

YZV 231E

10.01.2022

Probability Theory & Stats

GU.

Recap:

MLE : Maximum Likelihood Estimation:

model w/ unknown parameters

$$\hookrightarrow X \sim P_X(x_i | \theta)$$

pick  $\theta$  that makes data likely:

$$\left. \begin{array}{l} p_X(x_1 | \theta_1) \\ p_X(x_f | \theta_2) \end{array} \right\}$$

a family of parameterized models

$$\max_{\theta} p_X(x, \theta)$$

Compare to Bayesian

MAP estimation

$$\max_{\theta} \widehat{p}_{\text{post}}(\theta | x)$$

↑↑  
posterior

$$\stackrel{\text{prior}}{=} \max_{\theta} \frac{\widehat{p}_{\text{prior}}(\theta) \cdot p_{\text{data}}(x | \theta)}{p_X(x)}$$

→ Sample Mean Estimate of  $\theta$

$$\text{r.v. } \widehat{\theta}_n = \frac{x_1 + \dots + x_n}{n} \rightarrow$$

→  $(1 - \alpha)$  confidence interval

$$P(\widehat{\theta}_n^- \leq \theta \leq \widehat{\theta}_n^+) \geq 1 - \alpha, \forall \theta$$

↑  
e.g. 95%

Note: MLE & MAP appear the same when we have a uniform prior in MAP.

Construction of CI : Easy w/ CLT :

Confidence Interval for the sample mean  $\hat{\theta}_n$

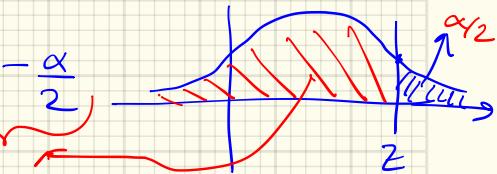
$$P\left(\hat{\theta}_n - \frac{z \cdot \sigma}{\sqrt{n}} \leq \theta \leq \hat{\theta}_n + \frac{z \cdot \sigma}{\sqrt{n}}\right) \approx 1-\alpha$$

where  $z$  is s.t.

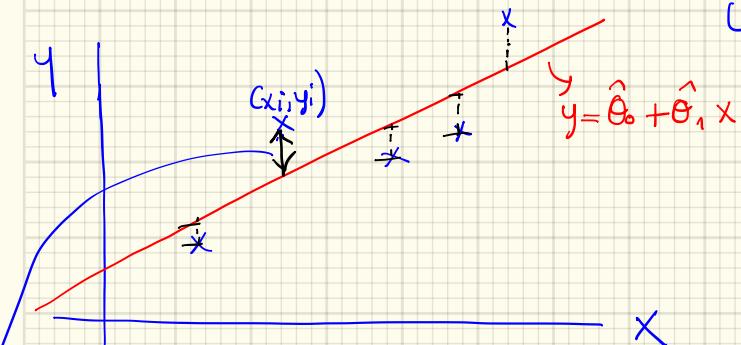
$$\Phi(z) = 1 - \frac{\alpha}{2}$$

pick an  $\alpha \rightarrow$  fix  $z$ .

need an estimate of the variance  $\sigma^2$ .



# REGRESSION :



Let  $X$  be your TUT exam score  
 Let  $y$  be your ITU GPA

Is there a relation between the two?

- Data  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

• Model :  $y \sim \boxed{\theta_0} + \boxed{\theta_1} x$

*unknown parameters  
that define our model.*

Find the "best model" to explain the data:  
 (always "optimal")  
 ask : w.r.t. which criterion ?

minimize residual error :  $y_i - \theta_0 - \theta_1 x_i$

$$\min_{\theta_0, \theta_1} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

residual error



Cost : sum of squared error in predictions .

→ Probabilistic interpretation:

$$\rightarrow y_i = \theta_0 + \theta_1 x_i + w_i, \quad w_i \sim \mathcal{N}(0, \sigma^2), \text{i.i.d.}$$

Do  $m \in \mathbb{E}$ ;  $P_{y|X}(x_i, y_i; \theta)$  : likelihood

$$\begin{aligned} w &\rightarrow C e^{-w_i^2/2\sigma^2} \\ &\rightarrow C \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2 \right\} \end{aligned}$$

take a log.

w<sub>i</sub>'s are indep.  $(e^{-w_1^2/2\sigma^2} \cdot e^{-w_2^2/2\sigma^2} \cdots)$

Maximize w.r.t.  $\theta_0 \vee \theta_1$ .

After you take log → this is the same cost as in  $(*)$

$\Rightarrow$  Linear Regression  $\equiv \text{MLE}$  where  $w_i \sim \mathcal{N}(0, \sigma^2)$ , i.i.d.

# Linear Regression:

- Model :  $y \approx \theta_0 + \theta_1 x$

$$\min_{\theta_0, \theta_1} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

- Solution : set the derivatives of the cost fn. to zero:

$$\bar{x} = \frac{x_1 + \dots + x_n}{n} ; \quad \bar{y} = \frac{y_1 + \dots + y_n}{n}$$

exercise:  
obtain  $\theta_0$  &  $\theta_1$ .

$$\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}$$

$$\hat{\theta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Covariance  $(X, Y)$   
 $\approx E[(X - E[X])(Y - E[Y])]$

(\*)



## Interpretation of the Solution

Assume a model:  $y = \theta_0 + \theta_1 x + w$  w/  $x, w$  are independent, w/ zero mean.

$$E[y] = \theta_0 + \theta_1 E[x] + 0$$

$$\theta_0 = \underbrace{E[y]}_{\text{estimate w/ sample mean}} - \theta_1 \underbrace{E[x]}_{\text{we assumed we have } \bar{x}}$$

$$\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}$$

To obtain  $\theta_1$  estimate:

$$y \cdot x = \theta_0 \cdot x + \theta_1 \cdot x^2 + w \cdot x$$

$$E[y \cdot x] = 0 + \theta_1 \text{Var}(x) + E[w] \cdot E(x)$$

for zero mean r.v.s  $y = \text{Cov}(x, y) = \theta_1 \cdot \text{Var}(x)$

$$\theta_1 = \frac{\text{Cov}(x, y)}{\text{Var}(x)} \rightarrow \text{cov}(x, y) \approx \frac{1}{n} \sum_i (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{Var}(x) \approx \frac{1}{n} \sum_i (x_i - \bar{x})^2$$

→ You may want to include more variables that your ITU GPA may depend on: family income, high school gpa, years of education of parents..

## Multiple Linear Regression:

- Data:  $(x_i^1, x_i^2, x_i^3, y_i)$ ,  $i = 1, \dots, n$   
 $\begin{matrix} \downarrow \\ \text{fam. income} \end{matrix}$   $\begin{matrix} \downarrow \\ \text{high school gpa} \end{matrix}$   $\begin{matrix} \downarrow \\ \text{education years of parents} \end{matrix}$   $\begin{matrix} \downarrow \\ \text{itu gpa} \end{matrix}$
- Model:  $y \approx \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \theta_3 x^3$  : linear fn. of all the variables
- Formulation:  $\sum \text{multiple explanatory variables}$

$$\min_{\theta_0, \theta_1, \theta_2, \theta_3} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i^1 - \theta_2 x_i^2 - \theta_3 x_i^3)^2$$

↓ take deriv. w.r.t. parameters, set to 0 , you set a system of linear eqns.

→ In matrix notation :

Digression: In vector notation: set  $\underline{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}_{n \times 1 \text{ (4x1)}}$

$$\underline{x}_i = \begin{bmatrix} 1 \\ x_i^{(1)} \\ x_i^{(2)} \\ x_i^{(3)} \end{bmatrix}_{n \times 1 \text{ (4x1)}}$$

$$\min_{\underline{\theta}} \|\underline{y} - \underline{x}^T \underline{\theta}\|^2 = f(\underline{\theta}) \quad \text{min. w.r.t. } \underline{\theta}$$

$$\underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}_{m \times 1} ; \underline{X} = \begin{bmatrix} \underline{x}_1 & \underline{x}_2 & \dots & \underline{x}_m \end{bmatrix}_{n \times m}$$

$$f(\underline{\theta}) = (\underline{y}^T - \underline{\theta}^T \underline{X})(\underline{y} - \underline{x}^T \underline{\theta})$$

$$f(\underline{\theta}) = \underline{y}^T \underline{y} - \underline{y}^T \underline{x}^T \underline{\theta} - \underline{\theta}^T \underline{x} \underline{y} + \underline{\theta}^T \underline{X} \underline{x}^T \underline{\theta}$$

# data instances

w.r.t.  $\underline{\theta}$ ) derivative

$$\nabla_{\underline{\theta}} f = \begin{bmatrix} \frac{\partial f}{\partial \theta_0} \\ \frac{\partial f}{\partial \theta_1} \\ \vdots \\ \frac{\partial f}{\partial \theta_n} \end{bmatrix} = 0$$

$$-2 \underline{x} \cdot \underline{y} + 2 \underline{x} \underline{x}^T \underline{\theta} = 0$$

$$\underline{\theta} = (\underline{x} \underline{x}^T)^{-1} \underline{x}^T \underline{y}$$

$$\underline{x}^T = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_m^T \end{bmatrix}$$

$$\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$$

$$\underline{\theta}_{n \times 1} = (\underline{X}^T \underline{X})_{n \times n}^{-1} \underline{X}^T \underline{y}_{m \times 1}$$

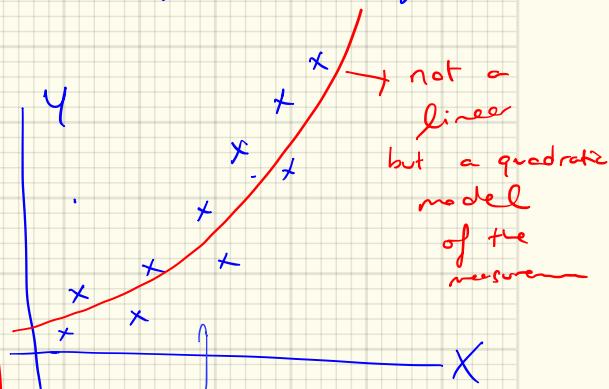
: Normal Equations  
give us the

$\underline{\theta}$  estimates  
for linear regression.

$$\underline{X} = \begin{bmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_m^T \end{bmatrix}_{m \times n}$$

---


$$Y \approx \underline{\theta}_0 + \underline{\theta}_1 \cancel{\cdot X} + h(X)$$



Nonlinear functions of the Data

Model :  $Y \approx \underline{\theta}_0 + \underline{\theta}_1 h(X)$

e.g.  $Y \approx \underline{\theta}_0 + \underline{\theta}_1 \cdot X^2$

Same  
Formulation

• Data points :  $(y_i, h(x_i))$

• Model :  $y = \theta_0 + \theta_1 h(x)$

Still a linear model

$$\min_{\theta} \sum_{i=1}^n \left( y_i - \frac{\theta_0}{\uparrow} - \frac{\theta_1}{\uparrow} h(x_i) \right)^2 \quad \left. \begin{array}{l} \text{eq. } x^3 \\ \text{fn. of } x. \end{array} \right\}$$

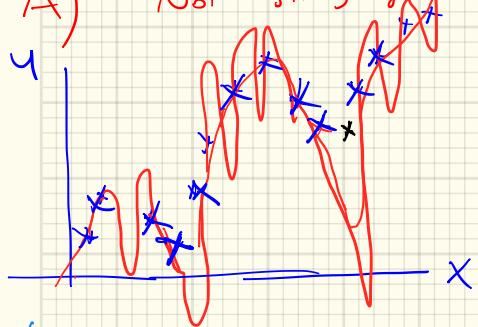
Q) In your linear regression, how do we choose  $h(x)$ 's?

A) Not straightforward : b/c e.g. you try to go over all data points by using a large # of parameters.

e.g. use an  $8^{th}$  degree polynomial

$$y_i = \sum_{j=0}^8 \theta_j x_i^j$$

$$= \theta_0 + \theta_1 x_i + \theta_2 x_i^2 + \dots + \theta_3 x_i^3 + \dots$$



$$(y_i - \sum_j \theta_j x_i^j)$$

error is small  
we have lots of parameters

Q: Is this a good model? No!

→ Choosing complex  $h(\cdot)$  = how many exploratory variables you should be using?

How much complex?

→ Open topic.

→ Start w/ simpler models, especially when you have few data points. → a few parameters

Good  
rule

---

Notes: In practice;  
→ confidence intervals for  $\theta_i$ : (not covered here)

→ "Standard error" estimates of  $\sigma$ : noise in the model.

$$Y \approx \theta_0 + \theta_1 X + \omega \rightarrow \sigma^2$$

uncertainty in the model

$\sigma^2$ : variance of the noise

→  $R^2$ : measure of the explanatory power of the model for linear regression.

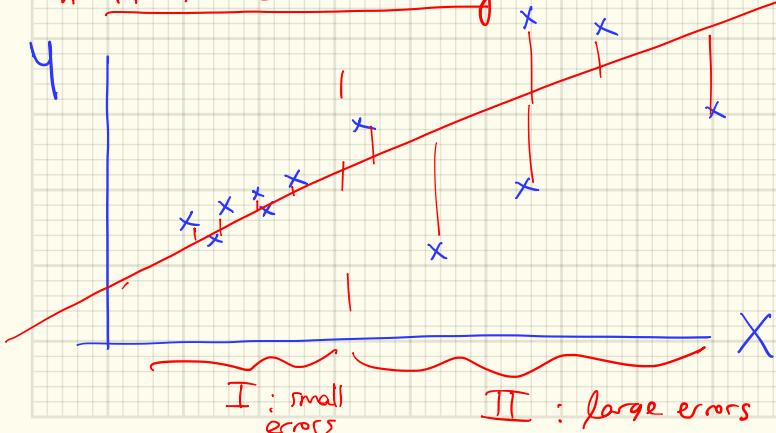
$$\frac{\text{Var}(Y|X)}{\text{Var}(Y)} \rightarrow$$

variance in  $Y$  conditioned on  $X$ .  
 naturally.  
 (uncertainty) in  $Y$  w/o knowing any other  
 information

e.g. "60% of students' GPA is explained by the T4T score."  
 ↓  
 in statistical studies.

### Using Linear Regression : Some pitfalls.

#### \* Heteroskedasticity:



A good fit in region I,  
 I) has small variance  
 II) has large variance.

X-Space has  
 very different variances  
 in different regions.

-Be careful (not covered here).

\* Multi-Collinearity: Multiple explanatory variables ; they are closely related w/ each other.

model

$$Y_{GPA} = \theta_0 + \theta_1 T4T_1 + \theta_2 T4T_2$$

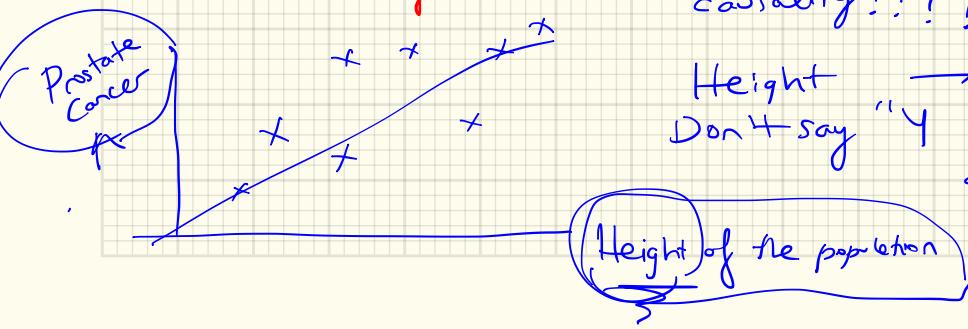
Your  $T4T_1$  &  $T4T_2$  (2 exam scores) are close to each other

$$\begin{aligned} Y_{GPA} &= \theta_0 + \underline{\theta_1 T4T_1} \\ \text{or} \quad &= \theta_0 + \underline{\theta_2 T4T_2} \end{aligned}$$

avoid such redundancy in explanatory variables.

\* Causality !!

People use linear regression to conclude causality!!!!



Height  $\rightarrow$  prostate  
Don't say "Y is caused by X  
according to your model".

good linear

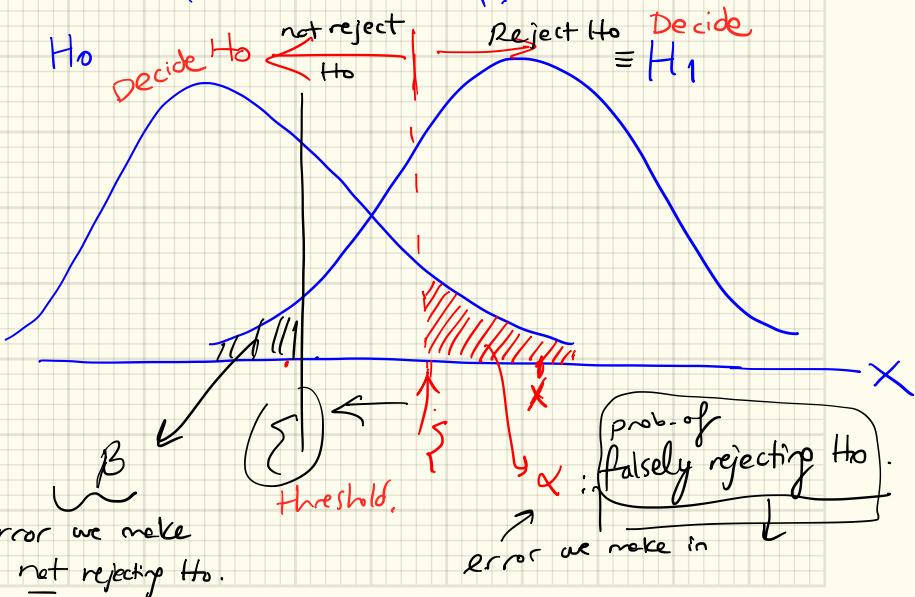
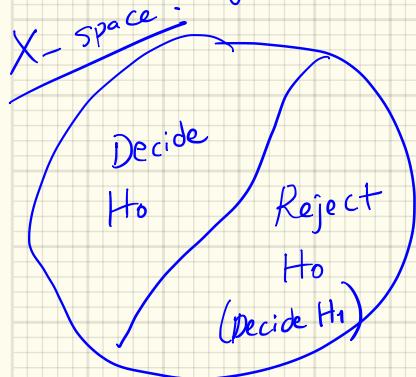
## Hypothesis Testing :

We have a null ( $H_0$ ) & and an alternate ( $H_1$ ) hypothesis

Coin: fair or not:  $H_0 : p = \frac{1}{2}$  vs.  $H_1 : p = \frac{3}{4}$

or  $p \neq \frac{1}{2}$   
 $p \neq \frac{1}{2}$  or  $p = \frac{3}{4}$

Testing coin's fairness ( $p = \frac{1}{2}$ ) against  
 $H_0$  vs  $H_1$ .



We want both

$\alpha \times \beta$  errors

to be as

small as

possible,

for a  
trade-off

$P_x(X|H_0)$

Decide  $H_0$

Decide  $H_1$

$P_x(X|H_1)$

$\alpha$   
move threshold to the right  
to make  $\alpha$  small

$\beta$

$\alpha \downarrow \beta \uparrow$

Q. How to set the threshold?

A. Likelihood Ratio Test (LRT) : look at posterior prob. of the hypothesis

Choose  $H_1$  if  $P(H_1 | X=x) > P(H_0 | X=x)$

pick the hypothesis that is more likely, given the data.

In a Bayesian setting (MAP), use Bayes rule.

$$\frac{P(X=x | H_1) P(H_1)}{P(X=x)} > \frac{P(X=x | H_0) P(H_0)}{P(X=x)} \Rightarrow \text{rewrite}$$

$$L(x) = \frac{P(X=x | H_1)}{P(X=x | H_0)} > \underbrace{\frac{P(H_0)}{P(H_1)}}_{\text{likelihood ratio}} \quad (\text{LRT})$$

likelihood ratio → compare to a threshold  $\xi$ .

- In a Bayesian setting: threshold  $\xi \propto$  prior probabilities of the two hypothesis.
- In a non-Bayesian setting, we don't have prior probabilities

$$\left\{ \frac{P(X=x | H_1)}{P(X=x | H_0)} \stackrel{\text{prob.}}{>} \xi \stackrel{\text{i.t.o. pdf}}{=} \frac{p_X(x; H_1)}{p_X(x; H_0)} > \xi \right.$$

if this ratio is large,

it is unlikely that my observation  $X$  occurred under  $H_0$ ,

∴ reject  $H_0$ .

## Simple Binary Hypothesis Testing:

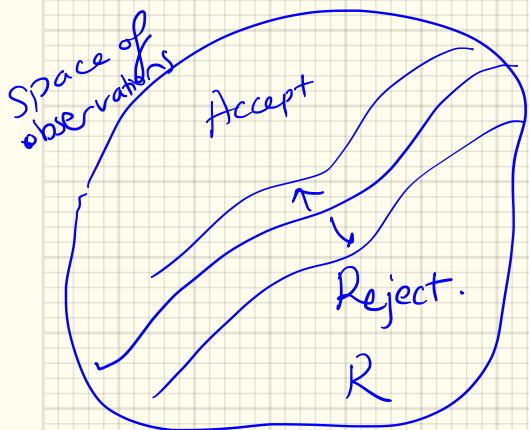
• Default hypothesis:

↳ Null Hypothesis  $H_0$ :  $X \sim p_x(x; H_0)$

— Alternative Hypothesis  $H_1$ :  $X \sim p_x(x; H_1)$

Want to check whether  $H_0$  is false or not.

} Whether to reject or not reject the null hypothesis.



Designing the hypothesis test :

1) Structure of the test : shape of the dividing curve

2) Given the shape , where to place the division ?

1) Choose a rejection region  $R$

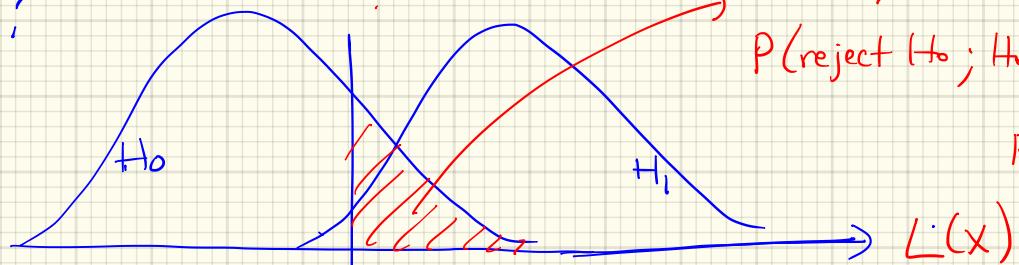
Reject  $H_0$  iff data  $X \in R$ .

→ Structure of the test : Likelihood Ratio test

$$1) \underbrace{L(x)}_{P_x(x; H_1)} = \frac{P_x(x; H_1)}{P_x(x; H_0)} > \xi$$

High value of  $L(x)$   $\equiv$  likelihood of  $H_1$   $>$  likelihood of  $H_0$

2)  $\xi$ ? How to choose  $\xi$ ?  $\xrightarrow{\text{Fix } \alpha}$ ; choose  $\xi$  so that



$$P(\text{reject } H_0; H_0) = \alpha$$

prob. of false  
rejection.

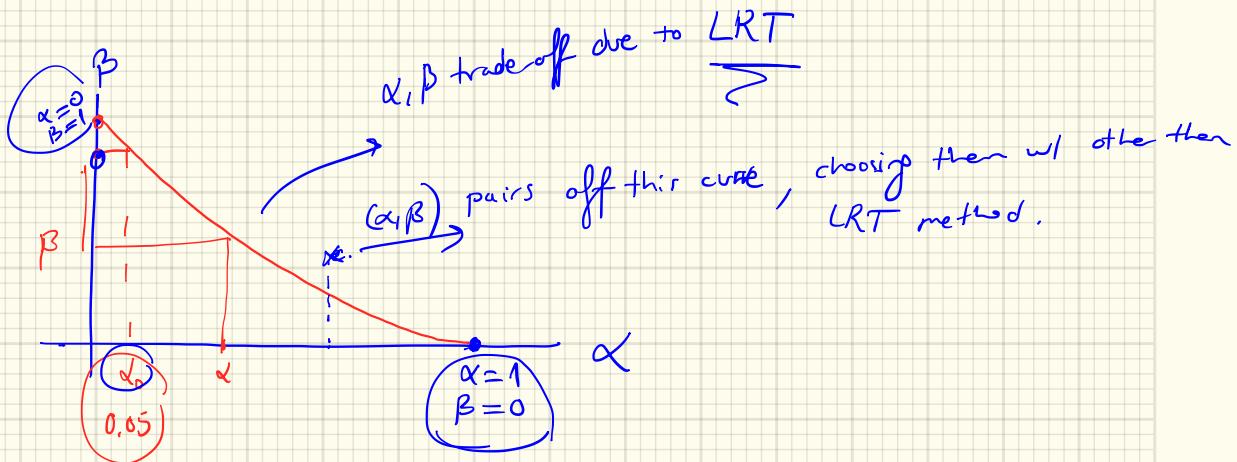
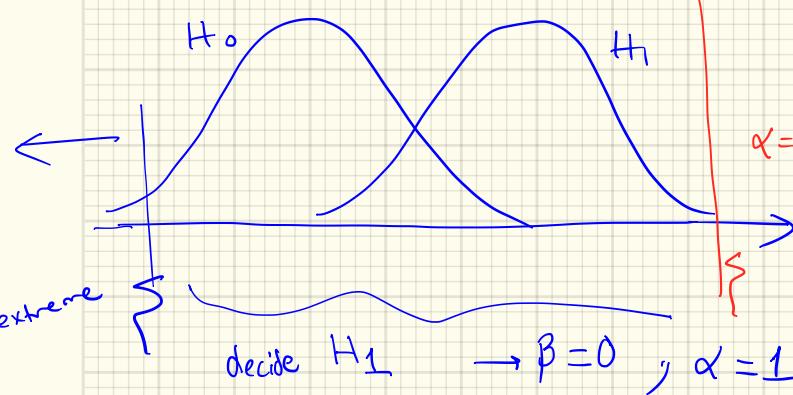
$$\begin{cases} L(x) < \xi \\ L(x) > \xi \end{cases}$$

e.g. set  $\alpha = 0.05 \rightarrow 5\%$

Note :  $(\alpha, \beta)$  trade-off.

$H_0$

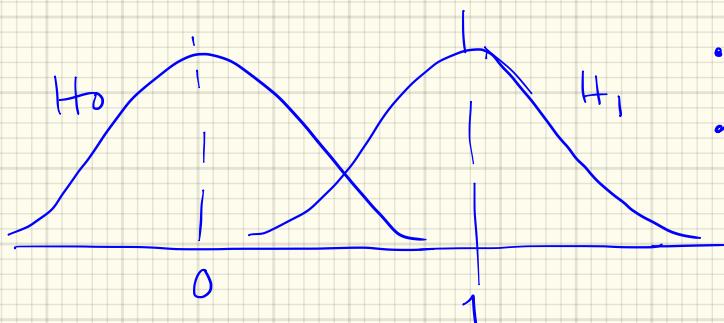
Extreme cases for  $\alpha$  &  $\beta$  probabilities,  
moving  $\Sigma$  to the right & left



Ex: <sup>Hypothesis</sup> Test on Normal Means :

n data points i.i.d.

You have 2 normal distributions w/ different means :



- $H_0 : X_i \sim N(0, 1)$
- $H_1 : X_i \sim N(1, 1)$

1) • Likelihood ratio test, <sup>Rejection region</sup>  $\leftarrow H_1$

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

$\leftarrow$   $\rightarrow$

$H_0$       Do some algebra to simplify to = exercise

1) LRT test: Reject  $H_0$  if :  $\sum_i X_i > \xi'$

a test "statistic":

↳ summarizes our measurements into a single number

$$(\xi' = \log s + \frac{n}{2})$$

(If this statistic is large, it's evidence to reject  $H_0$ !)

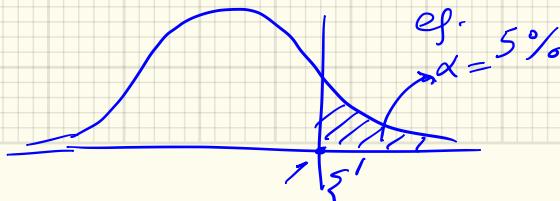
2) How to choose  $\xi'$ ? We set prob. of false rejection to a certain probability, e.g. 5%.

$$P\left(\sum_{i=1}^n X_i > \xi'; H_0\right) = \alpha$$

→ Sum of  $X_i$ 's have a Normal distrib. (blc  $X_i$ 's have a normal distrib.)

∴ Use Normal tables

$$\xi' = 1.96$$



if  $\sum_i X_i > 1.96$ ; reject  $H_0$ .  
 $< 1.96$ ; Not reject  $H_0$ .

Ex: Test on Normal Variances :

- n data points, i.i.d:  $H_0: N(0, 1)$
- $H_1: N(0, 4)$

} same mean but different variances

Intuitively, if  $X$ 's have wide-spread (variance) , choose  $H_1$   
 small " " choose  $H_0$ .

- LRT ; Rejection region :

$$\frac{\text{density Under } H_1}{\text{density Under } H_0} = \frac{\left(\frac{1}{\sqrt{2\pi}}\right)^n \exp\left\{-\sum_i X_i^2 / 2 \cdot 4\right\}}{\left(\frac{1}{\sqrt{2\pi}}\right)^n \exp\left\{-\sum_i X_i^2 / 2\right\}} > \{$$

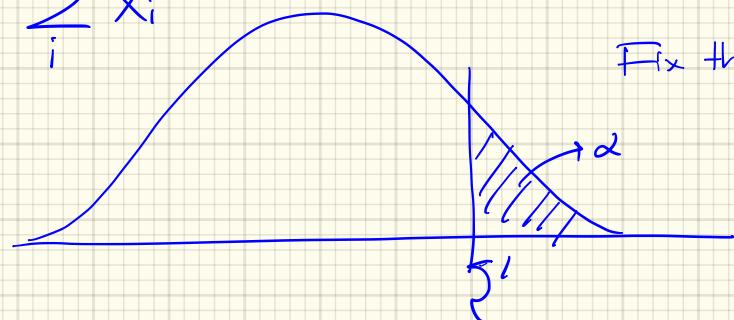
↳ Do algebra to simplify to:

Reject  $H_0$  if  $\sum_i X_i^2 > \xi'$



Find  $\sum_i^{} x_i^2$  s.t.  $P\left(\sum_i^{} x_i^2 > \sum_i^{} x_i^2 ; H_0\right) = \alpha$

Distribution of  $\sum_i^{} x_i^2$  is known:  $\chi^2$  (Chi-Squared distribution)  
 Derived distrib.  
 $\sum_i^{} x_i^2 \longrightarrow$  (Tables are available)



Fix the tail prob. of Chi-Squared  
 distrib.  
 e.g. to 95<sup>th</sup> percentile

→ From  $\chi^2$ -tables, read off  $\sum_i^{} x_i^2$  that corresponds to 0.95.

→ Note: Your "statistic"  $\sum_i^{} x_i^2$

## Composite Hypothesis:

e.g. you have a coin,

your question: is it fair / unfair?

- You make  $n$  tosses

→ You get  $S = 474$  heads in  $n = 1000$  tosses?

Is the coin fair?

~~binary hypothesis  
(not  $p=0.5$  vs  $p=0.6$ )~~

$$H_0 : p = \frac{1}{2}$$

(fair)

vs.

$$H_1 : p \neq \frac{1}{2}$$

(unfair)

$$\left. \begin{array}{l} p = 0.51 \\ 0.52 \\ 0.55 \\ 0.57 \\ 0.59 \end{array} \right\}$$

Many possible alternatives.

→ expected value:  $\frac{n}{2}$ : half heads, half tails ) here 500

(i)

Pick a statistic

: { H H T H T T T H . . . -

we did it

$$\rightarrow S = \# \text{Heads}$$

statistic

(ii) Pick shape of rejection region:  $|S - \left(\frac{n}{2}\right)| > \{$

500

iii) Pick a Significance Level  $\alpha$ : (e.g.  $\alpha = 0.05$ )

iv) Pick a threshold  $\xi$  (critical value) s.t.

$$P(\text{reject } H_0 ; H_0) = \alpha \underset{\approx}{\sim} \text{prob of outliers.}$$

Using CLT: # heads, i.e.  $S$  statistic is normal:

$$P(|S - 500| \leq \xi ; H_0) = 0.95$$

Use normal tables

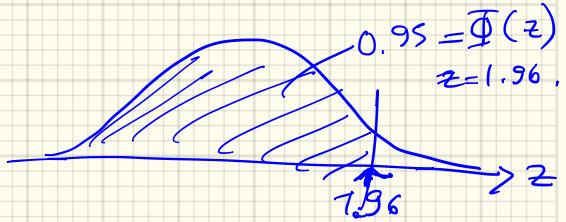
Recall:

$$-1.96 \leq \frac{S - 500}{\sqrt{\text{Var}(S)}} \leq 1.96$$

$$\sqrt{n \cdot \sigma^2} = \sqrt{1000 \cdot \frac{1}{4}} = 250$$

$$S - 500 \leq 1.96 \cdot \sqrt{250}$$

$$\xi = 31$$



$$\rightarrow |S - 500| < 1.96 \sqrt{250} \hat{=} 31 \Rightarrow \{ = 31$$

For our ex :  $\underline{S = 474} \rightarrow |S - 500| = 26 < \{ = 31$

$\rightarrow$  Do Not Reject  $H_0$ . (at the 5% level)

$\exists$  5% chance that the data we got is outlier.

---

Note : We tend to say

$H_0$  is not rejected rather than  $H_0$ : accepted.

=

$H_0$  is the default hypothesis: we do not reject it until we see evidence contrary to it.

If we some evidence to the contrary , we reject  $H_0$ .

Ex: Is your die fair?

Null hypothesis (fair die) : is a pmf:  $P(X=i) = p_i = \frac{1}{6}$ ,  $i=1, \dots, 6$

•  $N_i$ : observed occurrences for each  $i$ .

• Roll your die n times; count # 1's  $\rightarrow N_1$   
# 2's  $\rightarrow N_2$   
:  
# 6's  $\rightarrow N_6$

Q: We observe  $N_i$ 's; is my die fair?  
or is my pmf (above) valid?

• Choose form of rejection region:

Reject  $H_0$  if  $T = \sum_i \frac{(N_i - \underbrace{n \cdot p_i}_{\text{expected}})^2}{n \cdot p_i} > \xi \Rightarrow$

Under  $H_0$ :  
I expect  $N_i$ 's  
 $\approx \frac{n}{6} = N_i$ .  
 $= n \cdot p_i$ ,  $p_i = \frac{1}{6}$

• Choose  $\Sigma$  so that prob of false rejection 5%

$$P(\text{reject } H_0 ; H_0) = 0.05$$

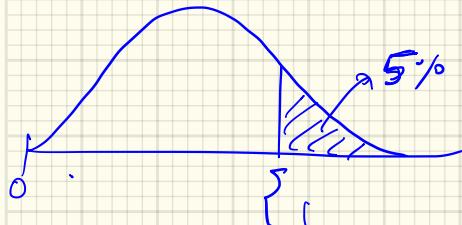
specifically

$$\equiv P(T > \Sigma ; H_0) = 0.05$$

→ We need distribution of  $T = \sum_i \frac{(N_i - n \cdot p_i)^2}{n \cdot p_i}$  derived distribution.

For large  $n$ ,  $T \sim$  a chi-squared distrib.

$P_T(+)$  under  $H_0$ .



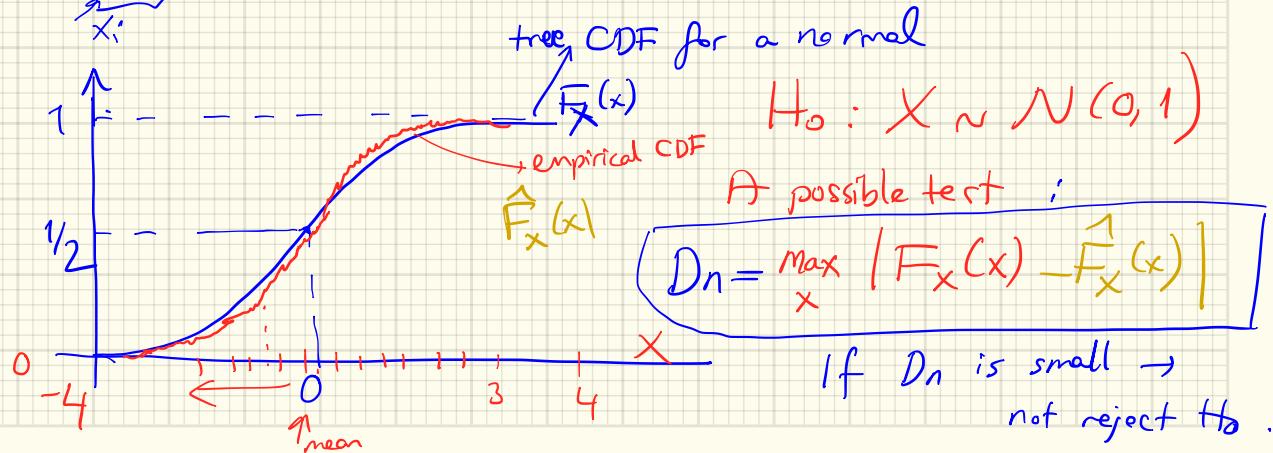
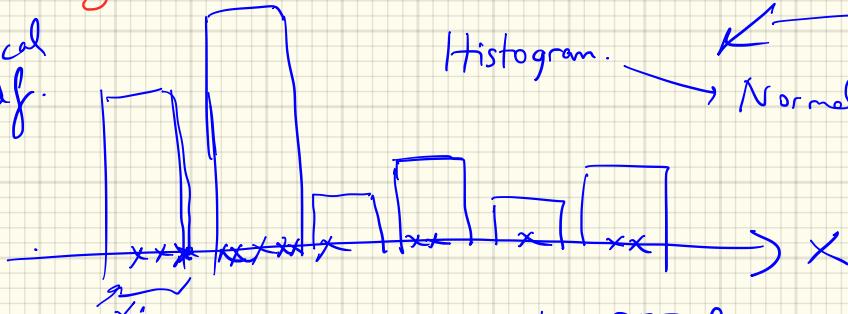
For distrib. of  $T \rightarrow$  there are tables for the  $T$ . Use them to find  $\Sigma$ .

Decide whether to  
Reject  $H_0$  or  
Not Reject  $H_0$ .

Now what happens

→ eg. When your data comes from a certain normal? what could be a good statistic to use in a test:

Kolmogorov - Smirnov Test: Form CDF (empirical)  $\hat{F}_X$  from empirical pdf.



$$\rightarrow P\left(D_n \geq \frac{1.36}{\sqrt{n}}\right) \approx 0.05$$

b/c this test is used frequently; people calculated distrib. of r.v.  $D_n$ .

From tabulated prob.-values of  $D_n$ :

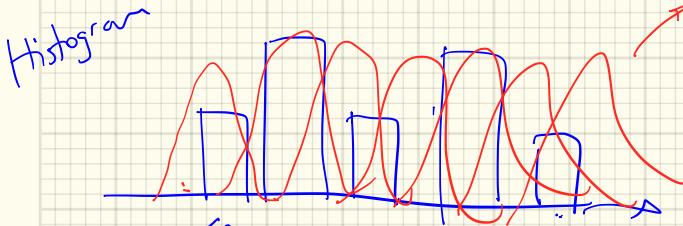
$$\xi = \frac{1.36}{\sqrt{n}}$$

$n$ : # data points .

If  $D_n \geq \xi$  → Reject  $H_0$  .

Note: Kernel Density Estimation:

Method to estimate unknown PDF from a histogram.  
opf smoothed Histogram estimates thru



weighted Gaussians.

$$h(\cdot)$$

+ Huge High-dim spaces of data  
+ Huge # parameters

→ Machine Learning  
Data mining

Deep Learning