

Interestori

course page: <http://www.csail.mit.edu/~rcazeigh/ScalabilityInMechanisms>

razaeigh.scazdiin@mech.mit.edu : subject AMS035

Koller ~~vs~~ vs Koller:
Theoretical Foundation for Decision Lists
(HT19)

Sets, maps & Numbers

Sets

A set is a collection of distinct elements.

e.g. {0, 03}. We give names to sets, e.g. A = {0, 03}.

$A = \{0, 0\}$ is not a set (multiset)

$\emptyset = \{\}$ is the empty set.

An element belongs to (or doesn't belong to) a set and we write 'ε' or '∉' respectively. e.g. 0 ∈ {0, 03} and 0 ∉ {0, 03}, etc. etc. depending on whether the element is in the set or not.

Set operations

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

$$\text{e.g. } \{0\} \cup \{03\} = \{0, 03\}$$

$$\{0, 03\} \cup \{0\} = \{0, 0, 03\}$$

$$\{0, 03\} \cup \{03\} = \{0, 03\}$$

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 called domain with exactly one element in the set range (codomain).

formally

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

Inverse image/ preimage

The inverse image of a function $f: X \rightarrow Y$ is $f^{-1}: Y \rightarrow X$ where $f^{-1}(y) = \{x \mid x \in X \text{ and } f(x) = y\}$ or more generally, for any $B \subseteq Y$

$$f^{-1}(B) = \{x \in X \mid f(x) \in B\} \quad \text{and} \quad f^{-1}(y) = \{x \in X \mid f(x) = y\}$$

Note

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

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 Ω .

An event is a subset of the sample space.

Probability is a function

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

- 1) $\forall \text{ event } A, \quad 0 \leq P(A) \leq 1$
- 2) $P(\Omega) = 1$
- 3) $A \cap \emptyset = \emptyset \Rightarrow P(A \cup \emptyset) = P(A) + P(\emptyset)$
- 4) $P(\bigcup_i A_i) = \sum_i P(A_i)$ where A_i are pairwise disjoint

Motivation for axioms:

idea of long-term relative frequency of independent experiments (experimental law)
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 $N(A,n)$ be the number of times A occurs in the first n trials.

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

$$\text{axiom: } 0 \leq \frac{N(A,n)}{n} \leq 1$$

$$\text{(i) } \frac{N(\Omega_1, n)}{n} = \frac{n}{n} = 1$$

"Something" happens every time

$$\text{(ii) } A \cap B = \emptyset \Rightarrow N(A \cup B, n) = N(A, n) + N(B, n)$$

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

Ex 1 Tossing a fair coin. Can construct reels by discrete probability.

$$\Omega = \{HT\} \quad \begin{matrix} HT \\ HT \\ HT \\ HT \end{matrix} \quad \text{Bernoulli random variable. Bernoulli}(\theta)$$

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

Ex 2 N'th Lotto (40 balls)

Label of the 1st ball reel as ω_1 :

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



De Moivre random variable,

$$\text{DeMoivre}(\frac{1}{40}, \frac{1}{40}, \dots, \frac{1}{40}) = \text{DeMoivre}(\theta_1, \theta_2, \dots, \theta_n)$$

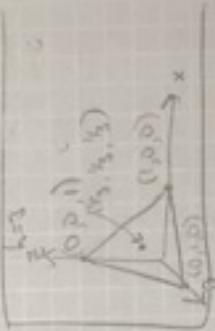
Here we have $\theta_1 = \theta_2 = \dots = \theta_n = \frac{1}{40}$, so Equi-probable DeMoivre variable

$$\begin{aligned} \text{So, what is } P(\text{"an number"}) &= P(\{2, 4, \dots, 40\}) = \\ &= P(\{1\} \cup \{3\} \cup \dots \cup \{40\}) = P(\{1\}) + P(\{3\}) + \dots + P(\{40\}) = \\ &\quad \text{repeatedly} \\ &< 20 \cdot \frac{1}{40} = \frac{1}{2} \end{aligned}$$

Properties:

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

$$2. A, B \text{ a.s. } P(A \cup B) = P(A) + P(B) - P(A \cap B)$$



The domain of P is called a sigma field or sigma algebra,
It's denoted by \mathcal{F}_{Ω} or \mathcal{F}_{Ω} or just \mathcal{F} if Ω is clear from context.
we see that

- $\Omega \in \mathcal{F}$,
- $A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}$,
- $A_1, A_2, \dots, A_n \in \mathcal{F} \Rightarrow \bigcup_{i=1}^n A_i \in \mathcal{F}$ by Kolmogorov.

for Bernoulli experiment (win or not) $\Omega = \{H, T\}$

The triple $(\Omega, \mathcal{F}_A, P)$ is called the probability space.
 If P is a set of probabilities, then $(\Omega, \mathcal{F}_A, P)$ is called a statistical experiment.

Recall

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

The product experiment

$$\Omega = \prod_{i=1}^n \Omega_i, \quad \text{then } P(\{\omega_1, \omega_2, \dots, \omega_n\}) = \prod_{i=1}^n P_i(\omega_i)$$

Conditional probability

$$P(A|B) := \frac{P(A \cap B)}{P(B)} \text{ provided } P(B) > 0$$

is conditional probability a prob.? Yes, basically B is new Ω .

Constructing random graph from tossing unfair coins.

Toss n coins with $P(H) = \theta$ and n^k times.

Graph:

Let $V = \{v_1, \dots, v_n\}$ be class of vertices and let $E \subseteq V^2$ and $|E| = k$ (no vertices, k edges). Let $\{v_1, \dots, v_n\}$ where $A_{ij} = \begin{cases} 1, & \text{if } H \\ 0, & \text{if } T \end{cases}$ in coin tosses so, $E(|E|) = |V|^2 \theta$

$$A = \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1n} \\ A_{21} & A_{22} & \dots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & \dots & A_{nn} \end{bmatrix}$$

Exercise

$$6.11 \quad f(t) = \begin{cases} 0 & t \leq 0 \\ e^{-0.1t} & t > 0 \end{cases}$$

$$\text{Q}(t) = 0.1t$$

$$\text{Q}(t) = 1 - e^{-0.1t} \quad \text{hours}^{-1} \quad \text{hours}^{-1} = 1 - F(t) = 1 - \int_0^t \lambda e^{-\lambda s} ds = 1 - e^{-\lambda t}.$$

$$6.10 \quad e^{-0.1t} = 0.0107 \quad \Rightarrow \quad t = \frac{\ln(0.0107)}{-0.1} = 10.0 \text{ years}$$

6.12

$$\frac{6}{576} \cdot \text{P}(0) = \frac{6}{576}, \quad \text{P}(1) = \frac{511}{576}, \quad \text{P}(2) = \frac{376}{576}, \quad \text{P}(3) = \frac{251}{576}, \quad \text{P}(4) = \frac{154}{576}, \quad \text{P}(5) = \frac{93}{576}, \quad \text{P}(6) = \frac{51}{576}, \quad \text{P}(7) = \frac{26}{576}, \quad \text{P}(8) = \frac{13}{576}, \quad \text{P}(9) = \frac{6}{576}$$

$$\text{P}(0) \in \lambda_0 = \lambda, \quad \lambda = \sqrt{2 \pi / 576}$$

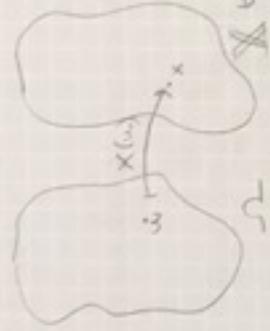
$$\lambda = \frac{576}{2 \pi}, \quad \text{No. of units} \sim P_0(\lambda)$$

$$\text{approx} = [576 \cdot \text{P}(P_0(\lambda) + i) \text{ for } i \text{ in } (0, 1, 2, 3, 4, 5)] = [216.7, 216.4, 216.5, 216.6, 216.3, 216.0]$$

which is very close to actual.

Prelude to Decision theory

Estimation in parametric models



Data Random variable ($R.V.$)
prob. dist. function
function ($f(x)$) of
distribution

Problem

Let X be a RV with values in \mathbb{X} and $\text{Law}(X)$ is law of X .

Assume $\text{L}(X)$ is known up to a finite dimensional parameter θ from the parameter space Θ (bold theta):

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

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

Typically (wLOG)

$$X = (X_1, \dots, X_n), X_i \in \mathbb{X}_i$$

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 RV. X (def.)

$$\{g(\theta) | \theta \in \Theta\}$$

Def.
An estimator T of $g(\theta)$ we admit any $T: \mathbb{X} \rightarrow g(\Theta)$.

Concretely:

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

$$\prod_{i=1}^n = g(\Theta)$$

Indicator function

$$X(\omega) \mathbf{1}_{\{\omega\}} := \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases}$$

$\text{Law}(X) \sim \text{Bernoulli}(\theta) = \mathbf{1}_{\{\omega : \omega \in A\}}$

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

Note that $P_\theta(x_{i=1}^n) = \prod_{i=1}^n P_\theta(x_i; \theta) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} = \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i}$.

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

Thus, the prob. of data (i.e. $P_\theta(x_{i=1}^n)$) only depends on the statistic T .

Now consider another statistic: (sample mean)

$$\bar{T}(x) = \frac{1}{n} \sum_{i=1}^n x_i = \bar{X}_n.$$

Then Θ becomes $\Theta^k([0, 1]^n) \stackrel{k \text{ int}}{=} \Theta^n([0, 1]) = \Theta^{\text{int}}([0, 1])$.

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

$$\text{Def } \hat{x}^{(n)}$$

A sequence $\bar{T}_n := \bar{T}_n(\hat{x}_1, \hat{x}_2, \dots, \hat{x}_n)$ of estimators (each based on a sample of size n) for a parameter θ is called (asymptotically) consistent: if $\forall \delta, \exists \epsilon, \exists n_0 : P_{\theta, n}(|\bar{T}_n(x^{(n)}) - \theta| > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$ or, in shorter notation:

$$\bar{T}_n = \bar{T}_n(x^{(n)}) \xrightarrow{n \rightarrow \infty} \theta \text{ if } \text{Law}(x^{(n)}) = P_\theta$$

\Leftrightarrow (Bon. product ex.)

The estimator $T_n(x_1, \dots, x_m) = \bar{X}_m$ is consistent for θ

Proof

Since x_1, \dots, x_m are iid $\text{Be}(\theta)$ RVs, 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 on 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 (incurring) of the unknown parameter θ is estimated by t . this estimation error is measured by $\text{Loss}(T, \theta)$ under $\text{T}(x^{(n)})$

Natural choices for loss when $\theta \in \mathbb{R}$

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

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

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

Note: $\text{Loss} \in Q.N. \subset \text{Loss}(\text{T}(x), \theta)$ needs to account for random variables.

Def

The risk of an estimator T at parameter θ is

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

\uparrow
Risk function of T

Note Risk might exist
since the expectation
might not exist

↳ class of estimators

$$Q(T^n, \theta) = \min_{T \in \mathcal{T}^n} Q(T, \theta) \text{ for any fixed } \theta \in \Theta$$

We find an estimator T^* that minimizes the whole risk function simultaneously. If such a T^* can be found, then it is called a uniformly best estimator (in \mathcal{T}).

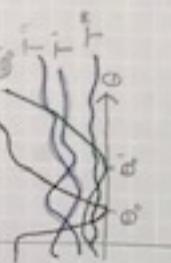
In general, such a UBE won't exist.

Argument

For each $\theta_0 \in \Theta$, consider $T_{\theta_0}(x) = \theta_0$.

so $Q(T_{\theta_0}, \theta_0) = 0 \Rightarrow$ best if θ_0 is true.

If T^* was UBE, then it would have to compete with T_{θ_0} .



Def

Unbiased estimators

Consider an estimator T s.t. E_T exists.

Then $E_\theta(T) - g(\theta)$ is called bias of the estimator.

If $E_\theta(T) = g(\theta)$ for any $\theta \in \Theta$, then the estimator T is called unbiased

for $g(\theta)$

Def

A statistic S is called sufficient for θ if $P_\theta(x \in B | S(X)=s)$ is independent of θ , for all values of s and all events B .



In words, conditional distribution of the data X given that $S(X)$ takes any value does not depend on the parameter θ .

For discrete experiments

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

Let's clarify again.

Originally, the form of X depended on θ ($L_{\text{exp}}(X) = P_\theta(x)$).

After the value of the sufficient statistic $S(X)$ is known, then the conditional law $P_\theta(\cdot | S(x)=t)$ is no longer dependent on θ .

Since we are interested in making inference about θ , the conditional law is uninteresting for our purpose, so we can disregard it.

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

Remark (Exercise to prove $S(x)$ is sufficient)

The data itself is sufficient, i.e.

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

Proof

In $\bigotimes_{i=1}^n \mathcal{B}(\theta)$ e.g., the sample \bar{X}_n is a sufficient statistic.

Proof

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

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

$$P_\theta(x_1=x_1, \dots, x_n=x_n | n\bar{X}_n=t) = \frac{P_\theta(x_1=x_1, \dots, x_n=x_n, n\bar{X}_n=t)}{P_\theta(n\bar{X}_n=t)} \quad (\times)$$

If θ results $\sum_{i=1}^n x_i=t$, then θ is $\frac{(t^t (1-\theta)^{n-t})}{\binom{n}{t} \theta^t (1-\theta)^{n-t}} = \frac{1}{\binom{n}{t}}$ which clearly is indep. of θ .

If for this X_i , $\sum_{i=1}^n X_i$ is θ , then the numerator is 0 which is
indep. of θ , so θ is indep. of θ . \square

Say we want to estimate θ in this example and we limit ourselves
to estimators that are functions of the sufficient statistic \bar{X}_n :
 $T(x) = h(\bar{X}_n)$

Additionally, suppose we limit ourselves further to unbiased estimators:

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

Proof

Under sufficiency and unbiasedness restrictions on allowed estimators
of θ in our $\hat{\Theta}$ Bel(θ) exp., the only possible estimator is \bar{X}_n .

Proof $E(\bar{X}_n)^{\theta}$

$$\begin{aligned} \theta &= E_{\theta} h(\bar{X}_n) - \theta = E_{\theta} (h(\bar{X}_n) - \bar{X}_n) = \underbrace{\sum_{k=0}^{\infty} \binom{n}{k} \left(\frac{1}{n} \right)^k \left(\frac{n-1}{n} \right)^{n-k}}_{C_k} \theta^k (1-\theta)^{n-k} \\ &= (1-\theta)^n \sum_{k=0}^{\infty} C_k r^k \text{ for } r = \frac{\theta}{1-\theta} \end{aligned}$$

Now, if $\theta \in [0, 1]$, then $r \in [0, \infty)$, hence the above polynomial can only
be zero if every C_k is also zero.

Thus, $0 = \binom{n}{k} \left(\frac{1}{n} \right)^k \left(\frac{n-1}{n} \right)^{n-k} \Rightarrow h\left(\frac{k}{n}\right) = \frac{1}{n}$ 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 $h: \Omega \rightarrow \mathbb{R}$:
 $(E_\theta(h(T(X))) = 0 \text{ for all } \theta \in \Theta) \Rightarrow$

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

Intuitively, completeness means there is "no superfluous information" or "no redundancy" in the complete statistic T .

Consider an event $T(X) \in \mathcal{B}$ and suppose $\Pr_\theta(T(X) \in B) = \alpha$ for a fixed θ . By taking $h(1) = \mathbf{1}_{\{T(X) \in B\}}$, completeness $\Rightarrow \Pr_\theta(\mathbf{1}_{\{T(X) \in B\}} = 1) = 1$ for all $\theta \in \Theta$ which means that α is either 0 or 1. Thus, for any event $T(X) \in \mathcal{B}$ which has a non-trivial probability ($\neq 0, 1$), this prob. must depend on θ .

$\forall h \otimes B \in (\Theta)$

The statistic $T(X) = (X_1, \bar{X}_n)$ is sufficient but not complete?

$h(X_1, \bar{X}_n) = X_1 - \bar{X}_n$ has $E(h) = 0$ but $h(X_1, \bar{X}_n)$ is not almost surely 0.

Prop In $\mathbb{R}^{\mathcal{B}(\Theta)}$, $T(X) = \bar{X}_n$ is sufficient and complete.

Proof

Suppose for some function h ,

$$E_\theta(h(\bar{X}_n)) = 0$$

This means that $\sum_{k \in \mathcal{O}} h\left(\frac{k}{n}\right) \binom{n}{k} \theta^k (1-\theta)^{n-k} = 0 \quad \forall \theta \in [0, 1]$.

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

Limits of R.V.s : Chapter 9 in CSE Book pdf 151-156

§1 - conv. of sequence of R.V.

X_1, \dots, X_n conv to ∞ R.V. X :

* In distribution

$X_n \rightsquigarrow X$

* In probability

$X_n \xrightarrow{p} X$

-Markov's ineq., Chebyshev's ineq.

- ③ Suppose T is a sufficient statistic with values in a set \mathbb{T}
and $S: \mathbb{T} \rightarrow \mathbb{S}$ is a one-to-one mapping with values in \mathbb{S} (ie,
there exists an inverse mapping $S^{-1}: \mathbb{S} \rightarrow \mathbb{T}$ such that $S(S^{-1}(t)) = t$ for
each $t \in \mathbb{T}$).
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 $\bigotimes_{i=1}^n \text{Bernoulli}(\theta)$ experiment. Prove this
claim.

Problem Set Week 3

- ① Show that the estimator \bar{X}_n is consistent for $\theta = \bigotimes_{i=1}^n \text{Bernoulli}(p)$ exponent
- ② Suppose the data X in a statistical experiment can take values in a countable set \mathbb{X} (i.e., $\text{Low}(X)$ 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 ... and there ... Shows ...
- ④ In Su ...
- ⑤ Let X_1, \dots, X_n be independent and identically distributed with $\text{Low}(\lambda) \sim \text{Poisson}(\lambda)$, $\lambda \in (0, \infty)$ 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 L.R.V. X :

- In distribution

$X_n \rightsquigarrow X$ (\Leftrightarrow) $\forall t$ where F_X cont: $\lim_{n \rightarrow \infty} F_{X_n}(t) = F_X(t)$ (pointwise convergence of dist. function)

- In probability

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

- Markov's Ineq. / Chebyshev's Ineq.

$$\forall \varepsilon > 0: P(X > \varepsilon) \leq \frac{E(X)}{\varepsilon}, \quad P(|X| > \varepsilon) \leq \frac{E(|X|)}{\varepsilon}$$

Problem Set week 3 $P(|X - E(X)| > \varepsilon) \leq \frac{V(X)}{\varepsilon^2}$

① Show that the estimator \bar{X}_n is consistent for θ in $\mathbb{R} \setminus \{0\}$ experiment

② Suppose the data X in a stat. experiment can take values in a countable set \mathcal{S} (i.e. $\text{Law}(X)$ is discrete).

In class it was claimed that the data itself is 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 set Π and $S: \Pi \rightarrow \mathcal{S}$ is one-to-one (bijective) mapping with values in \mathcal{S} (i.e. $\exists S^{-1}: \mathcal{S} \rightarrow \Pi$ s.t. $S^{-1}(S(t)) = t \quad \forall t \in \Pi$). Show that the statistic $S(T(X))$ is sufficient.

④ In class it was claimed that the statistic $T(X) = (\bar{X}_1, \bar{X}_n)$ is sufficient for the $\mathbb{R} \setminus \{0\}$ experiment. Prove this claim.

⑤ Let X_1, \dots, X_n be independent and identically distributed (i.i.d) with $\text{Law}(X_i) \sim P_0(\lambda), \lambda \in (0, \infty)$ is unknown. Show that the sample mean \bar{X}_n is a sufficient statistic for this data.

Weak Law of Large numbers

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

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

$$\text{Let } \varepsilon > 0, P(|\bar{X}_n - E(\bar{X}_n)| > \varepsilon) = \frac{V(\bar{X}_n)}{\varepsilon^2} = \frac{V(X_i)}{\varepsilon^2}.$$

Chub. law. $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n$.

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

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

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

Thus $\bar{X}_n \rightarrow E(\bar{X}_i)$ \square

Central Limit theorem (CLT)

If X_1, X_2, \dots, X_n and $E(X_i), V(X_i)$ 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) $\bar{X}_n - E(\bar{X}_n) \rightsquigarrow X - E(X_i) \sim \text{Normal}(0, \frac{V(X_i)}{n})$

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

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

iv) $\lim_{n \rightarrow \infty} P\left(\frac{\bar{X}_n - E(\bar{X}_n)}{\sqrt{V(\bar{X}_n)}} \leq z\right) = \lim_{n \rightarrow \infty} P(Z_n \leq z) = P(Z \leq z) = \phi(z)$

Approximation

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

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

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 θ in $\hat{\Theta}$ Bernoulli(θ) exp.

Sol.

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

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

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

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 θ . $B \subseteq \mathbb{X}$

$$P_\theta(X \in B | S(X) = s) = P_\theta(X \in B | X = s) = \left\{ \begin{array}{ll} \frac{P_\theta(\{X \in B\} \cap \{X = s\})}{P_\theta(X = s)}, & P_\theta(X = s) > 0 \\ 0 & , P_\theta(X = s) = 0 \end{array} \right.$$

$$P_\theta(\{X \in B\} \cap \{X = s\}) = \left\{ \begin{array}{ll} P_\theta(X = s), & s \in B \\ 0, & s \notin B \end{array} \right., \text{ so}$$

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

Hence

$$P_\theta(X \in B | T(X) = t) = P_\theta(X \in B | X = t) = \mathbf{1}_B(x) = \begin{cases} 1, & x \in B \\ 0, & x \notin B \end{cases} \quad \text{For all event } B \in \mathcal{F}_{\mathbb{X}}$$

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

Show that S is sufficient.

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

$$\frac{P_\theta(\{X \in B\} \cap \{S(T(X)) = s\})}{P_\theta(S(T(X)) = s)} = \frac{P(\{X \in B\} \cap \{T(X) = t\})}{P_\theta(T(X) = t)} \quad \text{where } t = S^{-1}(s)$$

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

□

Prof.

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 : $R(T, \theta) = E_\theta((\bar{X}_n - \theta)^2) \leq E_\theta((T(X) - \theta)^2) = R(T, \theta)$ 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 statistic 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

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

is defined as the expectation of $h(Y)$ under the cond. distribution of Y given $U(Y)=u$: $P(Y=y | U(Y)=u) = \frac{P(Y=y, U(Y)=u)}{P(U(Y)=u)}$

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

In special case $h(y) = 1_B(y)$ for some $B \subseteq \mathbb{Y}$, then

$$E(1_B(Y) | U(Y)=u) = \sum_{y \in B} P(Y=y | U(Y)=u) = P(Y \in B | U(Y)=u).$$

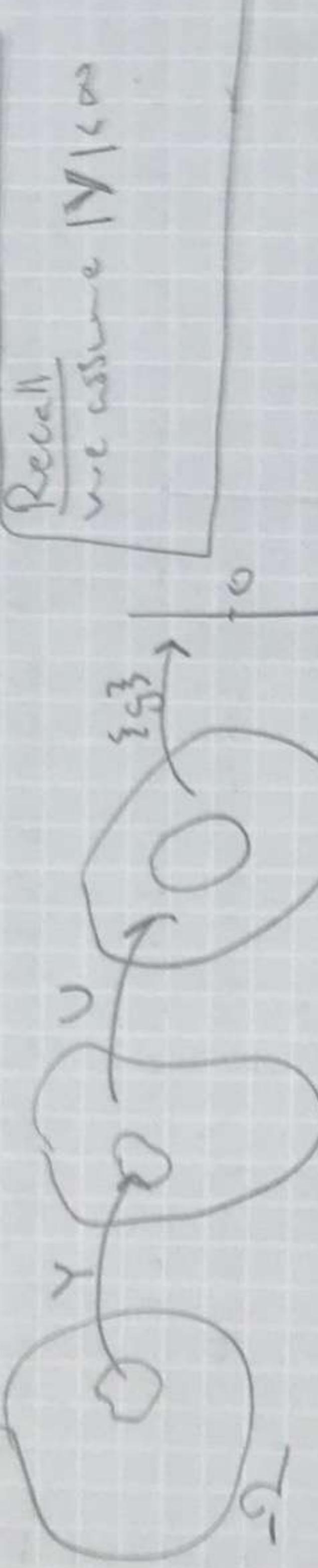
Def

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 $g: U \rightarrow \mathbb{R}$.



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

proof

exercise 0

i) 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 $1 = P(X=Y)$

ii) 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

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

exercise 1

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

Rough

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.

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

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 $\underbrace{E(H - \Pi(H))^2}_{E[(H - E(H|U))^2]} = \min_{z \in \mathcal{M}_U} E((H - z)^2)$

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" Π 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 θ in $\hat{\Theta}(\theta)$ 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 $\theta \in \Theta$ 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 θ

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

Before, we showed that $g(\bar{X}_n) = \bar{X}_n$ is the only unbiased
estimator of θ (\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 $Z \in \mathcal{N}_{\Theta}$,
 $Y = X$, $h(Y) = T$, $U = \bar{X}_n$, $E(h(Y) | U) = g(\bar{X}_n)$
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 estimate's variance:

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

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

uniform minimum variance unbiased estimator (UMVUE)

Bayes estimators

In the Bayesian approach to estimation of the param. θ 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 $\Theta \ni \theta$

e.g. In the $\text{Beta}(\alpha, \beta)$ 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 θ the prior $g(\theta)$ can be used to reduce the risk function $R(T, \theta)$ for each $\theta \in \Theta$ to a single number, by integration:

$$B_r(T) = R(T, \theta) g(\theta) d\theta = \int_0^1 R(T, \theta) g(\theta) d\theta$$

Suppose we have quadratic loss of T under risk R.

Def

A Bayes estimator T_B of θ 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

$$\begin{aligned} \arg \min_T B_r(T) &:= T \text{ s.t. } \\ B_r(T) &= \min_T B_r(T) \end{aligned}$$

The name "Bayesian" comes from Bayes formula:
 $\{D_1, \dots, D_k\}$ partition Ω , then $P(A) = \sum_{i=1}^k P(A|D_i)P(D_i)$

Motivation: Consider the case when P is the joint distribution of (X, U) where X is data and U is a RV. that takes k possible values in $\Theta := \{\theta_1, \dots, \theta_k\}$. Then $A = "X \in A"$ and $D_i := "U = \theta_i"$

Posterior prob of $\rightarrow P(U = \theta_i | X \in A) = \frac{P(X \in A | U = \theta_i)P(U = \theta_i)}{\sum_j P(X \in A | U = \theta_j)P(U = \theta_j)}$

The Bayesian approach views $\{P_\theta : \theta \in \Theta\}$ as a family of conditional distributions given θ , counts a prior dist.

In \mathbb{R}^n $\text{Be}(\theta)$ exp, consider the family of prior densities for $\theta \in \Theta$:

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

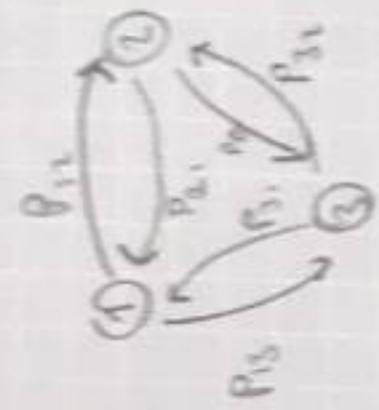
$$g_{\alpha, \beta}(\theta) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)}, \quad \theta \in [0, 1] \quad \text{where } B \text{ is beta function:}$$

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

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

Marrow Checks

If key "circ", then it implies



$$P = \begin{pmatrix} 0 & p_{12} & p_{13} \\ p_{21}, & 0 & p_{23} \\ p_{31}, & p_{32}, & 0 \end{pmatrix}$$

$$q_{ij} = q_{ji} p_{ij} \text{ if } i \neq j$$

$$P(+)=\alpha Q +$$

Note Details cont.

We consider the whole family $\mathcal{G}_{\alpha\beta}$, induced by $(\alpha, \beta) \in R^2_{>0}$ to allow a wide range of prior beliefs about Θ . It will become clear that the Bayesian approaches are very useful even to prove non-Bayesian prop. of estimators.

Prop. In $\bigotimes_{i=1}^n \mathcal{G}_i(\theta)$ wif, let g be an arbitrary prior density for $\theta \in [0, 1] = \Theta$, then the Bayes estimator is

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

Note

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

Proof

exercise

Def

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

$$\textcircled{2} : R(S, \theta) \leq R(\bar{\theta}, \theta) \quad \forall \theta \in \Theta \quad \text{implies}$$

$$\textcircled{3} : R(S, \theta) = R(\bar{\theta}, \theta)$$

Thus, admissibility of $\bar{\theta}$ means that there can be no estimator S which is uniformly at least as good [Θ holds] and strictly better for one $\theta_0 \in \Theta$, because then $R(S, \theta_0) < R(\bar{\theta}, \theta_0)$ and thus $\textcircled{3}$ is contradicted.

So, non-admissibility of $\bar{\theta}$ means that $\bar{\theta}$ can be improved by another estimator S .

Prop

Suppose that in $\textcircled{2}$ $\text{Be}(\theta)$ 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 $\bar{\theta}_B$ for this prior density is admissible for quadratic loss.

proof exercise

When trying these 2 exercises, show them to Rue2.

Q&A

Minimax Estimators

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

An estimator T_M is called Minimax if

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

some collection
of estimators

Prop

In $\bigotimes_{i=1}^n \mathcal{B}(\Theta)^{\text{exp}}$, the Bayes Estimator $T_{\alpha, p}$ for $\alpha = \beta = \frac{1}{2}$ is a minimax estimator

Proof exercise

Outline of course:

① Point Estimations

- Moment Estimations { solve sample moments = theoretical moments

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

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

→ CLT ↴

② Set Estimations

point est.

$\hat{\theta}_m$ set, $\hat{\theta}_n$ set, $\hat{\theta}_m$ prob. to contain θ_0 (true law)

③ Hypothesis testing

$\hat{f}_{y|x}$ Regression: $f_{y|x}$

Parametric Inference

Recall the D.T. setting!

Given a param. exp: $(\Omega, \mathcal{F}_n, \mathcal{P})$, $\mathcal{P} = \{\mathbb{P}_{\theta}; \theta \in \Theta, |\Theta| < \infty\}$
we are interested in some function $g(\theta)$ based on data x .
Typically, $x = (x_1, \dots, x_n)$, x_1, \dots, x_n has law θ is from \mathcal{P}
exp. x_1, \dots, x_n and $F_{\theta} = F_{\theta}(x_i; \theta)$ (or $f_{\theta} = f_{\theta}(x_i; \theta)$)
e.g. $x_1, \dots, x_n \sim \text{Beta}(\alpha)$

Ex1

$x_1, \dots, x_n \sim N(\mu, \sigma^2)$. Then $\Theta = \{\theta_1, \theta_2\} = \{(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0\}$
Suppose X_i is the outcome of a blood-alcohol breath-test and we
are interested in the fraction of testers with score > 1 .
What is $g(\theta)$? : $g(\theta) = P(X_i > 1) = 1 - P\left(\frac{X_i - \mu}{\sigma} < \frac{1 - \mu}{\sigma}\right) =$

$$= 1 - \phi\left(\frac{1 - \mu}{\sigma}\right)$$

Ex2 $X \sim \Gamma(\alpha, \beta)$, $\theta = (\theta_1, \theta_2) = (\alpha, \beta) \in \mathbb{R}_{>0}^2$, recall $f(x; \alpha, \beta) = \frac{1}{\beta^{\alpha} \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta}$.
Often X is used to model lifetime of items.

Then we are interested in an estimate of the mean lifetime:
 $g(\theta) = g(\alpha, \beta) = E(X; \alpha, \beta) = \alpha \beta$

Method of Moments Estimator (MOME)

Suppose the parameter θ has k components $(\theta_1(\theta), \dots, \theta_k(\theta))$, then for $1 \leq j \leq k$, define the j^{th} moment:

$$\alpha_j := \alpha_j(\theta) = E(X_j; \theta) = \int x^j dF_{\theta}(x)$$

and the sample moment:

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

Def MOME $\hat{\theta}_{\text{MOM}}(X) = \hat{\theta}_n$ is defined to be the value of θ s.t.

$$\begin{cases} \alpha_1(\hat{\theta}_n) = \hat{\alpha}_1 \\ \alpha_2(\hat{\theta}_n) = \hat{\alpha}_2 \\ \vdots \\ \alpha_k(\hat{\theta}_n) = \hat{\alpha}_k \end{cases}$$

i.e. $\hat{\theta}_n$ is the MOME estimate and is the solution θ 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 initialise 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 θ .

$$\begin{aligned} \text{sol. } \alpha_1(\hat{\theta}) &= E(X_1; \hat{\theta}) = \hat{\alpha}_1 \\ \hat{\alpha}_1 &= \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}_n \end{aligned} \quad \left. \begin{array}{l} \text{By equating we solve for } \hat{\theta}_n \\ \hat{\theta}_n = \bar{X}_n \end{array} \right\}$$

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

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

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

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

$$\left\{ \begin{array}{l} \hat{\mu}_n = \bar{X}_n \\ \hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \end{array} \right.$$

ex5

Show that MOME for parameter λ in a product $P_0(\lambda)$ esp

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

Example: homogeneous densities are possible (e.g. $\lambda = 1$)

Properties of MOME

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

- 1) $\hat{\theta}_n$ exists with $\text{prob} \rightarrow 1$ (equations are solvable for large n)
- 2) $\hat{\theta}_n \xrightarrow{P} \theta$ (asymptotic) consistency.
- 3) $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D} N(0, \Sigma)$, $\Sigma = \underbrace{g E_\theta(Y Y^\top)}_{\text{below}} g^\top$, where $Y = [X, X^\top, \dots, X^\top]^\top$,

$$g = (g_1, \dots, g_n), g_j = \frac{\partial}{\partial \theta} \alpha_j'(\theta)$$

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 $f(x_i; \theta) = \underbrace{f(x_i; \theta)}_{f_\theta(x_i)} = \underbrace{dF_\theta(x_i)}_{f_\theta(x_i)}$

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

$$L_n(\theta), \theta \rightarrow [0, \infty) =: R_{>0}$$

Def

log-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 θ random (because $x = (x_1, \dots, x_n)$ is random)

Def

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

Q1

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

$$\text{sol: } L_n(\theta) = \prod_{i=1}^n f(x_i; \theta) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} = \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} = \theta^{\sum x_i} (1-\theta)^{n-1}$$

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

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

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

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

$$\text{ex 2} \quad \text{Show that MLE } \hat{\lambda}_n \text{ of param } \lambda \text{ in a product Exp}(\lambda) \text{ if } \hat{\lambda}_n \text{ is } \frac{1}{\bar{x}_n}.$$

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 $\lambda \in (0, \infty)$

Q2 $X_1, \dots, X_n \sim N(\mu, \sigma^2)$. Find the MLE $\hat{\Theta}_n$ for the unknown parameter $\Theta = (\mu, \sigma^2)$

Sol: Likelihood function (ignoring constants)
for convenience
(as they are irr.)

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

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

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

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

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

$\hat{\theta}_{n,1} = \bar{X}_n$

$\hat{\theta}_{n,2} = S_n$

ex3 Tricky conceptually.

X_1, \dots, X_n iid Uniform(0, θ)

Find the MLE $\hat{\theta}_n$ of θ

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

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

The $f(x_i; \theta) = 0$ and thus $L_n(\theta) = 0$ if $x_i > \theta$

so, $L_n(\theta) = 0$ if $\max\{x_1, \dots, x_n\} > \theta$.

Now, consider any $\theta > x_{(n)}$. Then for each x_i we have

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

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

Properties of MLEs

- 1) The MLE is consistent: $\hat{\theta}_n \xrightarrow{P} \theta_0$, where θ_0 denote the true value of parameter θ : if the model is "good" or "reasonable".
- 2) The MLE is equivariant: if $\hat{\theta}_n$ is the MLE of θ , then $g(\hat{\theta}_n)$ is the MLE of $g(\theta)$ ("the function of interest")

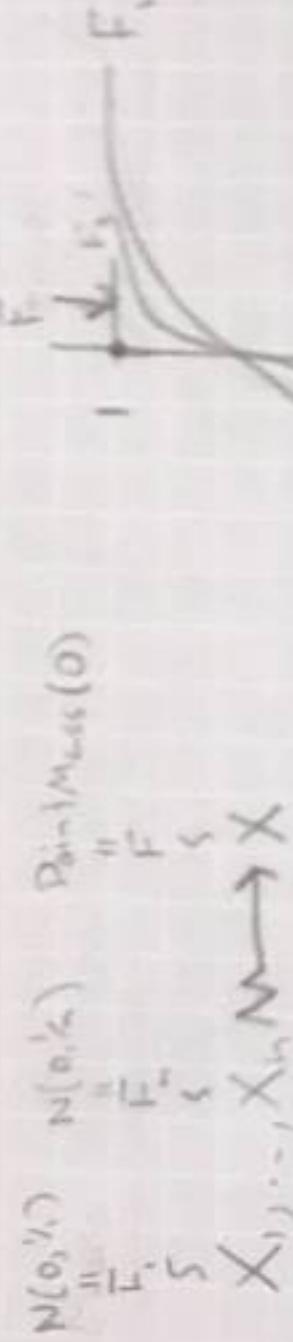
- 3) The MLE is asymptotically normal: $(\hat{\theta}_n - \theta_0) \xrightarrow{\sqrt{n}} N(0, \sigma_{\theta_0}^2)$
 $\sigma_{\theta_0}(\hat{\theta}_n) = \sqrt{\text{Var}(\hat{\theta}_n)}$ is the standard error

- 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 $f(x_i; \theta)$
 $f(x_1, \dots, x_n; \theta)$



Def. conv. in distribution

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

$$(X_n \xrightarrow{D} X)$$

$\lim_{n \rightarrow \infty} F_n(t) = F(t)$ (Point wise convergence of DFs)

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

Prop. ($L\bar{T}$)

Let $X_1, \dots, X_n \sim X_1$ and $E(X_1), V(X_1)$ exist.

$$\text{Then: } \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{D} X \sim N(E(X_1), \frac{V(X_1)}{n})$$

$$\therefore \bar{X}_n - E(X_1) \xrightarrow{D} \sqrt{n} (X - E(X_1)) \sim N(0, \frac{V(X_1)}{n})$$

$$\therefore \sqrt{n} (\bar{X}_n - E(X_1)) \xrightarrow{D} Z \sim N(0, 1)$$

$$\cdot Z_n = \frac{\sqrt{n}}{\sqrt{V(X_1)}} (\bar{X}_n - E(X_1)) = \frac{\bar{X}_n - E(\bar{X}_n)}{\sqrt{V(\bar{X}_n)}} \xrightarrow{D} Z \sim N(0, 1)$$

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

~~ex~~ Suppose collection X_1, \dots, X_n w/ $\#$ errors in program varred
 $\sim \text{Ind.}$. Suppose $X_i \sim \text{Po}(\lambda=5)$ and $X_i \text{ iid}$, $E(X_i) = 5$

Suppose w/ 125 and we want to make prob. statement of X_n
 $P(\bar{X}_n < 5.5) = P\left(\frac{\sum(X_{i,n} - E(X_{i,n}))}{\sqrt{V(X_{i,n})}} < \frac{\sqrt{5}(5.5 - E(X_{i,n}))}{\sqrt{V(X_{i,n})}}\right) = P(Z \leq \frac{\sqrt{5}(6.5 - 5)}{\sqrt{5}}) =$

$$= \phi(2.5)$$

Prob. back-fall

Def
Conv. in prob

$X_n \xrightarrow{P} X$ if $\forall \epsilon > 0$, $P(|X_n - X| > \epsilon) \xrightarrow{n \rightarrow \infty} 0$ ($X_n \xrightarrow{P} c$ if c is Point Mass(c))

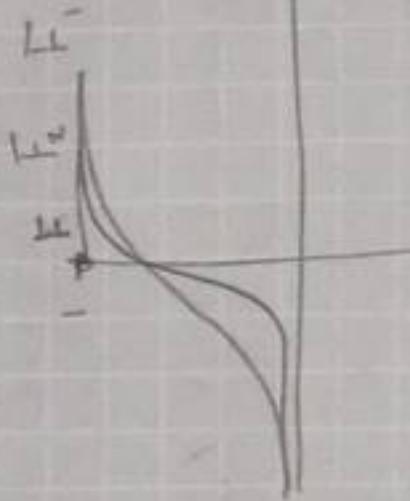
Def
Conv. in dist.

$X_n \rightsquigarrow X$ if $\forall t$: $F(t)$ 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 $E((X_n - X)^2) \xrightarrow{n \rightarrow \infty} 0$



$\text{ex} X_n \sim N(0, 1/n)$, $X \sim \text{Point Mass}(0)$.

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, $F_n(t) \rightarrow F(t)$ $\forall t \neq 0$ But this is ok since we only need to be for t where

$F(t)$ cont.
 conv. in prob. : $\forall \epsilon > 0$, $P(|X_n - 0| > \epsilon) = P(|X_n|^2 > \epsilon^2) = P((X_n)^2 > \epsilon^2) = \frac{E((X_n)^2)}{\epsilon^2} = \frac{E(X_n^2)}{\epsilon^2} = \frac{1}{\epsilon^2} \rightarrow 0$

Prop. (PNN)

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

b) $X_n \xrightarrow{P} X \Rightarrow X_n \xrightarrow{m} X$

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

So, conv. in q.m. \Rightarrow conv. in prob. \Rightarrow conv. in dist.

\curvearrowleft

if $X \sim \text{PointMass}(c)$

ex conv. in prob \neq conv. in q.m.

Let $U \sim U(0,1)$, $X_n = \sqrt{n} \mathbb{1}_{(0,1/n)}(U)$, $X_n \sim \text{PointMass}(0) = 0$

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

But $E(X_n^2) = \int_0^{1/n} (n u)^2 \, du = n \int_0^{1/n} u^2 \, du = n [u^3]_0^{1/n} = n \left(\frac{1}{n^3} - 0\right) = 1$ for all n so $X_n \not\xrightarrow{q.m.} X$

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"* $\Rightarrow \rightarrow \rightarrow$ but since $N(0,1)$ symmetric, still the same)

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

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

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

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

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

We typically don't know σ .

If we replace σ by its estimate $s_n := \sqrt{\frac{1}{n-1} \sum_{i=1}^{n-1} (X_i - \bar{X}_n)^2}$ it still works.

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

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

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

Back to properties of MLE

Prop (consistency)

Let θ_0 denote the true (and possibly unknown) value of θ .

$$M_n(\theta) := \frac{1}{n} \sum_{i=1}^n \ln \left(\frac{f(x_i; \theta)}{f(x_i; \theta_0)} \right), \quad M(\theta) = E_{\theta_0} \left(\ln \left(\frac{f(x_1; \theta)}{f(x_1; \theta_0)} \right) \right)$$

(maximizing (=) maximizing $\lambda_n(\theta)$)

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

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

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 θ and $\hat{\theta}_n$ be MLE of θ

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

Proof

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

$$\text{For any } \tau, L_\tau(\tau) = \prod_{i=1}^n f(x_i; h(\tau)) = \prod_{i=1}^n f(x_i; \hat{\theta}_n) = L_h(\hat{\theta}_n) \text{ where } \hat{\theta}_n = h(\tau).$$

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

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

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

Def
Fisher information is $I_n(\theta) = \text{Var}_{\theta} \left(\sum_{i=1}^n s(X_i; \theta) \right) = \sum_{i=1}^n I_{\theta}(X_i; \theta)$.
when $n=1$, $I(\theta) := I_1(\theta)$

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

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 thus construct normal-based asymptotic confidence intervals for θ .

Prop

$$\text{Let } C_n = [\bar{c}_n, \bar{c}_{\bar{n}}] = (\hat{\theta}_n - z_{\alpha/2} \hat{s}_{\theta_n}, \hat{\theta}_n + z_{\alpha/2} \hat{s}_{\theta_n}).$$

Theorem, $P_\theta(\theta \in C_n) \rightarrow 1 - \alpha$ as $n \rightarrow \infty$

Proof

Let ε be standard normal RV. Then,

$$P_\theta(\theta \in C_n) = P_\theta(\hat{\theta}_n - z_{\alpha/2} \hat{s}_{\theta_n} \leq \theta \leq \hat{\theta}_n + z_{\alpha/2} \hat{s}_{\theta_n}) = P_\theta(-z_{\alpha/2} \frac{\varepsilon}{\hat{s}_{\theta_n}} \leq \varepsilon \leq z_{\alpha/2}) \rightarrow P_{\varepsilon \sim N(0, 1)}(\varepsilon_{\alpha/2} \leq \varepsilon \leq -\varepsilon_{\alpha/2}) = 1 - \alpha$$

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

when you read an opinion poll in 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 $\hat{\theta}_n \pm 2\hat{s}_{\theta_n}$.

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

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

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

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

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

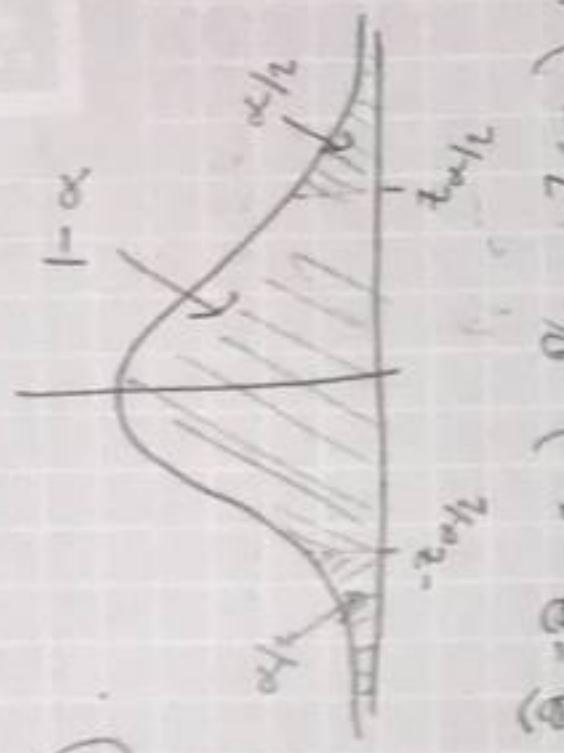
Ex Let $X_1, \dots, X_n \sim N(\theta, \sigma^2)$ where σ^2 is known.
 $s(x; \theta) = (x - \theta)/\sigma$ and $s'(x; \theta) = \frac{-1}{\sigma^2}$, so $I(\theta) = \frac{1}{\sigma^2}$.

The MLE of θ is $\hat{\theta}_n = \bar{X}_n$. $\bar{X}_n \sim N(\theta, \sigma^2/n)$.

In this case the Normal approx. is actually exact.

Exercise

Find 95% conf. int. for X_1, \dots, X_n and $P_\theta(\Lambda)$ and MLE



Optimality of MLE:

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

Another reasonable estimator is sample median $\tilde{\theta}_{n,1}$.
We know that MLE satisfies $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D} N(0, \sigma^2)$

We know that MLE satisfies $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D} N(0, \sigma^2)$
One can show that $\tilde{\theta}_{n,1}$ satisfies $\sqrt{n}(\tilde{\theta}_{n,1} - \theta) \xrightarrow{D} N(0, \sigma^2 \frac{1}{2})$
So, the median converges to the right value θ 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, $ARE(\tilde{\theta}_{n,1}, \hat{\theta}_n) = \frac{2}{1} = 0.63$, in words: if you use the median-based estimator $\tilde{\theta}_{n,1}$, then you are effectively using only a fraction of the samples (0.63 of).

Prop.

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

All of this MLE stuff is only as good as the Θ -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 $\hat{\tau}_n$?

Prop (The Delta Method)

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

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

Hence, if $C_n = (\hat{\tau}_n - \bar{x}_{\text{obs}}, \hat{\tau}_n - \bar{x}_{\text{obs}} \cdot \hat{se}(\hat{\tau}_n))$, then

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

ex Let $X_1, \dots, X_n \sim \text{Ber}(\theta)$ and let $\psi = g(\theta) = \lambda_n(\frac{\theta}{1-\theta})$

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

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

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

Exercise

$X_1, \dots, X_n \sim N(\mu, \sigma^2)$. Suppose μ is known but not $\sigma^2 > 0$.

We want to estimate $\psi = \lambda_n(\sigma)$.

hints

- Find MLE (log-likelihood etc.)

$$\hat{\psi} = \lambda_n(\hat{\sigma}_n)$$

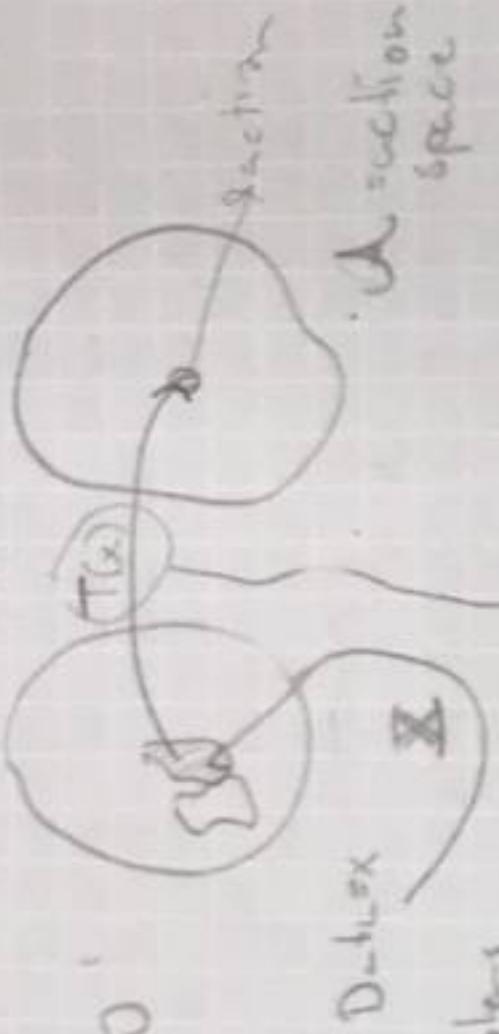
$$\text{find } \hat{se}(\hat{\sigma})$$

$$\text{find } \hat{se}(\hat{\psi})$$

$$\frac{\hat{\psi} - \psi}{\hat{se}(\hat{\psi})} = \frac{\sqrt{\lambda_n(\hat{\sigma}_n)}}{\sqrt{\lambda_n(\sigma)}} \sim N(0, 1)$$

$$\text{95\% ci. : } \hat{\psi} \pm \frac{2}{\sqrt{\lambda_n(\sigma)}}$$

Generally in decision theory:



Decision problems

estimation problem

$\hat{\mu} = \Theta$, T is $\hat{\Theta}_n$ (estimator)

$\sqrt{\Theta}|_{\infty}$

parametric

$X_1, \dots, X_n \sim L_\theta$

Point estimators:
• MOME
• MLE $\hat{\Theta}_n$

$\hat{F}_n \rightarrow F_n$ (Givens-Law)

Set

estimators: $\{C_n, \bar{C}_n, C_n^*, \text{conf. interval}\}$

$\cdot [\bar{F}_n, \hat{F}_n]$ (interval)
(in function space)

(Dvoretzky-Kiefer-Wolfowitz)
(ineq.)

Another Decision Problem: (hypotheses testing)

$\mathcal{A} = \{\emptyset, \Omega\}$, β , T is Test statistic

Rejection region
(Rejection region)^c



Karl Popper
"A scientific hyp.
is one that is
falsifiable"

H_0 - null hypothesis
(or H_n)

H_A - alternative
hypothesis

false pos.

(fail to reject)

Retain Null | Reject Null

Type I error |

Type II error |

$1 - \beta$ |

Power

Hyp. testing Problem

Should H_0 be rejected?

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

So, the summing of outcome of h.y. testing:

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

ex) The Null hypothesis: "The disease rate is the same in two groups".
Alt. hypothesis: "The disease rate is not the same".

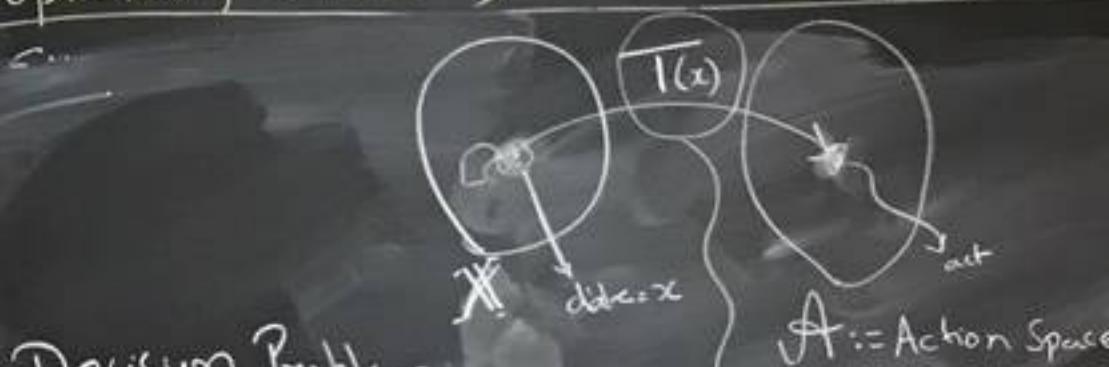
Formally

We partition index set into two disjoint sets Θ_0, Θ_1 ,
and we want to tell
 $H_0: \theta \in \Theta_0$ versus $H_1: \theta \in \Theta_1$
↑
Null hyp. Alt. hyp.
 $\{\theta_0\}$ $\{\theta_1\}$

Def. Rejection region

$Q \subset \mathcal{X}$ s.t. if $x \in Q$, we reject H_0 .

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



Decision Problems

Estimation Problem

$$A = \hat{\Theta}, T = \hat{\theta}_n$$

parametric / non-parametric

Point Estimator

MOME

MLE $\hat{\theta}_n$

Set Estimators

$$[\underline{c}_n, \bar{c}_n], C_n$$

Another Decision Problem is Hypothesis Testing Problem.

$A = \{0, 1\}$, T is called Test Statistic

Hypotheses

- Null Hypothesis H_0
- Alternative Hypothesis H_1

Rejection Region

Acceptance Region

Statistical Posable Science

Karl Popper

Falsifiability

Science? vs. Nonsense?

Demarcation Problem

Hyp Testing Problem is A Scientific hypothesis is one that is falsifiable

Should H_0 be rejected?

ex The Null hypothesis: "The disease rate is the same in two groups".
 Alt. hypothesis: "The disease rate is not the same".

For clarity

We partition index sets into two disjoint sets Θ_0, Θ_1 and we want to test

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

\uparrow

Null hyp.

$$\Theta = \{\bar{\Theta}_0\} \cup \{\bar{\Theta}_1\}$$

Def Rejection region

$R \subset \mathbb{X}$ 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	H_1 true	Type I error
H_0 false ($\Rightarrow H_1$ true)		Type II error	✓

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 h.p. testing is to find appropriate T and c .

Q2

The power function of a test with rejection region Ω is

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

The size of a test is $\alpha := \int_{\Omega} \phi(\theta) d\theta$

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

(Types) of hypothesis:

- Simple hypothesis: $H_0: \theta = \theta_0, H_1: \theta \neq \theta_0$
- One-sided test: $H_0: \theta \leq \theta_0, H_1: \theta > \theta_0$ or $H_0: \theta \geq \theta_0, H_1: \theta < \theta_0$
- Composite hypothesis: $H_0: \theta \in \Theta_0, H_1: \theta \notin \Theta_0$

Types of tests:

- Two-sided tests: $H_0: \theta = \theta_0, H_1: \theta \neq \theta_0$
- 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 $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, where σ 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 $\Omega = \{\bar{X}_1, \dots, \bar{X}_n : T(\bar{X}_1, \dots, \bar{X}_n) = \frac{1}{n} \sum_{i=1}^n x_i > c\}$

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

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

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

So, we reject the null hypothesis if $\bar{X}_n > \Phi^{-1}(1-\alpha) \frac{\sqrt{n}}{\sigma}$ for a size α 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 $H_1(\theta \in H_1)$ among all size α tests. Such a test (if it exists) is called the most powerful size α test. (Generally very hard to find)

A common test is the so-called Wald Test:

Let $\theta \in \mathbb{R}$, $\hat{\theta}_n$ be an estimator of θ with realization (estimator) $\hat{\theta}_n$, and let \hat{s}_{θ_n} be the estimated standard error of $\hat{\theta}_n$.

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 α 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 α , i.e.

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

Proof

Under $\underbrace{H_0: \theta = \theta_0}_{(\text{if } H_0 \text{ is true})}$, $\frac{(\hat{\theta}_n - \theta_0)}{\hat{s}_{\theta_n}} \xrightarrow{D} N(0, 1)$ (by assumption), hence

the probability of rejecting H_0 when H_0 is true is $P_{\theta_0}(|W| > z_{\alpha/2}) = P_{\theta_0}\left(\left|\frac{\hat{\theta}_n - \theta_0}{\hat{s}_{\theta_n}}\right| > z_{\alpha/2}\right) \rightarrow \lim_{n \rightarrow \infty} P\left(\left|\frac{\hat{\theta}_n - \theta_0}{\hat{s}_{\theta_n}}\right| > z_{\alpha/2}\right) = \alpha$

Remark

Alternatively, $W = \frac{(\hat{\theta}_n - \theta_0)}{\hat{s}_{\theta_n}}$ where \hat{s}_{θ_n} is the standard error computed with $\theta = \theta_0$. Directly, So, one can use \hat{s}_{θ_n} or s_{θ_0} 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 θ 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- α asympt. conf. interval $\hat{\theta}_n \pm z_{\alpha/2} s_{\hat{\theta}_n}$

The size α Wald test rejects $H_0: \theta = \theta_0$ vs. $H_1: \theta > \theta_0$, iff $\theta_0 + C_n - [\bar{C}_n, \bar{C}_n] \in [\hat{\theta}_n - z_{\alpha/2} s_{\hat{\theta}_n}, \hat{\theta}_n + z_{\alpha/2} s_{\hat{\theta}_n}]$.

Thus, testing the hypothesis is equivalent to checking if the null value θ_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, and the second Y $\frac{X}{n} \leq Y \leq \frac{X+1}{n}$

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

sol Let $X \sim \text{Bin}(n, \theta)$, $Y \sim \text{Bin}(n, \hat{\theta})$. We want to test $H_0: \theta = \hat{\theta}$ or $\delta = 0$ where $\delta := \theta - \hat{\theta}$. So, $H_0: \delta = 0$, $H_1: \delta \neq 0$ to conduct Wald test.

The MLE of δ is $\hat{\delta} = \hat{\delta}_{\min} = \hat{\theta}_n - \hat{\theta}_{\hat{n}}$ by equivariance prop of MLE.

The estimated std. error of $\hat{\delta} - \hat{\delta}_0 - \hat{\theta}_n$ is $\hat{s}_{\hat{\delta}_n} = \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} + \frac{\hat{\theta}_{\hat{n}}(1-\hat{\theta}_{\hat{n}})}{n}$

The size α Wald test is to reject H_0 when $|W| > z_{\alpha/2}$, where $W = \frac{\hat{\delta} - 0}{\hat{s}_{\hat{\delta}_n}} = \frac{\hat{\theta}_n - \hat{\theta}_{\hat{n}}}{\hat{s}_{\hat{\delta}_n}}$.

Suppose you observe 18 bad prediction out of 117 trials for alg. 1 and 21 for alg. 2.

Does the size $\alpha = 0.05$ Wald test reject the null hypothesis that $\delta = 0$

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 λ , s_n of $\hat{\lambda}_n$, $(1-\alpha)$ c.i. for λ . Then give a size $\alpha = 0.05$ Wald Test to reject $H_0: \lambda = \lambda_0$

P-values

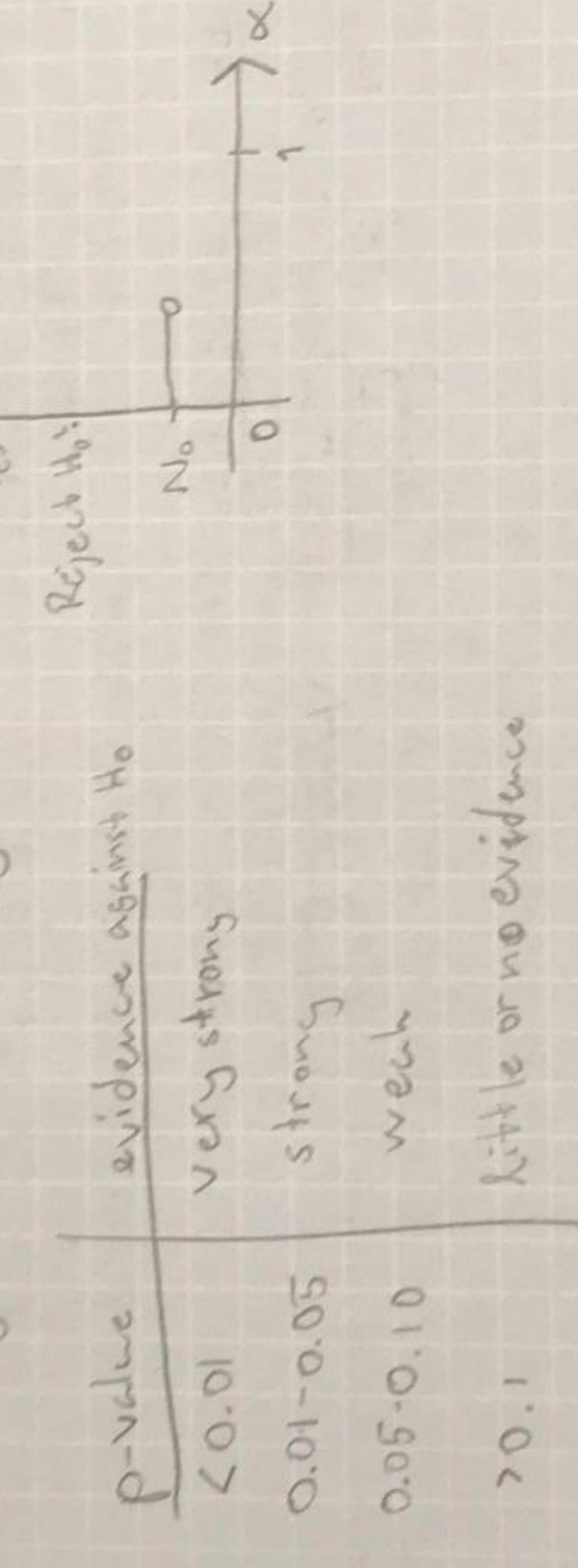
Instead of just reporting "reject H_0 " or "fail to reject H_0 " at a given size α , we could make a more informed report of the test. If a test rejects at level α , it will also reject at α' s.

But there will be a smallest α at which the test rejects H_0 .

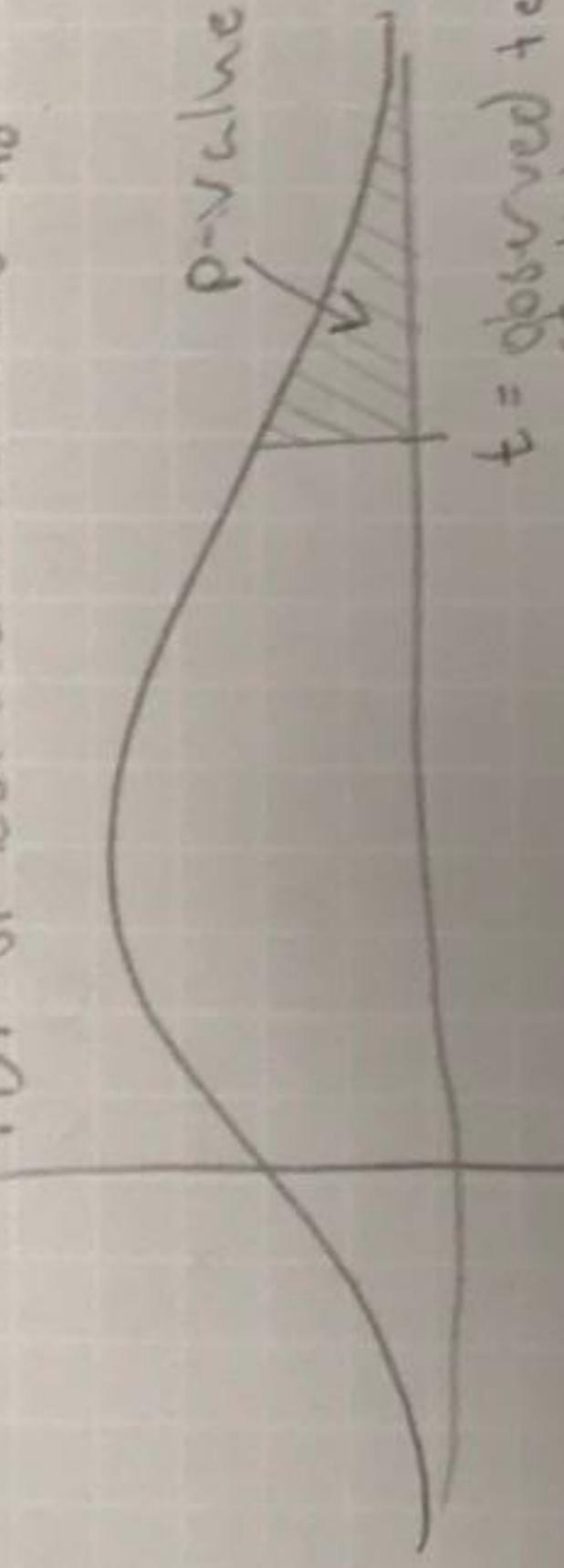
Def Suppose for every $\alpha_{(0)}$, we have a test of size α with rejection region R_α . Then, p-value: $\inf_{\alpha_{(0)}} \{\alpha : T(x_1, \dots, x_n) \in R_\alpha\}$. i.e. p-value is the smallest α at which we reject H_0 .

Remarks

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



↑ PDF of test statistic under H_0



To find p-value for a particular test statistic T , we find α s.t. the observed statistic is just at the boundary of rejection region R
 $t = \text{observed test statistic}$
 \Rightarrow p-value is tail area $P(T > t)$
 (if T is of the form $T > c$)

Pearson's χ^2 test for multinomial data

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

$$\sum_{k=1}^3 = \sum_{k=1}^3 k = n$$

\Rightarrow binomial

$$P((X_1, \dots, X_k) = (x_1, \dots, x_k); n, \theta_1, \theta_2, \dots, \theta_k) = \binom{n}{x_1, x_2, \dots, x_k} \theta_1^{x_1} \theta_2^{x_2} \dots \theta_k^{x_k} = \frac{n!}{x_1! x_2! \dots x_k!} \theta_1^{x_1} \theta_2^{x_2} \dots \theta_k^{x_k}$$

Multinomial:

$$2 \leq k < \infty$$

$$P((X_1, \dots, X_k) = (x_1, \dots, x_k); n, \theta_1, \theta_2, \dots, \theta_k) = \binom{n}{x_1, x_2, \dots, x_k} \theta_1^{x_1} \theta_2^{x_2} \dots \theta_k^{x_k} = \frac{n!}{x_1! x_2! \dots x_k!} \theta_1^{x_1} \theta_2^{x_2} \dots \theta_k^{x_k}$$

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

$$\text{We can check that the MLE of } \theta = (\theta_1, \dots, \theta_k), \hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k) = \left(\frac{x_1}{n}, \frac{x_2}{n}, \dots, \frac{x_k}{n} \right)$$

Let $\theta_0 = (\theta_{01}, \dots, \theta_{0k})$ be some fixed vector and suppose we want to test

$$H_0: \theta = \theta_0 \text{ vs. } H_1: \theta \neq \theta_0$$

$$\text{Def Pearson } \chi^2 \text{-statistic is}$$

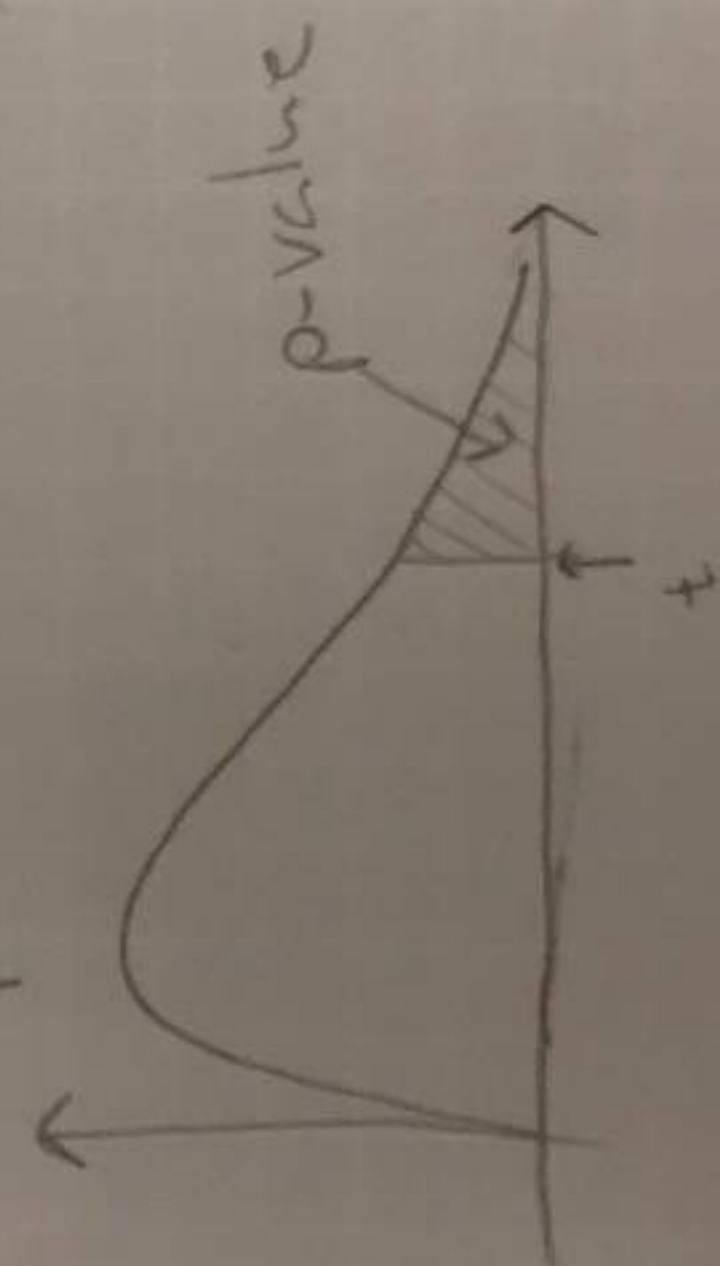
$$T = \sum_{j=1}^k \frac{(x_j - n\theta_{0j})^2}{n\theta_{0j}} = \sum_{j=1}^k \frac{(x_j - E_{\theta_0}(x_j))^2}{E_{\theta_0}(x_j)}$$

x_j under H_0 .

Prop.

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

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

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

-0.47.

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

The χ^2 -distribution (Prelude to Pearson's χ^2 -test.)

Def Let $Z_1, \dots, Z_k \stackrel{iid}{\sim} N(0, 1)$ and $Y = \sum_{i=1}^k Z_i^2$. Then Y has χ^2 distr.

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

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

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

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

Permutation test

This is a non-parametric test (in ANB testing context).

Say we want to test if two sets of samples come from the same distribution or not. $X_1, \dots, X_m \sim F_X \in \{\text{all DFs}\}$ and $Y_1, \dots, Y_n \sim F_Y \in \{\text{all DFs}\}$

$H_0: F_X = F_Y \text{ vs } H_1: F_X \neq F_Y$

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 T and denote them T_1, \dots, T_N .

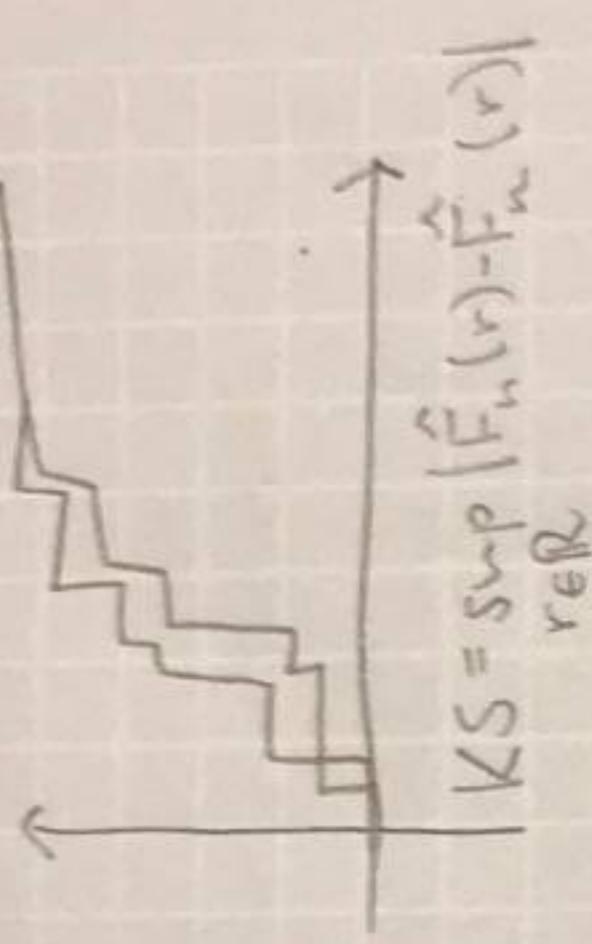
Under H_0 , all the T_i 's have equal prob. So the distribution

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

This P_0 is called the permutation distribution of T .

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

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



car (Toy example)

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

Let $T(X_1, X_2, Y) = |X_1 - \bar{X}_2| = 2$

Permutation	T	Prob under H_0
$(1, 9, 3)$	$2 = t_{\text{obs}}$	$\frac{1}{6} = \frac{1}{3!}$
$(4, 1, 3)$	2	$\frac{1}{6}$
$(1, 3, 9)$	7	$\frac{1}{6}$
$(3, 1, 9)$	7	$\frac{1}{6}$
$(3, 9, 1)$	5	$\frac{1}{6}$
$(9, 3, 1)$	5	$\frac{1}{6}$

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 LN.

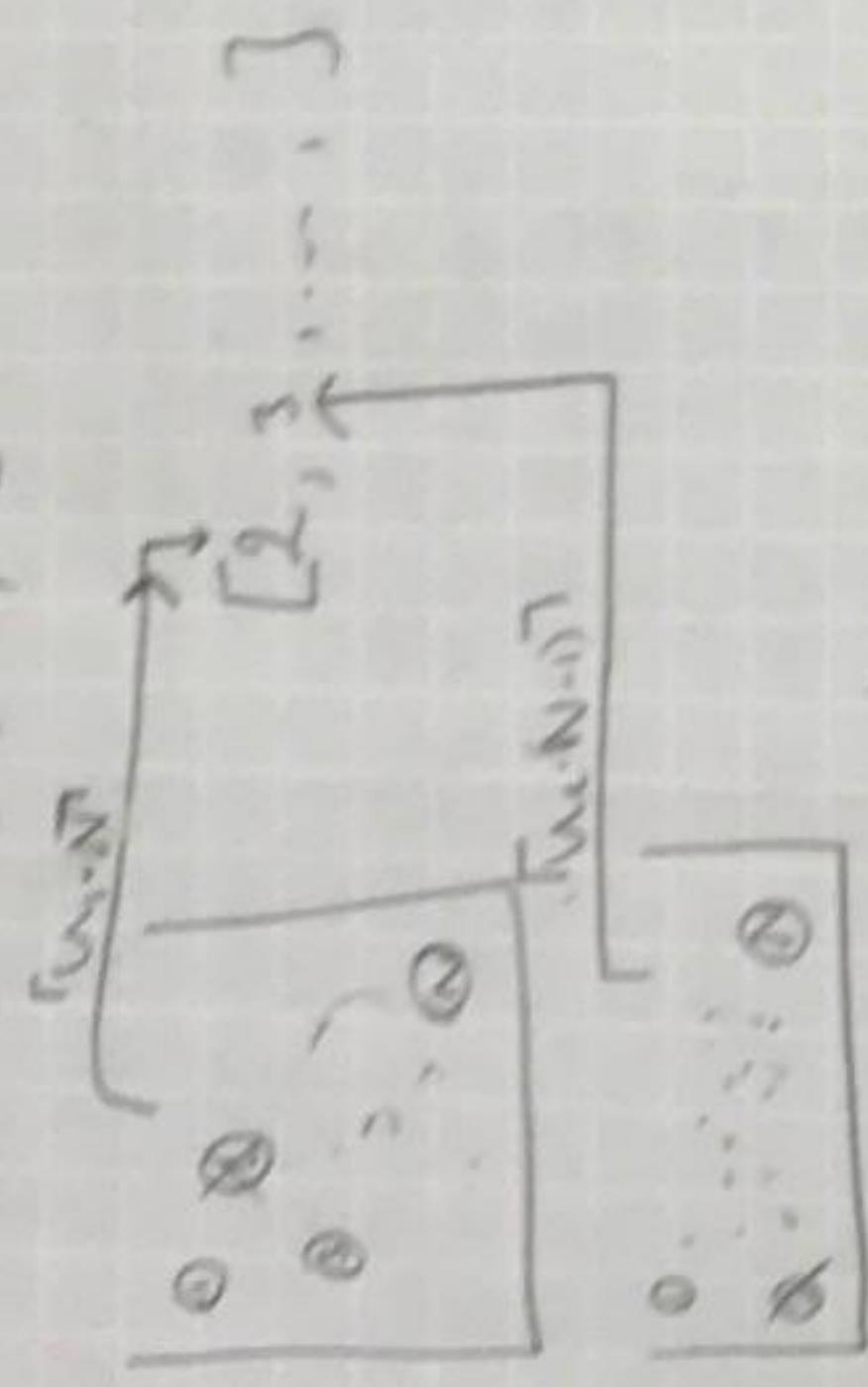
Alg. for Permutation test

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:

$$u_{1:n}, u_{2:n}, \dots, u_{n:n} \sim \text{Uniform}(0,1)$$

independent $\sim U(1, 2, \dots, N)$



Likelihood Ratio Test Statistic (LRTS)

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

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:
$$\Lambda_n := \Lambda_n := 2 \ln \left(\frac{\sup_{\theta \in \Theta} L_n(\theta)}{\sup_{\theta \in \Theta_0} L_n(\theta)} \right) = 2 \ln \left(\frac{L_n(\hat{\theta}_n)}{L_n(\hat{\theta}_{0,n})} \right)$$
 where $\hat{\theta}_n$ is MLE and
$$\hat{\theta}_{0,n}$$
 is MLE for θ restricted to Θ_0 .

Prop.

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

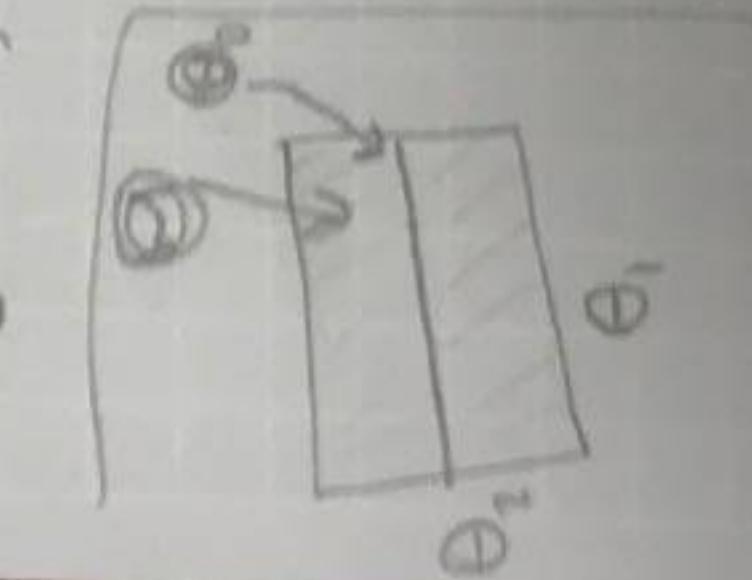
$$|\Theta| = r, |\Theta_0| = q$$

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

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

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

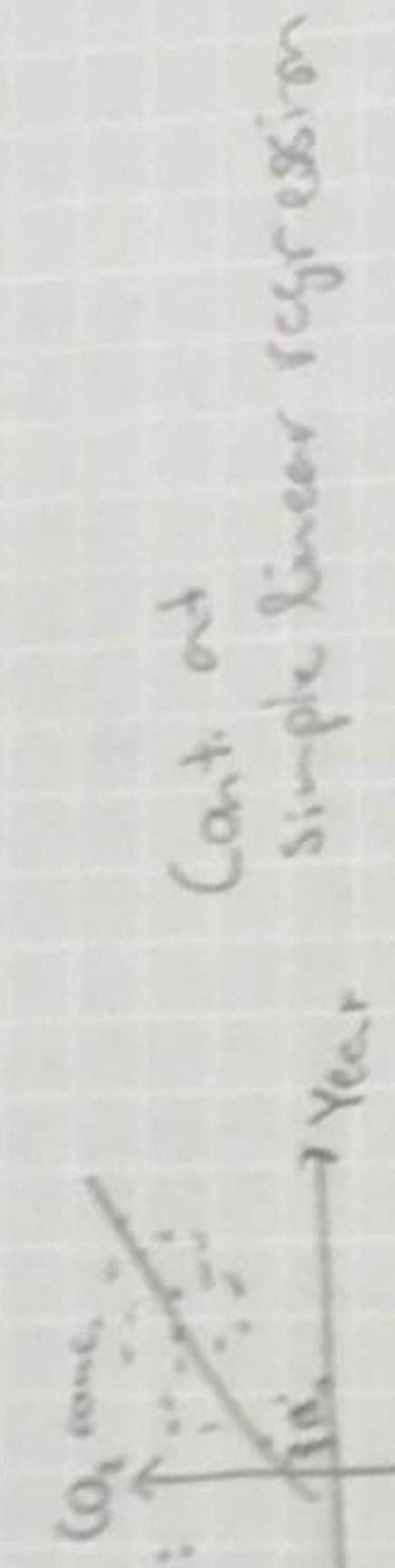
↑ observed LRTS



For instance, if $\theta = (\theta_0, \dots, \theta_6)$ and we want to test $\theta_4 = \theta_5 = \theta_6 = 0$, then Λ_θ has limiting distn given by $X_{i,3}^T X_{i,4}^T X_{i,5}^T X_{i,6}^T Q_V$.

Consider the 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

Sugimura notebook from day 1:


Non-parametric Estimation

Ques: Let $X_1, \dots, X_n \stackrel{iid}{\sim} F$ DF, $\{F\}$

$$\text{Def: } \hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(X_i \leq x) : \quad \begin{array}{c} \uparrow \\ \vdots \\ \uparrow \\ x_{(1)} \ x_{(2)} \ x_{(3)} \end{array}$$

Prop

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

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

$$\text{So, MSE (mean squared error)} = \frac{E(x)(1-F(x))}{\hat{F}_n(x) \xrightarrow{n} F(x)} \rightarrow 0$$

Prop. (Givensko-Cantelli lemma)

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

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

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

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

Drost, Dikw #, for any fixed but unknown $F \in \{all DFs\}$,

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

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: $r(x) := E(Y|X=x) = \int y f(y|x) dy$.

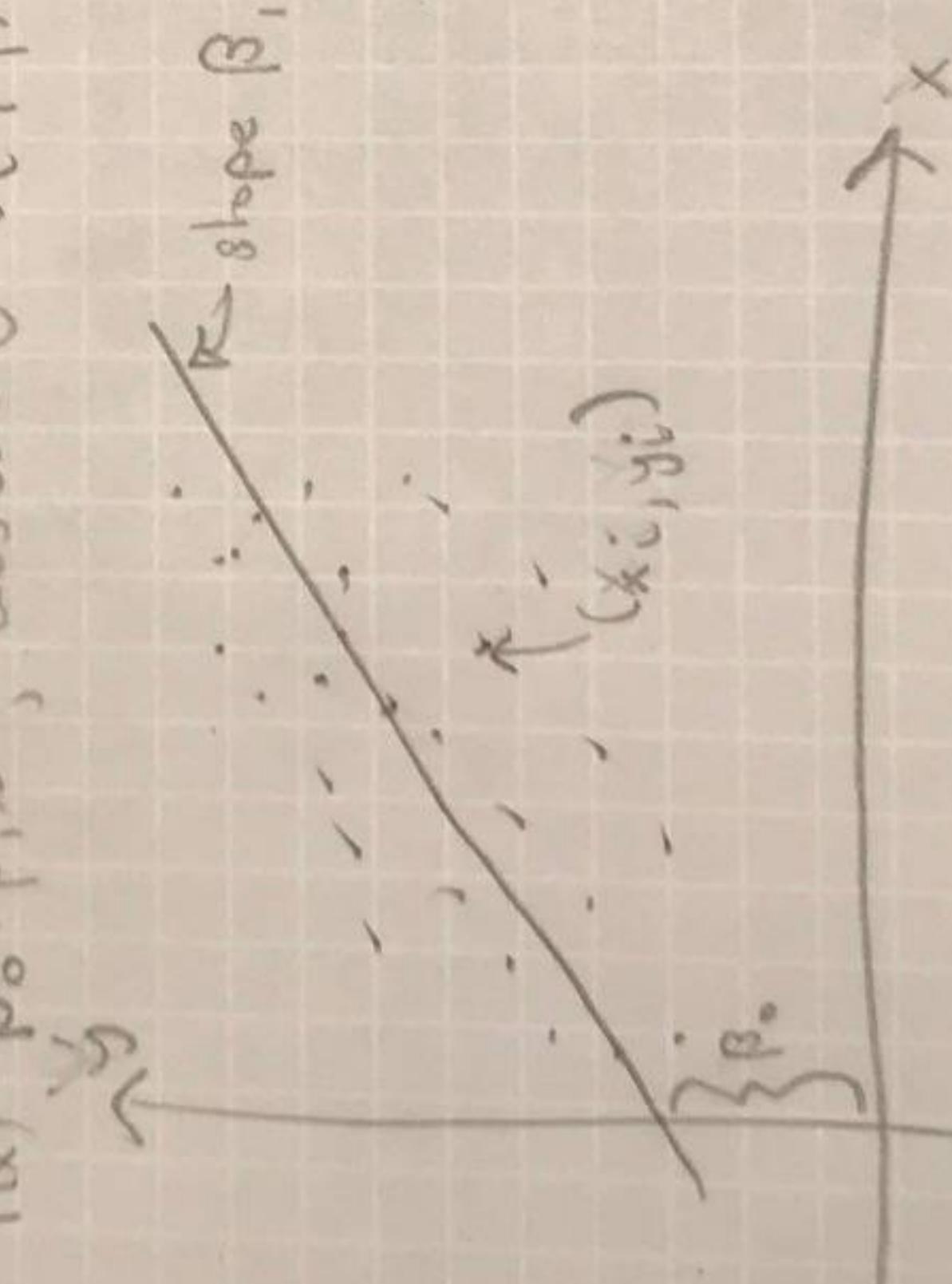
Our goal is to estimate $r(x)$ from the data of the form

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

Simplest version is called simple linear regression:

X_i, Y_i : real-valued.

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



① Def SLR Model

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

We want estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ for β_0, β_1 from the data.

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

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

$\hat{y}_i = \hat{r}(x_i)$ 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.

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x}_n)(y_i - \bar{y}_n)}{\sum (x_i - \bar{x}_n)^2}, \quad \hat{\beta}_0 = \bar{y}_n - \hat{\beta}_1 \bar{x}_n$$

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

Prop.

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