

INTRODUCTION TO STATISTICS

LECTURER: PEDRO IZQUIERDO LEHMANN
pizquier1@jh.edu

LECTURE 1

TODAY

- WHAT IS STATISTICS?
- STATISTICAL MODEL
- SAMPLE
- ASSUMPTIONS
- FINITE SAMPLE VS ASYMPTOTIC
- OUTLINE OF 6 LECTURES.

WHAT IS STATISTICS?

- STATISTICS \neq MATH.

IT'S A DISCIPLINE THAT USES MATH TO BUILD MODELS OF REALITY. AND DRAW CONCLUSIONS BASED ON EXPERIENCE.

- MAIN ASSUMPTION: SOME PHENOMENON OF INTEREST IS GENERATED BY SOME PROBABILITY LAW. $\rightarrow f_0$

EX 0: THROWING WEIGHTED COINS:

$$f_0 = \begin{cases} \text{HEADS} = 1 & \text{WITH PROB } p_0 \\ \text{TAILS} = 0 & \text{WITH PROB } 1-p_0 \end{cases}$$

$$p_0 \in (0, 1)$$

WE DON'T KNOW f_0 ! BUT..

I) WE KNOW (WE ASSUME) SOME STRUCTURE OF f_0 :

$f_0 \in \mathcal{F}$ (\mathcal{F} IS A STRUCTURED SET OF DISTRIBUTIONS)

\mathcal{F} IS CALLED A STATISTICAL MODEL.

EX: TOSSING COINS:

$f_0 \in \mathcal{F} = \left\{ f = \begin{cases} 1 & \text{WITH PROB } p \\ 0 & \text{WITH PROB } 1-p \end{cases} \mid p \in (0,1) \right\}$

II) WE HAVE A SAMPLE OF OBSERVATION OF THE PHENOMENON OF INTEREST:

x_1, \dots, x_n (DATA POINTS)

GIVEN I), THESE DATA POINTS ARE THE REALIZATION OF i.i.d RANDOM VARIABLES

$x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} f_0$ THIS CAN BE RELAXED!

GOAL : CHARACTERIZE f_0 FROM

I) $f_0 \in \mathcal{F}$ (STATISTICAL MODEL)

II) $X_1, \dots, X_n \sim f_0$ (SAMPLES)

Ex: TOSSING COINS

- ESTIMATE p_0 , e.g.

$$\hat{p}_0 = \frac{1}{n} \sum_{k=1}^n x_k$$

- CONFIDENCE INTERVALS
- HYPOTHESIS TESTING : CAN WE CONFIDENTLY SAY THAT $p_0 \in S \subseteq (0, 1)$.

STRUCTURE (ASSUMPTIONS)

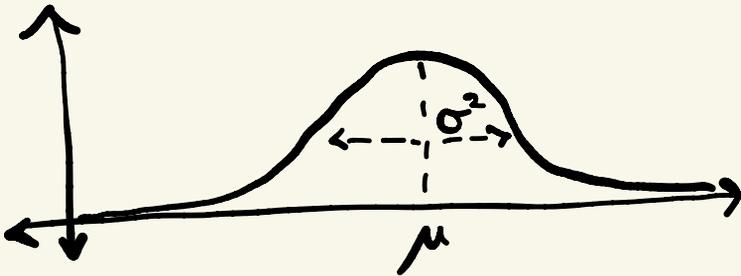
I) STATISTICAL MODEL

$$a) \mathcal{F} = \{ f_{\theta} : \theta \in \Theta \subseteq \mathbb{R}^d \}$$

PARAMETRIC STATISTICAL MODEL.

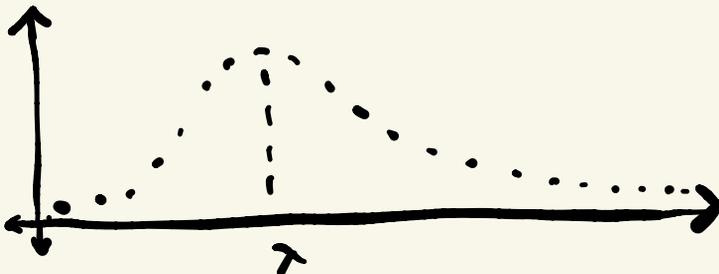
- POPULATION ATTRIBUTE (e.g. HEIGHT)

$$f_{\theta} \in \mathcal{F} = \{ \mathcal{N}(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0 \}$$



- NUMBER OF GOALS IN A SOCCER MATCH

$$f_{\theta} \in \mathcal{F} = \{ \text{POISSON}(\lambda) : \lambda \in (0, \infty) \}$$



b) NON PARAMETRIC MODELS

\mathcal{F} IS SUCH THAT $\dim(\mathcal{F}) = \infty$.

EX : SOME POPULATION ATTRIBUTE FOR WHICH WE HAVE **LITTLE INTUITION**,

$f_0 \in \mathcal{F}$ { ALL THE K-LIPSCHITZ CONTINUOUS DISTRIBUTIONS }

II) SAMPLES.

- $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} f_0$

INDEPENDENT AND IDENTICALLY DISTRIBUTED STANDARD BUT STRONG!

- $X_1, \dots, X_n \sim f_0$ **IDENTICALLY DISTRIBUTED AND EXCHANGABLE**

FOR ANY INDEX PERMUTATION $\sigma: [n] \rightarrow [n]$:

$(X_{\sigma(1)}, \dots, X_{\sigma(n)})$ HAS THE SAME DISTRIBUTION **INCREASINGLY POPULAR!**

- $(X_1, \dots, X_n) \sim f_0 \rightarrow$ **VECTOR VALUED DISTRIBUTION**
WEAK BUT NOT VERY VERSATILE!

PAY ATTENTION TO OUR ASSUMPTIONS!

- ASSUMPTIONS DON'T HOLD:
 - MATHEMATICAL MODEL DOES NOT REPRESENT OUR PHENOMENON OF INTEREST
 - CONCLUSIONS ARE DERIVED FROM MEANING.
- ASSUMPTION HOLDS APPROXIMATELY
 - THE ABOVE ISSUE IS JUST APPROXIMATELY PROBLEMATIC
- WE NEED ASSUMPTIONS TO DO MATH.

TRADE OFF!

FINITE SAMPLE VS ASYMPTOTIC

- IDEALLY WE WANT FINITE SAMPLE RESULTS:

$$\rightarrow n < \infty.$$

(BUT REAL SAMPLE IS FINITE!)

EX: HOEFFDING BOUND: IF

$X_1, \dots, X_n \stackrel{i.i.d.}{\sim} f_0$, $0 \leq X_i \leq 1$; THEN

$$\forall \epsilon: \mathbb{P}\left(\left|\frac{1}{n} \sum_{k=1}^n X_k - \mathbb{E}(X_1)\right| \geq \epsilon\right) \leq e^{-2n\epsilon^2}$$

- WE MIGHT NOT HAVE SUFFICIENT \mathcal{F} STRUCTURE.
- MATH CAN BE TOO HARD / INTRACTABLE.

- WE CAN STUDY ASYMPTOTICS!

$$\rightarrow n = \infty, \quad n \rightarrow \infty$$

(NO REAL SAMPLE IS INFINITE)

EX (WEAK) LAW OF LARGE NUMBERS

$$\forall \epsilon: \mathbb{P}\left(\left|\frac{1}{n} \sum_{k=1}^n X_k - \mathbb{E}(X_1)\right| \geq \epsilon\right) \xrightarrow{n \rightarrow \infty} 0$$

LIFE (MATH) IS EASIER!

OUTLINE FOR NEXT LECTURES

- 1) WHAT'S STATISTICS?
- 2) STATISTICS & ESTIMATORS
- 3) CONFIDENCE INTERVALS
- 4) HYPOTHESIS TESTING
- 5) FINITE SAMPLE & ASYMPTOTIC PARAMETRIC METHODS
- 6) LINEAR REGRESSION

TOPICS WE ARE NOT INTRODUCING BUT YOU MIGHT LIKE

- NON PARAMETRIC STATS
- HIGH DIMENSIONAL STATS
- BAYESIAN STATS
- STATISTICAL LEARNING



LECTURE 2

LAST TIME

- WHAT IS STATISTICS?
- SAMPLE
- STATISTICAL MODEL
- ASSUMPTIONS
- FINITE SAMPLE v/s ASYMPTOTICS

TODAY

- STATISTICS & ESTIMATORS
- CONSISTENCY
- BIAS & VARIANCE
- EFFICIENCY
- SUFFICIENCY.

STATISTICS & ESTIMATOR

I) DENOTE $\vec{x} = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$ AS OUR DATA

SAMPLE
↑

II) ASSUME OUR SAMPLE IS THE REALIZATION OF

$$\vec{X} = (X_1, \dots, X_n) \sim \underline{f_0} \in \mathcal{F} \quad \begin{array}{l} \text{STATISTICAL} \\ \text{MODEL.} \end{array}$$

- WE CALL STATISTIC OR SAMPLE STATISTIC TO ANY QUANTITY THAT CAN BE COMPUTED AS A FUNCTION OF THE SAMPLE

$$T(\vec{x}) \text{ FOR SOME } T: (\mathbb{R}^d)^n \rightarrow \mathbb{R}^m$$

- FROM II), THE STATISTIC IS A REALIZATION OF THE RANDOM VECTOR

$T(\vec{x})$ ALSO CALL THIS A STATISTIC!

- IF OUR STATISTICAL MODEL IS PARAMETRIC

$$\mathcal{F} = \{f_\theta : \theta \in \Theta\}, \quad \Theta \subseteq \mathbb{R}^d$$

THEN A STATISTIC THAT ESTIMATES THE PARAMETER θ IS CALLED AN ESTIMATOR OR PARAMETER ESTIMATOR, DENOTED $\hat{\theta}$.

EX: $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} f_{\theta_0} \in \mathcal{F} = \{f_\theta : \theta \in \Theta\}$
 $\theta_0 = \mathbb{E}(X_1)$

ESTIMATORS FOR THE MEAN θ_0 ARE

$$\hat{\theta}_M = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{1ST-MOMENT ESTIMATOR}$$

$$\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \mathcal{L}_n(\theta; \vec{x})$$

$\left[\prod_{i=1}^n f_{\theta_0}(X_i) \right]$

 \downarrow i.i.d.
 LIKELIHOOD FUNCTION

MAXIMUM LIKELIHOOD ESTIMATOR

CONSISTENCY

- WE SAY AN ESTIMATOR $\hat{\theta}$ OF A PARAMETER θ IS CONSISTENT IF

$$\hat{\theta} \xrightarrow{n \rightarrow \infty} \theta \quad \text{IN PROBABILITY.}$$

$$\Leftrightarrow \lim_{n \rightarrow \infty} \mathbb{P}(|\hat{\theta} - \theta| > \varepsilon) = 0$$

- CONSISTENCY IS THE "LEAST" THAT WE CAN EXPECT FROM AN ESTIMATOR!

Ex : 1ST MOMENT ESTIMATOR FOR MEAN μ

$$\bullet \hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$i) \mathbb{E}(\hat{\mu}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i) = \mu \quad \begin{array}{l} \text{i.i.d.} \\ \text{ASSUME} \\ \text{THIS IS FINITE} \end{array}$$

$$ii) \text{Var}(\hat{\mu}) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{1}{n} \text{Var}(X_0) \quad \begin{array}{l} \text{i.i.d.} \\ \uparrow \end{array}$$

$$iii) \mathbb{P}(|\hat{\mu} - \mu| > \varepsilon) \stackrel{\text{MARKOV}}{\leq} \frac{\mathbb{E}[|\hat{\mu} - \mu|^2]}{\varepsilon^2}$$

$$\stackrel{i)}{=} \frac{\text{Var}(\hat{\mu})}{\varepsilon^2} \stackrel{ii)}{=} \frac{\text{Var}(X_0)}{n \varepsilon^2} \xrightarrow{n \rightarrow \infty} 0 \quad \text{CONSISTENT!}$$

OBS : THIS IS A PROOF FOR THE WEAK LAW OF LARGE NUMBERS!

BIAS & VARIANCE

- FOR AN ESTIMATOR $\hat{\theta}$ OF A PARAMETER θ WE DEFINE ITS BIAS AS

$$\text{BIAS}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \theta$$

- WE WANT $\text{BIAS}(\hat{\theta}) = 0$ **UNBIASED ESTIMATOR**

EX: 1ST MOMENT ESTIMATOR FOR MEAN μ

$$\mathbb{E}(\hat{\mu}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(X_i) \stackrel{\text{i.i.d.}}{=} \mathbb{E}(X_1) = \mu$$

UNBIASED!

EX: NAIVE ESTIMATOR FOR VARIANCE σ^2

$$\hat{\sigma}_N^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu})^2$$

$$\mathbb{E}(\hat{\sigma}_N^2) = \left(1 - \frac{1}{n}\right) \sigma^2 \xrightarrow{\frac{n-1}{n}}$$

BIASED!

EX: UNBIASED ESTIMATOR FOR THE VARIANCE σ^2

$$\hat{\sigma}^2 = \left[\frac{n}{n-1}\right] \hat{\sigma}_N^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu})^2$$

$$\mathbb{E}(\hat{\sigma}^2) = \left[\frac{n}{n-1}\right] \mathbb{E}(\hat{\sigma}_N^2) = \left[\frac{n}{n-1}\right] \left[\frac{n-1}{n}\right] \sigma^2$$

UNBIASED!

BUT $\text{Var}(\hat{\sigma}^2) = \text{Var}\left(\left[\frac{n}{n-1}\right] \hat{\sigma}_N^2\right)$

$$= \underbrace{\left[\frac{n}{n-1}\right]^2}_{> 1} \text{Var}(\hat{\sigma}_N^2)$$

$> \text{Var}(\hat{\sigma}_N^2)$

• BIAS/VARIANCE TRADE OFF.

$0 < \|\hat{\theta} - \theta\|_2 := \mathbb{E}[(\hat{\theta} - \theta)^2]$ MEAN SQUARED ERROR

$$= \mathbb{E}[(\hat{\theta} - \mathbb{E}(\hat{\theta}) + \mathbb{E}(\hat{\theta}) - \theta)^2]$$

$$= \mathbb{E}[(\hat{\theta} - \mathbb{E}(\hat{\theta}))^2] + (\mathbb{E}(\hat{\theta}) - \theta)^2$$

$(a+b)^2 = a^2 + b^2 + 2ab$

$$+ 2 \mathbb{E}[(\hat{\theta} - \mathbb{E}(\hat{\theta}))(\mathbb{E}(\hat{\theta}) - \theta)]$$

$$= \text{Var}(\hat{\theta}) + \text{Bias}(\hat{\theta})^2$$

TRADE OFF!

LECTURE 3

LAST TIME

- STATISTICS, ESTIMATORS
- CONSISTENCY
- UNBIASEDNESS
- BIAS / VARIANCE TRADE OFF

TODAY

- ESTIMATORS AGAIN
- EFFICIENCY
- SUFFICIENCY
- CONFIDENCE INTERVALS

EFFICIENCY

- WE KNOW UNBIASED ESTIMATORS WITH ZERO VARIANCE DON'T EXIST.
- HOW SMALL / GOOD CAN BE THEIR VARIANCE ?

CRAMER RAO LOWER BOUND

- ANY UNBIASED ESTIMATOR $\hat{\theta}$ OF θ_0 (UNDER SOME REGULARITY CONDITIONS) :

$$\text{Var}(\hat{\theta}) \geq \mathbf{I}(\theta_0)^{-1}$$

FISHER'S INFORMATION

$$\mathbf{I}(\theta_0) = \mathbb{E} \left[\frac{\partial}{\partial \theta} \log \mathcal{L}_n(\theta_0, \vec{X})^2 \right]$$

TRUE PARAMETER OF STATISTICAL MODEL θ

- WE SAY $\hat{\theta}$ IS EFFICIENT OR THAT IS THE MINIMUM VARIANCE UNBIASED ESTIMATOR (M.V.U.) IFF:

$$\text{Bias}(\hat{\theta}) = 0 \quad \text{AND} \quad \text{Var}(\hat{\theta}) = I(\theta)^{-1}$$

Ex: $\hat{\theta}_{\text{MLE}} = \underset{\theta \in \Theta}{\text{argmax}} \ln(\theta, \vec{X})$

IS ASYMPTOTICALLY UNBIASED AND EFFICIENT (UNDER SOME REGULARITY) CONDITIONS:

$$\sqrt{n}(\hat{\theta}_{\text{MLE}} - \theta) \xrightarrow{n \rightarrow \infty} N(0, I(\theta)^{-1})$$

\downarrow CONVERGENCE IN DISTRIBUTION \downarrow UNBIASED \downarrow EFFICIENT

→ AS LONG AS WE HAVE ENOUGH SAMPLES WE SHOULD USE THE M.L.E.

SUFFICIENCY

• $\vec{X} = (X_1, \dots, X_n) \sim f_{\theta} \in \mathcal{F} = \{f_{\theta} : \theta \in \Theta\}$

• IN SPARIT, A STATISTIC $T(\vec{X})$ IS SUFFICIENT FOR A PARAMETER θ IF ITS USE INVOLVES NO LOSS OF INFORMATION ABOUT θ .

• MATHEMATICALLY, THIS TRANSLATES TO: FOR ANY $t \in \text{RANGE}(T)$:

$\vec{X} | T(\vec{X})$ DISTRIBUTION IS INVARIANT ON $\theta, \theta \in \Theta$

EX: $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} f_p \in \mathcal{F} = \{f_p : p \in (0, 1)\}$

$f_p = \begin{cases} 1 & \text{WITH PROB } p \\ 0 & \text{WITH PROB } 1-p \end{cases}$

$T(\vec{X}) = \sum_{j=1}^n X_j$ (HEADS COUNTER)

$$\mathbb{P}(\vec{X} = \vec{x} | T(\vec{X}) = t) = \frac{\mathbb{P}(\vec{X} = \vec{x}, T(\vec{X}) = t)}{\mathbb{P}(T(\vec{X}) = t)}$$

ASSUMING $T(\vec{x}) = t$ OTHERWISE THE PROB. IS ZERO.

$$= \frac{p_0^t (1-p_0)^{n-t}}{\binom{n}{t} p_0^t (1-p_0)^{n-t}} = \frac{1}{\binom{n}{t}} \rightarrow \text{INVARIANT ON } p_0!$$

SUFFICIENT!

MINIMAL SUFFICIENCY

- WE CAN HAVE SUFFICIENT STATISTICS THAT STORE UNNECESSARY INFORMATION ABOUT THE PARAMETER θ .

Ex: $T(\vec{x}) = \vec{x} \rightarrow [\vec{x} | T(\vec{x}) = \vec{x}] = \vec{x}$

→ SUFFICIENT FOR ANY PARAMETER θ
STORES ALL THE INFORMATION OF THE SAMPLE!

- IN SPIRIT, A SUFFICIENT STATISTIC $T(\vec{x})$ FOR A PARAMETER θ IS MINIMAL SUFFICIENT IF IT STORES NO UNNECESSARY INFORMATION ABOUT θ .

- MATHEMATICALLY, THIS TRANSLATES TO

∀ SUFFICIENT STATISTIC $\tilde{T}(\vec{x})$

∃ A FUNCTION S SUCH THAT

$$S(\tilde{T}(\vec{x})) = T(\vec{x})$$

- LEHMANN-SCHEFFÉ CRITERION

CONFIDENCE INTERVALS

IN A PARAMETRIC STATISTICAL MODEL,

$$\vec{X} = (X_1, \dots, X_n) \sim f_{\theta} \in \mathcal{F},$$

$$\mathcal{F} = \{f_{\theta} : \theta \in \Theta\}, \quad \Theta \subseteq \mathbb{R}^m$$

CONSIDER $\hat{\theta}$ FOR ESTIMATING θ , EVEN IF CONSISTENT, UNBIASED, EFFICIENT, ETC. STILL WE DON'T KNOW HOW PRECISE IS OUR ESTIMATION

SOLUTION: CONFIDENCE INTERVALS

WITH ADEQUATE STRUCTURE WE CAN FIND

$$T_L(\vec{x}), T_U(\vec{x}) \quad \text{s.t.}$$

$$\mathbb{P}(T_L(\vec{x}) \leq \theta \leq T_U(\vec{x})) = \gamma \quad (\approx 1)$$

Ex : $X_1, \dots, X_n \stackrel{i.i.d}{\sim} \mathcal{F}_{\mu_0} \in \mathcal{F}$, WHERE

$\mathcal{F} = \{ \mathcal{N}(\mu, \sigma^2); \mu \in \mathbb{R}, \sigma^2 > 0 \}$

FACT: $D = \frac{\sqrt{n}(\hat{\mu} - \mu_0)}{\hat{\sigma}^2} \sim t_{n-1}$ t-STUDENT DISTRIBUTION
THIS IS NOT AN STATISTIC!
↳ DEGREES OF FREEDOM

→ FOR ANY δ (≈ 0) WE CAN FIND $p \in \mathbb{R}$ SUCH THAT

$$\mathbb{P}(-p \leq D \leq p) \leq \delta$$

$$\rightarrow \mathbb{P}\left(\underbrace{\hat{\mu} - \frac{p\hat{\sigma}^2}{\sqrt{n}}}_{T_L(\vec{x})} \leq \mu_0 \leq \underbrace{\hat{\mu} + \frac{p\hat{\sigma}^2}{\sqrt{n}}}_{T_U(\vec{x})}\right) = 1 - \delta$$

→ WITH PROB. $1 - \delta$: $\mu_0 \in [T_L(\vec{x}), T_U(\vec{x})]$

- δ LARGER → LARGER INTERVAL
- n LARGER → SMALLER INTERVAL!

LECTURE 4

LAST TIME

- EFFICIENCY
- SUFFICIENCY
- CONFIDENCE INTERVALS

NOW

- HYPOTHESIS TESTING
- ERRORS (TYPE I & II)
- LIKELIHOOD RATIO TEST

TESTING HYPOTHESIS

- MODERN SCIENCE ASSUMPTION (POPPER)
HYPOTHESIS CAN NOT BE PROVEN TRUE
BUT CAN BE FALSIFIABLE: WE CAN REJECT
THEM IF EXPERIENCE DON'T MATCH IT.
- MANY POSSIBLE STATISTICAL APPROACHES
 - BAYESIAN
 - FREQUENTIST
 - FISHER
 - NEYMAN & PEARSON

MOST WIDESPREAD AND FAMOUS

NEYMAN & PERSON FREQUENTIST TEST

H_0 : NULL HYPOTHESIS (TO BE TESTED)

H_A : ALTERNATIVE HYPOTHESIS (H_0 NEGATION)

TO TEST H_0 :

0) COLLECT A SAMPLE $\vec{x} = (x_1, \dots, x_n)$

1) ASSUME A STATISTICAL MODEL \mathcal{F} FOR \vec{x}
WE SHOULD BE EXPRESS H_0 IN TERMS
OF A SUBSET $\mathcal{F}_0 \subseteq \mathcal{F}$

2) SELECT A STATISTIC $T_f(\vec{x})$, FOR ANY $f \in \mathcal{F}_0$
WE ASSUME WE KNOW f WHEN COMPUTING IT.
LARGER VALUES OF $T_f(\vec{x})$ SHOULD REFLECT
EVIDENCE AGAINST H_0 .

3) SELECT A SIGNIFICANCE LEVEL $\alpha \approx 0$

4) COMPUTE THE P-VALUE:

$$P(\vec{x}) = \sup_{f \in \mathcal{F}_0} \mathbb{P}_f(T_f(\vec{x}) \geq T_f(\vec{x}))$$

$\underbrace{f \in \mathcal{F}_0}_{\text{BEST CASE OVER } H_0}$ \downarrow RANDOM VARIABLE \downarrow CONCRETE OBSERVATION

UNDER H_0 , THE PROBABILITY OF T BEING MORE
DEVIATED THAN OUR OBSERVATION

$P(\vec{x}) \approx 0 \rightarrow$ OUR OBSERVATION IS DEVIATED UNDER H_0
 $\rightarrow H_0$ IS UNLIKELY TO FIT OUR EXPERIENCE

5) REJECT H_0 IFF $P(\vec{x}) \leq \alpha$.

Ex: $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} f_{\vec{\theta}} \in \mathcal{F}$
 $\vec{\theta} \rightarrow (\bar{\mu}, \bar{\sigma}^2)$

1) $\mathcal{F} = \{ \mathcal{N}(\mu, \sigma^2) ; \mu \in \mathbb{R}, \sigma^2 > 0 \}$

$H_0 : \bar{\mu} \leq c \leftrightarrow \mathcal{F}_0 = \{ \mathcal{N}(\mu, \sigma^2) : \mu \leq c, \sigma^2 > 0 \}$

$H_a : \bar{\mu} > c$

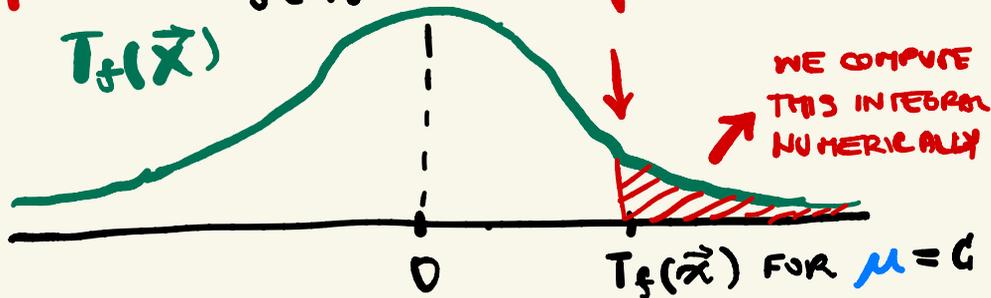
2) FOR ANY $f \in \mathcal{F}_0$, WITH PARAMETERS (μ, σ^2) , CONSIDER THE STATISTIC

$$T_f(\vec{x}) = \frac{\sqrt{n}(\hat{\mu} - \mu)}{\hat{\sigma}} \sim t_{n-1}$$

\downarrow
ASSUMING $\vec{X} \sim f$

3) CHOOSE $\alpha = 0.05 \approx 0$.

4) $P(\vec{x}) = \sup_{f \in \mathcal{F}_0} P_f(T_f(\vec{x}) \geq T_f(\vec{x}))$



FOR ANY OTHER $f \in \mathcal{F}_0$ WITH $\mu < c$, $T_f(\vec{x})$ IS TRANSLATED TO THE RIGHT IN THE ABOVE PLOT.
→ THE AREA UNDER THE CURVE IS NOT LARGER!
→ THE SUPRENUM IS ATTAINED FOR $\mu = c$.

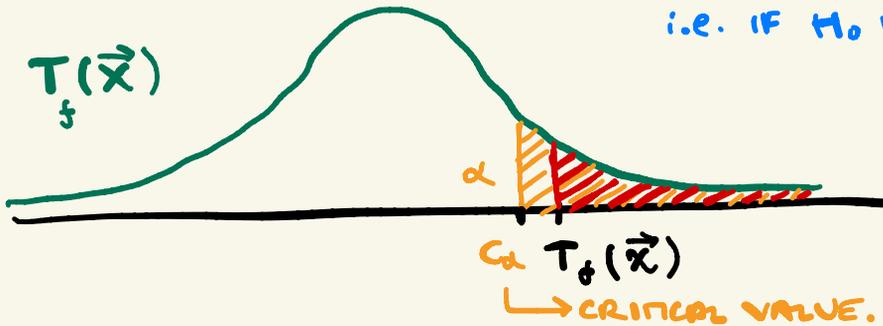
5) REJECT H_0 IFF $P(\vec{x}) \leq \alpha = 0.05$.

TESTING ERRORS

TYPE I: REJECTING H_0 WHEN IT IS TRUE

$$P(\text{TYPE I ERROR}) \leq \alpha$$

→ "=" IF $T_0 = \bar{x}_f$
i.e. IF H_0 IS SIMPLE



$$P(\vec{x}) \leq \alpha \iff T_f(\vec{x}) \geq C_\alpha$$

- BY CONSTRUCTION, TYPE I ERROR IS LOW!
- THE TEST IS "PLAYING FAIR" WITH H_0

TYPE II: NOT REJECTING H_0 WHEN IT IS FALSE.

$$P(\text{TYPE II ERROR}) =: \beta$$

"POWER" OF TEST := $1 - \beta$

- β QUANTIFIES HOW "SEVERE" IS THE TEST WITH H_0 .
- β IS NOT SET BY THE SCIENTIST
IT IS A PROPERTY OF ANY TEST WITH FIXED α .

LIKELIHOOD RATIO TEST

- CONSIDER A STATISTICAL MODEL \mathcal{F} AND A NULL HYPOTHESIS $\mathcal{F}_0 \subseteq \mathcal{F}$.
- THE LIKELIHOOD RATIO TEST IS DEFINED BY ITS STATISTIC

$$T(\vec{x}) = -2 \ln \left[\frac{\sup \{ f(\vec{x}) : f \in \mathcal{F}_0 \}}{\sup \{ f(\vec{x}) : f \in \mathcal{F} \}} \right]$$

• NEYMAN PEARSON LEMMA

IF $\mathcal{F} = \{f_0, f_A\}$, $\mathcal{F}_0 = \{f_0\}$

(SIMPLE NULL AND SIMPLE ALTERNATIVE)

THEN, FOR ANY FIXED SIGNIFICANCE α ,

THE LIKELIHOOD RATIO TEST IS THE

MOST POWERFUL TEST WITH SIGNIFICANCE

α .

LECTURE 5

LAST TIME

- TESTING
(NEYMAN-PEARSON)

TODAY

- MORE TESTING
- FINITE SAMPLE TEST
- ASYMPTOTIC TESTS
- LIKELIHOOD RATIO TEST

FINITE SAMPLE TESTS

- WITH ADEQUATE STRUCTURE, WE KNOW EXACTLY THE DISTRIBUTION OF THE TEST STATISTIC $T(\vec{x})$

→ THEN WE CAN EXACTLY^{*} COMPUTE THE **P-VALUE** OF TEST FOR ANY FINITE n .

* UP TO NUMERICAL ERRORS

EX: ONE SAMPLE TESTING

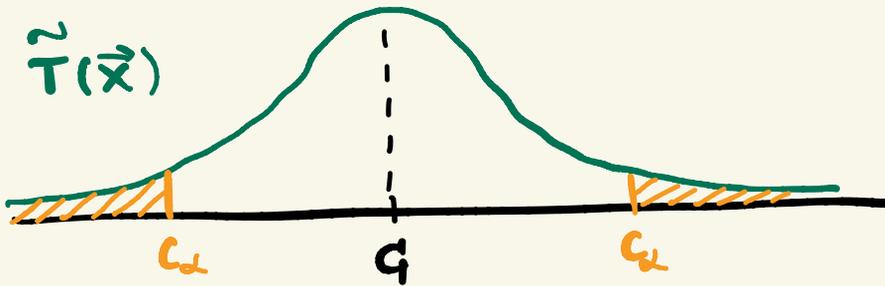
- $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} f_{\theta} \in \mathcal{F}_N$, WHERE

$$\mathcal{F}_N = \{ \mathcal{N}(\mu, \sigma^2); \mu \in \mathbb{R}, \sigma^2 > 0 \}$$

- $H_0: \bar{\mu} = c, H_A: \bar{\mu} \neq c$

- $T(\vec{x}) = |\tilde{T}(\vec{x})|, \tilde{T}(\vec{x}) = \frac{\sqrt{n}(\hat{\mu} - c)}{\hat{\sigma}}$
 $\sim t_{n-1}$, b.c. \mathcal{F}_N

- FOR ANY $\alpha, C_{\alpha} = t_{n-1}(\alpha/2)$.



- REJECT IFF $T(\vec{x}) \geq C_{\alpha}$

EX: TWO SAMPLES TESTING.

- 1ST SAMPLE: $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} f_{\theta_X} \in \mathcal{F}_N$,
- 2ND SAMPLE: $Y_1, \dots, Y_m \stackrel{i.i.d.}{\sim} f_{\theta_Y} \in \mathcal{F}_N$

(TWO SAMPLES, EACH FROM A DIFFERENT DISTRIBUTION OF THE SAME STATISTICAL MODEL. ASSUME SAME VARIANCE)
 $\sigma_X = \sigma_Y$

- $H_0: \mu_X = \mu_Y$, $H_A: \mu_X \neq \mu_Y$.

$$T(\vec{X}, \vec{Y}) = \left| \frac{\hat{\mu}_X - \hat{\mu}_Y}{\underbrace{\sqrt{\hat{\sigma}_P^2 \left(\frac{1}{n} + \frac{1}{m} \right)}}_{\sim t_{n+m-2}}} \right|, \text{ WHERE}$$

BECAUSE OF \mathcal{F}_N

$$\hat{\sigma}_P^2 = \frac{(n-1)\hat{\sigma}_X^2 + (m-1)\hat{\sigma}_Y^2}{n+m-2} \quad \left(\text{POOLED } t\text{-TEST} \right. \\ \left. \text{FOR EQUAL VARIANCES} \right)$$

- FOR ANY α , $C_\alpha = t_{n+m-2}(\alpha/2)$
- REJECT IFF $T(\vec{X}, \vec{Y}) \geq C_\alpha$

"LARGE SAMPLE" TEST (ASYMPTOTIC RESULT)

- ASSUME $X_1, \dots, X_n \stackrel{i.i.d}{\sim} f_0 \in \mathcal{F}$ ANYTHING
- WE **DON'T KNOW** THE DISTRIBUTION OF ANY STATISTIC WE CAN THINK OF.
- WE CAN CONSIDER THE STATISTIC

$$T_f(\vec{x}) = \frac{\sqrt{n}(\hat{\mu} - \mu)}{\sigma} \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, 1)$$

CONVERGENCE IN DISTRIBUTION, BY C.L.T.

WHERE $\mu = \mathbb{E}_f(X_1)$, $\sigma = \sqrt{\text{Var}_f(X_1)}$

- WE KNOW THE **ASYMPTOTIC** DISTRIBUTION OF $T_f(\vec{x})$.
- OTHER DISTRIBUTIONS THAT ARISE FROM ASYMPTOTICS: χ^2 , RATIOS OF χ^2 'S.

EX : ONE SAMPLE TESTING OF PROPORTIONS.

- $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} f_p \in \mathcal{F}_p$, WHERE

$$\mathcal{F}_p = \{f_p : p \in (0,1)\}, \quad f_p = \begin{cases} 1 & \text{w.p. } p \\ 0 & \text{w.p. } 1-p \end{cases}$$

- $H_0 : \bar{p} = p_0$, $H_A : \bar{p} \neq p_0$.

- $T_f(\vec{x}) = |\tilde{T}_f(\vec{x})|$

$$\tilde{T}_f(\vec{x}) = \frac{\sqrt{n}(\hat{p} - p_0)}{\sqrt{p_0(1-p_0)}} \stackrel{n \rightarrow \infty}{\sim} \mathcal{N}(0,1)$$

$(\hat{p} = \bar{\mu})$

- FOR ANY α , $C_\alpha = z(\alpha/2)$

- REJECT IFF $T_f(\vec{x}) \geq C_\alpha$

Ex : TWO SAMPLE PROPORTIONS.

- $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} f_{P_x} \in \mathcal{F}_P$
 $Y_1, \dots, Y_n \stackrel{\text{i.i.d.}}{\sim} f_{P_y} \in \mathcal{F}_P$

- $H_0: P_x = P_y$, $H_A: P_x \neq P_y$

- $T(\vec{x}, \vec{y}) = |\tilde{T}(\vec{x}, \vec{y})|$ (WALD STATISTIC WITH POOLING)

$$\tilde{T}(\vec{x}, \vec{y}) = \frac{\hat{p}_x - \hat{p}_y}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n} + \frac{1}{m}\right)}} \stackrel{n \rightarrow \infty}{\sim} \mathcal{N}(0,1)$$

- FOR ANY α , $C_\alpha = z(\alpha/2)$

- REJECT IFF $T(\vec{x}, \vec{y}) \geq C_\alpha$

(NUMERICAL METHODS)

LECTURE 6

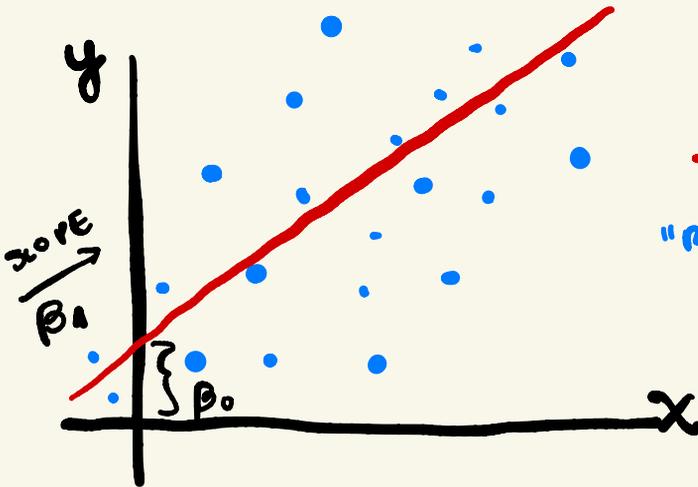
LAST TIME

- NEWMAN-PEARSON TESTING
- FINITE SAMPLE
- ASYMPTOTICS

NOW

- LINEAR REGRESSION
- MODEL
- ESTIMATION
- FURTHER TOPICS

WHAT IS LINEAR REGRESSION?



- DATA
- LINEAR REGRESSION

"RESPONSE" \uparrow "PREDICTOR" \uparrow

$$y_i = \beta_0 + \beta_1 x_i + \epsilon$$

β_0 INTERCEPT

GEOMETRICALLY: "BEST" LINEAR APPROXIMATION TO OUR DATA

STATISTICALLY?

STATISTICAL MODEL

- $x^1, \dots, x^n \in \mathbb{R}^d$ (PREDICTOR)
 INDICES
- $y_1, \dots, y_n \in \mathbb{R}$ (RESPONSE)

MAIN ASSUMPTION

$$y_i = \beta_0 + \sum_{k=1}^d x_k^i \beta_k + \epsilon_i \quad \forall i \in \{1, \dots, n\}$$

WHERE $\epsilon_1, \dots, \epsilon_n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2) \rightarrow$ THIS CAN BE RELAXED!

STATISTICAL MODEL COMPACT FORM

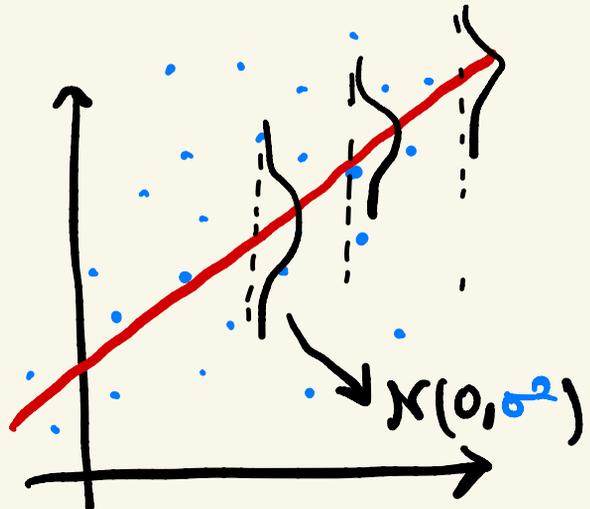
- $\vec{Y} \sim f_{\theta} \in \mathcal{F} = \{f_{\theta} : \theta = (\beta, \sigma^2) \in \mathbb{R}^{d+1} \times \mathbb{R}_{++}\}$

- $f_{\theta} = X\beta + \epsilon$

- $\epsilon \sim \mathcal{N}(0, \sigma^2 I_d)$

- $X = \begin{bmatrix} 1 & x^1 \top \\ \vdots & \vdots \\ 1 & x^n \top \end{bmatrix}$

\uparrow
 $\mathbb{R}^{n \times (d+1)}$



PARAMETER ESTIMATION

- WE FOCUS ON THE β PARAMETER. CONSIDER

$$T(X, y) = \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{2} \|X\beta - y\|_2^2$$

$$\text{LET } h(\beta) = \frac{1}{2} \|X\beta - y\|_2^2$$

$$\nabla h(\beta) = \frac{2}{2} X^T (X\beta - y)$$

OPTIMAL
LEAST
SQUARES
ESTIMATOR
(O.L.S)

OPTIMALITY CONDITIONS FOR CONVEX FUNCTIONS GIVE

$$\nabla h(\hat{\beta}) = 0 \leftarrow \hat{\beta} = (X^T X)^{-1} X^T y$$

OBS: $X^T X$ IS INVERTIBLE IFF THE COLUMNS OF X ARE LINEARLY INDEPENDENT.

IF THEY AREN'T, REPLACE INVERSE WITH PSEUDOINVERSE. (MATRIX ANALYSIS)

$$\hat{\beta} = (X^T X)^{\dagger} X^T y$$

LINEAR ESTIMATOR
W.R.T y !

LEAST SQUARES IS BLUE

BEST ✓
LINEAR ✓
UNBIASED ✓
ESTIMATOR

$$\bullet \mathbb{E}(\hat{\beta}) = \mathbb{E}((X^T X)^{-1} X^T \vec{y})$$

↳ SOURCE OF RANDOMNESS

$$= \mathbb{E}((X^T X)^{-1} X^T (X \beta + \epsilon))$$

$$= \mathbb{E}((X^T X)^{-1} (X^T X) \beta)$$

$$+ \mathbb{E}((X^T X)^{-1} X^T \epsilon)$$

LINEARITY
OF EXPECTATION

$$= \mathbb{E}(\beta) + (X^T X)^{-1} X^T \mathbb{E}(\epsilon)$$

$$= \beta \rightarrow \text{UNBIASED ESTIMATOR}$$

GAUSS MARKOV THEOREM

ANY ESTIMATOR $\tilde{\beta}$ THAT IS A LINEAR
FUNCTION OF y AND IS UNBIASED
HAS LARGER VARIANCE THAN $\hat{\beta}$:

$$\underbrace{\mathbb{E}(\tilde{\beta} \tilde{\beta}^T)}_{\text{COVARIANCE MATRIX}} \succeq \mathbb{E}(\hat{\beta} \hat{\beta}^T) = I_d \sigma^2$$

SUMMARY

WE HAVE SEEN

- SAMPLES & STATISTICAL MODEL
- STATISTICS & ESTIMATORS
- CONFIDENCE INTERVALS & TESTING
- LINEAR REGRESSION

MORE OF THIS IN AMS:

- MATHEMATICAL STATISTICS (FOUNDATIONS) [FALL]
- STATISTICAL THEORY I (DISCUSSION) [FALL]
- STATISTICAL THEORY II (ASYMPTOTICS) [SPRING]

WHAT WE DIDN'T COVER

- BAYESIAN STATISTICS [FALL / SPRING]
- STATISTICAL LEARNING
 - ↳ • STATISTICAL PATTERN RECOGNITION [SPRING]
 - MACHINE LEARNING I [FALL]
 - MACHINE LEARNING II [SPRING]
- NON PARAMETRIC STATISTICS
 - ↳ • APPLIED STATISTICS AND DATA ANALYSIS II [S]
- HIGH DIMENSIONAL STATISTICS / PROBABILITY
 - ↳ • MATHEMATICS OF DATA SCIENCE [FALL]
 - HIGH DIM. PROB. & APPROXIMATION [SPRING]