSS.

$$\hat{\beta}_1 = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}, \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum \hat{\varepsilon}_i^2$$

Prop.

If  $\varepsilon_i \sim N(0, \sigma^2)$ , then MLE = LSE