YZV 231E

26.12.2022

Probability Theory & Stats

Week 14

Gü.

Recap:

**MLE** : Maximum Likelihood Estimation :      Note: $\theta$ is not an r.v.

Model w/ unknown parameters $\underline{\theta}$ ;

$$X \sim P_X(x; \underline{\theta})$$

a family of parameterized models

$$P_X(x; \theta_1)$$
$$P_X(x; \theta_2)$$
$$\vdots$$
$$P_X(x, \theta_N)$$
$$\vdots$$

some possibilities.

MLE picks $\theta$ that makes the data

most likely : $\arg\max_{\underline{\theta}} P_X(x; \underline{\theta}) = \hat{\underline{\theta}}$.

Compare to Bayesian approaches to estimation:   $\theta$ is an r.v.

→ **MAP** : $\arg\max_{\theta} P_{\Theta|X}(\theta|x) = \hat{\theta}$      $P_{\Theta|X}(\theta|x) = \dfrac{P_{X|\Theta}(x|\theta) \, P_\Theta(\theta)}{P_X(x)}$

likelihood   prior

find $\theta$ most likely under the posterior distrib.

Note: MLE & MAP appear the same when we have a uniform prior, but in principle they are very different.

→ **LMS** $\doteq E[\Theta|X] = \hat{\theta}$

estimator   $\Theta|X$

— **Sample Mean Estimator** of $\theta$

r.v. $\boxed{\hat{\Theta}_n = \dfrac{X_1 + \cdots + X_n}{n}}$ : $\boxed{\text{point}}$ estimator.

$\rightarrow$ Properties :    Unbiased   ,    Consistent   , "small" MSE

of an estimator    $E[\hat{\theta}] = \theta$      $\hat{\theta} \underset{\text{in prob}}{\longrightarrow} \theta$     $\approx \mathrm{Var}(\hat{\theta}) + (\mathrm{Bias})^2$

Bias $= \hat{\theta} - \theta$.

— $(1-\alpha)$ Confidence Interval (CI)

$$P\left( \hat{\Theta}_n^- \leq \theta \leq \hat{\Theta}_n^+ \right) \geq 1-\alpha \quad , \quad \forall \theta.$$

      $\uparrow$                $\uparrow$         $\underbrace{\phantom{1-\alpha}}_{\text{g. } 95\%}$

    r.v.            r.v.

Construction of the CI : w/ $\boxed{\text{CLT}}$   pick an $\alpha \rightarrow$ fix $z$.

Confidence Interval for the Sample mean $\hat{\Theta}_n$

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

           $\underbrace{\phantom{xxx}}_{\hat{\Theta}_n^-}$              $\underbrace{\phantom{xxx}}_{\hat{\Theta}_n^+}$     $\underbrace{\phantom{xx}}_{95\%}$   $\alpha = 0.05$

where $z$ is s.t $\Phi(z) = 1 - \dfrac{\alpha}{2} = 0.975$
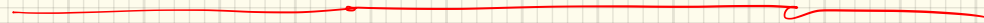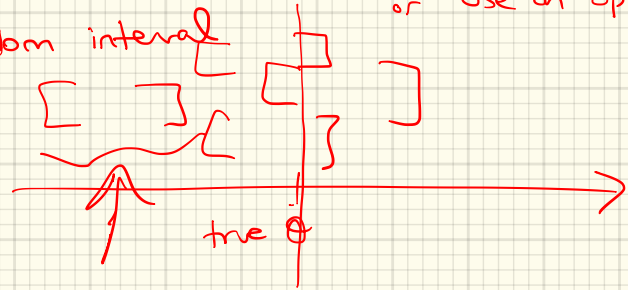
— need an estimate of the variance $\sigma^2$ ; either use sample variance, or use an upper bound if any.

CI: $\left[\hat{\Theta}_n^-, \hat{\Theta}_n^+\right]$ : random interval

r.v.   r.v.

true $\theta$
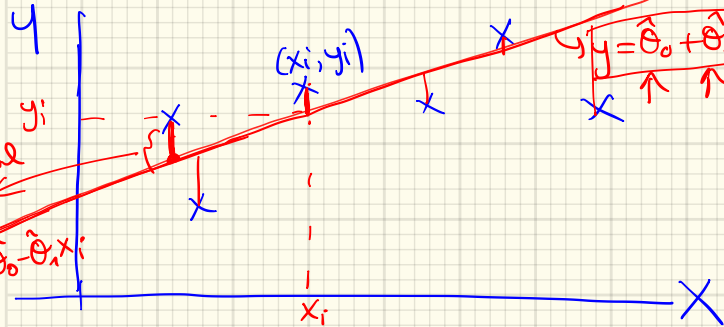
# REGRESSION:

Let $X$ be your **TYT exam score**

Let $Y$ be your **ITU GPA**

Q. Is there a <u>relation</u> between the two?



$$Y = \hat{\theta}_0 + \hat{\theta}_1 x$$

residual error$_i$

$y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i$

**Goal**: Find the "best" model to explain the data

Always ask: "optimal" or "best" w.r.t. which <u>criterion</u>? (measure)

**Data**: $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n) = \{x_i, y_i\}_{i=1}^{n}$

**Model**: $y_i \approx \boxed{\theta_0} + \boxed{\theta_1} x_i$

↑ unknown ↑ parameters that define our model.

Minimize the Residual error: $y_i - \theta_0 - \theta_1 x_i$:

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

residual error

(A) Cost = sum of squared errors in predictions.

Probabilistic Interpretation:

model the GPA score

$$\rightarrow \quad y_i = \theta_0 + \theta_1 x_i + \boxed{w_i} \quad , \quad \underline{w_i \sim N(0, \sigma^2)}$$

random noise

$w_i$'s are indep. $\forall i$.

choose a specific probabilistic model.

Want to do $\boxed{MLE}$ estimation:

Write a $\underline{likelihood}$ fn. $\quad P_{y,x \mid \theta}(x_i, y_i; \theta)$

$\sim$ probability

likelihood of $w$

$$w \sim c\, e^{-w_i^2/2\sigma^2}$$

$\downarrow$ write this for all samples $y_1, \ldots, y_n$.

a sum b/c $w_i$'s are independent.

normalization coeff.

likelihood of $y$

$$P(y; x, \theta) \sim c \exp\left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^{n} \underbrace{(y_i - \theta_0 - \theta_1 x_i)^2}_{w_i^2} \right.$$

maximize w.r.t. $\theta_0, \theta_1$: you can take a logarithm

$\rightarrow$ this is the same cost as in (A) previous page.

∴ Linear Regression $\equiv$ MLE where $\underline{w_i \sim N(0, \sigma^2)}$

i.i.d.

# Linear Regression:

model: $\quad y \approx \theta_0 + \theta_1 x$

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

Solution: set the derivatives of the cost function to zero:

exercise, derive these $\hat{\theta}_0$ & $\hat{\theta}_1$ expressions.

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

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

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

Covariance $X, Y$

$$\approx E\left[ (X - E[X])(Y - E[Y]) \right]$$

(A)

or through a probabilistic interpretation:

Our model: $Y = \theta_0 + \theta_1 X + W$ , $\quad$ X & W are independent

w/ zero mean.
(for simplicity)

$$E[Y] = \theta_0 + \theta_1 E[X] + 0$$

$$\theta_0 = E[Y] - \theta_1 E[X]$$

$$\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \cdot \bar{x} \qquad \checkmark \quad \text{given } \hat{\theta}_1 \text{ already}$$

Now, obtain $\hat{\theta}_1$ (estimate):

$$Y \cdot X = \theta_0 \cdot X + \theta_1 X^2 + W \cdot X$$

$$E[Y \cdot X] = \theta_0 \cdot E[X] + \theta_1 \cdot E[X^2] + \underbrace{E[W \cdot X]}_{0} \quad \overset{\nearrow E[W] \cdot E[X]}{}$$

for zero mean r.v.s

$$= \text{Cov}(X, Y) = \theta_1 \cdot \text{Var}(X)$$

$$\rightarrow \theta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \longrightarrow$$

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

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

# → Multiple Linear Regression:

eg. include more variables that may affect your ITU GPA:
your high school GPA, years of education of parents, ?
? family income, ---or such

Data: $(X_i^{(1)}, X_i^{(2)}, X_i^{(3)}, Y_i)$

↳ high school gpa    tyt.    educ.yrs of parents    ↳ itu gpa

not squared!
superscript, not to the power!

Model: $\boxed{y_i \sim \theta_0 + \theta_1 x_i^{(1)} + \theta_2 x_i^{(2)} + \theta_3 x_i^{(3)}}$ : linear fn. of all the variables

multiple explanatory variables in our model.

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

↓ take derivatives w.r.t. $\theta_0, \theta_1, \theta_2, \theta_3$, set to 0, you get a system of linear equations.

In Matrix notation : you can get a closed form solution.

Digression: In vector notation : set $\underline{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}_{\substack{n\times 1 \\ 4\times 1}}$  $\underline{x_i} = \begin{bmatrix} 1 \\ x_i^{(1)} \\ x_i^{(2)} \\ x_i^{(3)} \end{bmatrix}_{\substack{n\times 1 \\ 4\times 1}}$

m data points $\{\underline{x_1}, \underline{x_2}, \dots \underline{x_m}\}$ :

$$\min_{\underline{\theta}} \| \underline{y} - \underline{\underline{X}}\,\underline{\theta} \|^2$$

vector — matrix vector

$$\underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}_{m\times 1}$$

m: #data instances

$$\underline{\underline{X}} = \begin{bmatrix} - \underline{x_1^T} - \\ - \underline{x_2^T} - \\ \vdots \\ - \underline{x_m^T} - \end{bmatrix}_{m\times n}$$  $1\times n$

$$f(\underline{\theta}) = \left( \underline{y}^T - (\underline{\underline{X}}\,\underline{\theta})^T \right)\left( \underline{y} - \underline{\underline{X}}\,\underline{\theta} \right)$$

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

take deriv. w.r.t. $\underline{\theta}$

$$\nabla_{\underline{\theta}} f = \begin{bmatrix} \dfrac{\partial f}{\partial \theta_0} \\ \dfrac{\partial f}{\partial \theta_1} \\ \vdots \\ \dfrac{\partial f}{\partial \theta_n} \end{bmatrix} = 0$$
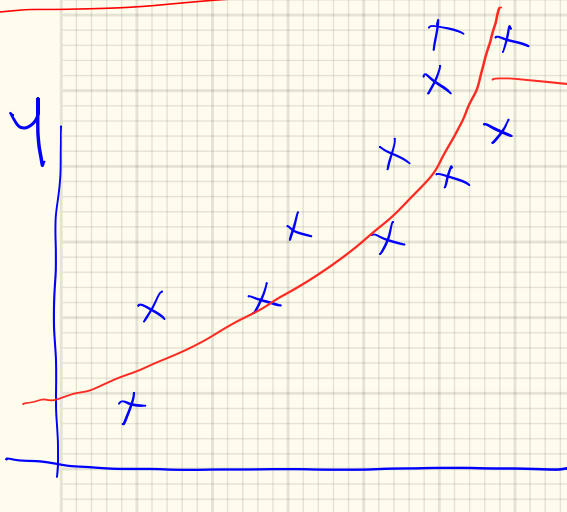
$$-2\,\underline{y}^T \underline{\underline{X}} + 2\,\underline{\underline{X}}^T \underline{\underline{X}}\,\underline{\theta} = 0$$
$$= \underline{\underline{X}}^T \underline{y}$$

$$\boxed{\underline{\theta} = \left( \underline{\underline{X}}^T \underline{\underline{X}} \right)^{-1} \underline{\underline{X}}^T \cdot \underline{y}}$$

a closed form solution for multiple linear regression.

$$\underline{\theta}_{n \times 1} = \left( \underline{\underline{X}}^T \underline{\underline{X}} \right)^{-1}_{n \times n} \underline{\underline{X}}^T_{n \times m} \underline{y}_{m \times 1} \quad : \quad \text{Normal Equations}$$

gives us the $\underline{\theta}$ vector estimates for linear regression.

$$\underline{\underline{X}} = \begin{bmatrix} \underline{x_1}^T \\ \vdots \\ \underline{x_m}^T \end{bmatrix} \begin{array}{l} \text{row} \\ \text{vector} \\ \text{for each} \\ \text{data instance} \end{array}$$

$m \times n$.

not a linear
but a quadratic
model of the measurements.

$$y \approx \theta_0 + \theta_1 \underbrace{( \,\times\!\!\!\times\, )}_{h(x)}$$

Now, use : nonlinear functions of
the data.

still, this is a linear regression.

Model: $\quad y \approx \theta_0 + \theta_1 \boxed{h(x)}$

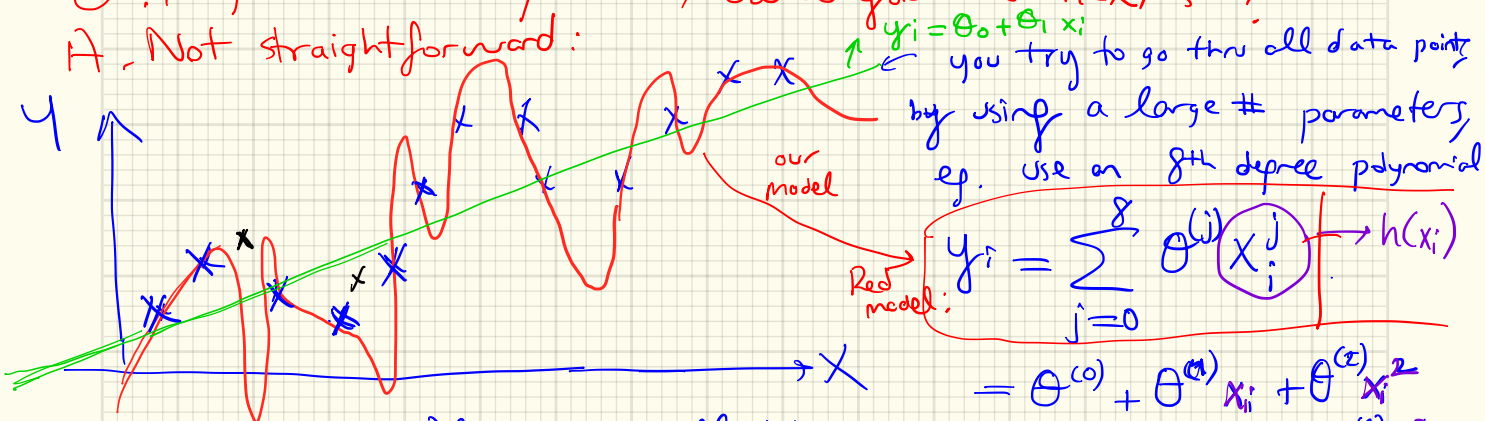eg. $\quad y \approx \theta_0 + \theta_1 x^2 \Leftarrow$

→ Same formulation

• Data points : $(y_i, h(x_i))$

• model : $y_i = \theta_0 + \theta_1 h(x_i)$ , $\forall i = 1, \cdots m$

$$\min_{\theta} \sum_{i=1}^{m} \left( y_i - \underline{\theta_0} - \underline{\theta_1} h(x_i) \right)^2$$

$(\theta_0, \theta_1, \ldots)$

fr. of $X$: e.g. $\boxed{x_i^3}$ → $3^{rd}$.

$\left.\begin{array}{c}\end{array}\right\}$ still a linear model.

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

A. Not straightforward :



$y_i = \theta_0 + \theta_1 x_i$

← you try to go thru all data points by using a large # parameters, eg. use an 8th degree polynomial

our model

Red model : $y_i = \sum_{j=0}^{8} \theta^{(j)} \left( x_i^{j} \right)$ → $h(x_i)$

$= \theta^{(0)} + \theta^{(1)} x_i + \theta^{(2)} x_i^2$

$+ \theta^{(3)} x_i^3 + \cdots + \theta^{(8)} x_i^8$

$\left( y_i - \sum_{j=0}^{8} \theta^{(j)} x_i^{j} \right)^2$ : error is small b/c we have lots of parameters.

Q : Is the red model a good model ? No!

→ Your model cannot generalize to a new data point well!

→ Overfitting problem!

→ Choosing complex $h(.)$ vs simple $h(.)$ ?

( Q. How complex ? How many explanatory variables ?

( open & extensive research topic,

— When you have a ⟨few data points⟩, avoid using too many parameters in your model.

→ Good rule: Start w/ simpler models ≡ a few parameters
especially when you have a few data points . Gradually increase later w/ more data etc.

⭐

Notes | For these $\theta_i$ , people also report confidence interval
(not covered in this class)

— $R^2$ : measure of explanatory power of the model in your linear regression.

— Standard error estimates of $\sigma^2$

$y = \theta_0 + \theta_1 X + W \longrightarrow \sigma^2$ ; variance of the noise

uncertainty in the model .

related to $R^2$

$$\boxed{\frac{Var(Y|X)}{Var(Y)}} < 1 \quad \text{naturally}$$

adding knowledge of $Y$, how much of the randomness in $Y$ is reduced.

If this is small, including $X$ as my explonatory variable for the target $Y$ variable helps /improves our predictions on $Y$.

60% of the student's ITU GPA is explained by the TYT score.

## Some Pitfalls in Using Linear Regression :

### * Heteroskedasticity :



$\sqrt{}$ A linear ~~model~~ model that you fit to this data.

A good fit in Region I.

model
I ) has small variance
II) has large variance.

Be careful w/ this problem. (not covered) in this class.

- eg.
- Need to incorporate varying $Var(W)$.

I ; small errors

II ; large errors

**✳ Multi _ Collinearity:** Multiple explanatory variable; they are closely w/ each other.

eg. model $y_{GPA} = \theta_0 + \theta_1 \, TYT_1 + \theta_2 \, TYT_2$

$\underbrace{TYT_2}_{AYT}$

Your TYT & AYT (2 exam scores) are close to each other.
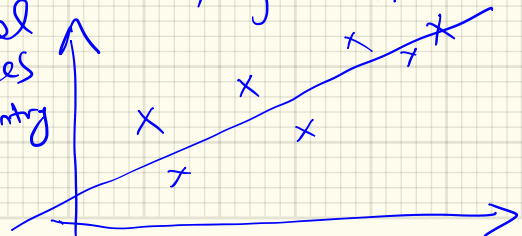
→ Correlated → Redundancy.

$y_{GPA} = \theta_0 + \theta_1 \, TYT_1$

$y_{GPA} = \theta_0 + \theta_2 \, TYT_2$

⎫ avoid such redundancy in explanatory variables.
b/c they create sensitivity of the model to small changes in the data.

**✳ Causality** ; Do not use (linear) regression to conclude causality!



#Nobel prizes in a country ↑ ... Chocolate Consumption

Never say that $y$ is CAUSED by $X$ according to your linear regression model

# Hypothesis Testing:

We have a null $H_0$ $\cancel{X}$ and an alternate $H_1$ hypothesis.
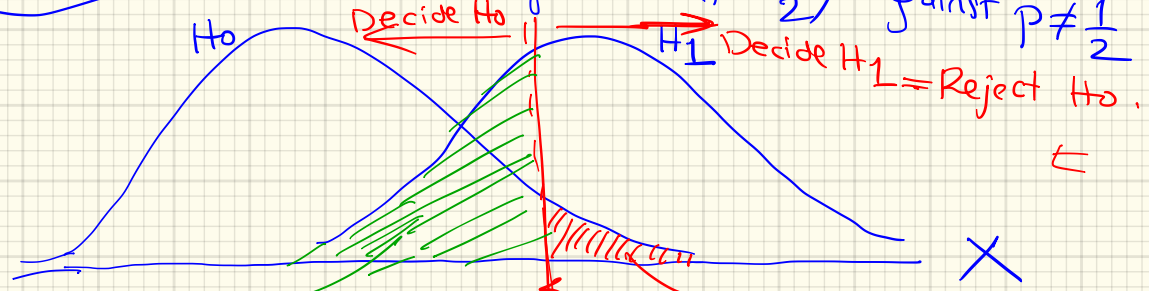
Coin; Fair or not;

$$H_0 : p = \frac{1}{2} \quad vs. \quad H_1 : p \neq \frac{1}{2}$$

single alternate hypothesis

or

$$H_1 : \begin{array}{l} p = 0.55 \\ p = 0.6 \\ p = 0.65 \end{array}$$

many alternate hypothesis

**X-space:**

Decide $H_0$

Reject $H_0$ (Decide $H_1$)

Test coin's fairness $\left(p = \frac{1}{2}\right)$ against $p \neq \frac{1}{2}$

Decide $H_0$ ← | → $H_1$ Decide $H_1$ = Reject $H_0$.

$H_0$

$\leftarrow$

Choose a threshold

$\alpha$; prob. of rejecting $H_0$ when it was true.

$\cancel{X}$

$\beta$: error I make for not rejecting $H_0$, actually when $H_0$ should have been rejected.

likelihoods

$P_X(X|H_0)$

$P_X(x|H_1)$

X

trade off:

$\beta \downarrow$ , $\alpha \uparrow$

$\beta \uparrow$ , $\alpha \downarrow$

$\beta$

$\xi_2$

$\xi_1$ $\ll$ move threshold to the right
to make $\alpha$ small

$\alpha$

move $\xi$ to the left
to make $\beta$ small.

* want both $\alpha, \beta$ to be small; however $\exists$ a trade-off.

Q. How to set the threshold $\xi$ ?

A [ Likelihood Ratio Test (LRT) ] ; Compare the
posterior prob. of
the hypothesis.

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

Pick the hypothesis which is more likely, given the data.

$$P(H_1 \mid X=x) > P(H_0 \mid X=x)$$

using Bayes

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

$$\frac{P(X=x \mid H_1)\, P(H_1)}{P(X=x)} > \frac{P(X=x \mid H_0)\, P(H_0)}{P(X=x)}$$

$$\Rightarrow L(x) \overset{\Delta}{=} \boxed{\frac{P(X=x \mid H_1)}{P(X=x \mid H_0)}} > \boxed{\frac{P(H_0)}{P(H_1)}} \quad (LRT)$$

likelihood ratio $\longrightarrow$ $\xi$ : threshold $\propto$ ratio of the prior prob. of the hypothesis. (Bayesian view)

compare to a threshold $\xi$.

$\uparrow$ Bayesian setting

— In a non-Bayesian setting: don't have prior probabilities

Still $\left\{ \dfrac{P_x(X=x \,;\, H_1)}{P_x(X=x \,;\, H_0)} > \xi \right.$   so we re-write i.t.o. pdf of $x$.

Q: If $\left( L \overset{this}{(x)} \right)$ ratio is large, $\boxed{Q.}$ Is it likely that my observations $X$ occurred under $H_0$?

No!

No ; it is unlikely that observations $X$ occured under $H_0$.

$\therefore$ Reject $H_0$.

— Threshold $\zeta$ trades-off 2 types of error:



Choose $\zeta$. s.t.

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

$$1 - CDF_X(\zeta) = \alpha.$$

We fix $\alpha$ , ep. $\alpha = 0.05 \rightarrow$ find $\zeta$ (threshold)

$\rightarrow$ then $\beta$ is already fixed.

Simple **Binary Hypothesis Testing**

Want to make a decision whether to **Reject** or **Not Reject** the null hypothesis

- (Default) Null Hypothesis $H_0$ : $X \sim P_X(x; H_0)$
  (r.v.)
- Alternative Hypothesis $H_1$ : $X \sim P_X(x; H_1)$

Designing the Hypothesis Test ( checking whether $H_0$ is false or not )

1) Structure of the test : <u>shape of the dividing curve</u>.

   ex. Likelihood Ratio Test :
   $$\frac{P_X(x; H_1)}{P_X(x; H_0)}$$

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



$X$-space of observations

Do Not Reject.

Reject $H_0$

1) LRT : Reject $H_0$ if :
$$L(x) = \frac{P_X(x; H_1)}{P_X(x; H_0)} > \zeta$$

2) How to choose $\zeta$ ?

**2)** $\xi$ **?**

Fix $\alpha$ $\longrightarrow$ choose $\xi$ so that

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



$H_0$

$H_1$

$\alpha = \text{prob. of false rejection}$

$X$.

$$L(x) < \xi :$$
Do not Reject $H_0$.

$$L(x) > \xi \text{ ; reject } H_0$$

eg. set $\alpha = 0.05$ $\quad (5\%)$ $\rightarrow$ that sets $\xi$.

Note:
(α, β) trade off

extreme cases for α, β probabilities

Ho ⟵

$\left( \alpha = 0, \beta = 1 \right).$

extreme

always decide $H_1$ → $\left( \beta = 0, \alpha = 1 \right)$

extreme

not rejecting Ho at all.

$\alpha = 0$
$\beta = 1$

$(\alpha, \beta)$ trade-off due to $\boxed{LRT}$

$(\alpha, \beta)$ → choosing $(\alpha, \beta)$ w/ method other than LRT.

$(\alpha, \beta)$ pair

$\alpha = 1$
$\beta = 0$

Note: Theoretically
For a given α value,
LRT minimizes
the probability β.

# Ex: Hypothesis Test on Normal Means

- $n$ data points, $X_i$: i.i.d. and normal

You have 2 normal distributions w/ different means



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

1) Likelihood Ratio Test ; Reject $H_0$ if:

$$\frac{P_X(x; H_1)}{P_X(x; H_0)} = \frac{\left(\frac{1}{\sqrt{2\pi}}\right)^n \exp\left\{-\sum_{i=1}^{n}(x_i - 1)^2/2\right\}}{\left(\frac{1}{\sqrt{2\pi}}\right)^n \exp\left\{-\sum_{i=1}^{n}x_i^2/2\right\}} \quad \gtrless$$

$\leftarrow H_1$

$\leftarrow H_0$

exercise : do some algebra to simplify to.

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

a test "statistic"

Summarizes our measurements into a single number

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

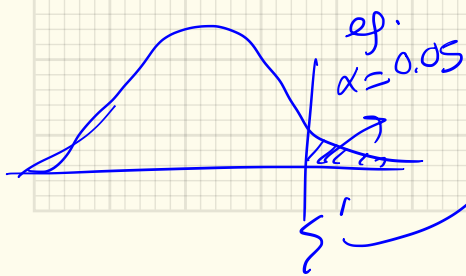Intuitive here, if $\sum_i X_i$ is large $\rightarrow$ evidence to reject Ho.

2) How to choose $\xi''$? Set prob of false rejection to a certain probability $\alpha$.

eg. 5%.

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

$X_i$'s normal $\rightarrow$ sum $X_i$s

$\sum X_i \rightarrow$ Normal distrib.

$\rightarrow$ Use Normal tables

eg.
$\alpha = 0.05$

$\xi' = 1.96$

$\xi'$

If $\sum_i X_i > 1.96$ ; Reject Ho.

$< 1.96$ ; Do Not Reject Ho.

Ex: Hypothesis Test on Normal Variances.

. $n$ data points $X_i$, i.i.d.  $H_0 : N(0, \boxed{1})$  } same mean but different variances.

$H_1 : N(0, \boxed{4})$

LRT : Rejection (of $H_0$) region:

$$\frac{\text{Density of data under } H_1}{\text{Density under } H_0} = \frac{\left(\frac{1}{\sqrt{2\pi}}\right)^n \exp\left(-\sum_i X_i^2 / 2(4)\right)}{\left(\frac{1}{\sqrt{2\pi}}\right)^n \exp\left(-\sum_i X_i^2 / 2(1)\right)} > \xi$$
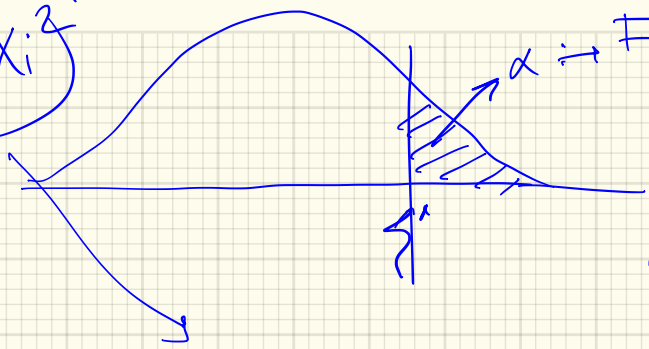
Do algebra to simplify to:

✱ | Reject $H_0$ if $\boxed{\sum_i X_i^2} > \xi'$ |

→ test statistic

✱ Find $\xi'$ s.t. $P\left(\sum_i X_i^2 > \xi' ; H_0\right) = \alpha$.

Distribution of $\sum_i X_i^2$ is known : $\chi^2$ ( Chi-squared distrib)

Recall: derived distrib.)

Tables are available

$\boxed{\sum_i X_i^2}$



$\alpha \rightarrow$ Fix the tail prob. of the $\chi^2$- distrib.
eg. to 95th percentile.

From $\chi^2$- tables,
read off the $\xi'$ that corresponds
to 0.95.

Note : Your "statistic" : $\sum_i X_i^2$

If $\sum_i X_i^2 > \xi'$ : Reject Ho

$\leq \xi'$ : Do not Reject Ho.

# Composite Hypothesis :  eg. coin → is it fair or unfair ?

→ You make $n$ tosses of the coin.

You get $S = 474$ Heads in $n = 1000$ tosses ?

Is the coin fair ?

$$H_0 : P = \frac{1}{2}$$
(fair)

vs $H_1 : P \neq \frac{1}{2}$
(unfair)

$\left. \begin{array}{c} P = 0.51 \\ 0.52 \\ : \\ : \end{array} \right\}$

Expected value : $\frac{n}{2}$ : half heads half tails

(i) Pick a statistic :    H H H T T H T T T - - - - -
                                                    1000

come up 3/

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

: Design/pick a statistic
≡ reasonable summary of your data

(ii) Pick shape of the rejection region :

1) ≡ Decide how to make your decision.

$$\left| S - \frac{n}{2} \right| > \xi$$   : Reject the hypothesis.

expected value
eg. 500

iii) Pick a significance level $\alpha$ (eg. $\alpha = 0.05$)
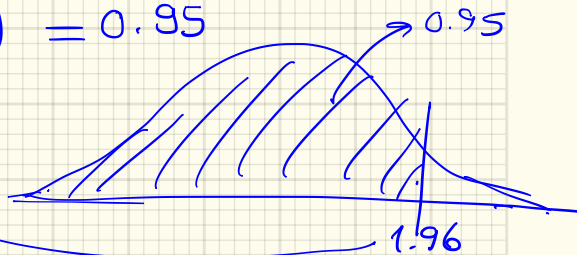
iv) Pick a threshold $\zeta$ s.t.

$$P(\text{reject } H_0 ; H_0) = \alpha \approx \text{probability of } \boxed{\text{outliers.}}$$

Using CLT : #heads , ie. S statistic is Normal.

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

0.95

From the Normal table:

$$\Phi(z) = 0.95 \rightarrow \dot{z} = 1.96.$$

1.96

normalize S

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

$\underbrace{\qquad}_{n \cdot \sigma^2}$

$$1000 \cdot \frac{1}{4} = 250$$

$\hookrightarrow$ use an upper bound on $\sigma^2$. (recall Bernoulli)

$$\boxed{S - 500 \leq (1.96) \, 250 \approx \boxed{31} = \zeta}$$

Test:

$\rightarrow$ $|S - 500| \leqslant 1.96\sqrt{250} \cong 31 = \xi$

For our ex. $S = 474 \rightarrow |S - 500| = 26 < \xi = 31$

$\rightarrow$ Do not Reject Ho (at the 5% level of error.)

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

Note: Say Ho is Not Rejected rather than (Ho: accepted)

Ho : default hypothesis $\longrightarrow$ we do not reject it until we see evidence contrary to Ho.

**Ex:** Is your die fair? $i = 1, \ldots, 6$.

$H_0$: is a pmf. : $\boxed{P(X = i) = p_i = \frac{1}{6}}$

Null Hypothesis → fair die.

• For each $i$ $(1, \ldots, 6)$: $N_i$ :# occurences for each

Roll your die $\boxed{n \text{ times}}$, count
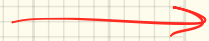
#1's → $N_1$
#2's → $N_2$
⋮
#6's → $N_6$

You Observe $N_i$'s : Is your die fair?

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

1a) Choose a form of Rejection region.

Reject $H_0$ if $T = \sum_i \frac{(\overbrace{N_i}^{observed} - \overbrace{n \cdot p_i}^{expected})^2}{n \cdot p_i} > \zeta$

(2) Choose $\zeta$ so that prob of false rejection 5%.

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

$$\downarrow P(T > \zeta ; H_0) = 0.05$$

We need distrib. of
Test statistic : $\longrightarrow$ $T = \sum_i \dfrac{(N_i - n p_i)^2}{n p_i}$ $\Longleftarrow$ derived distrib.

For large $n$, $T \sim$ a chi-squared distribution
(widely used in statistical tests)
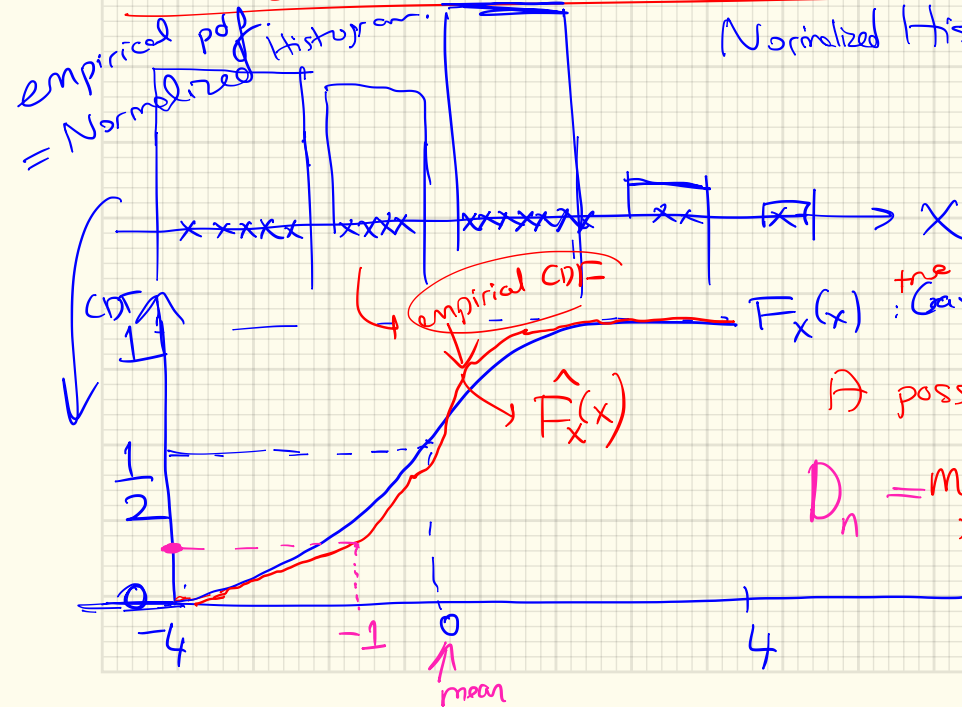
$\exists$ tables $\rightarrow$ set $\zeta$. (threshold)

$P_T(t)$



5%

$\zeta$           $t$

Reject $H_0$ if $T > \zeta$
Do not Reject $H_0$ if $T \leq \zeta$.

$\rightarrow$ Decide whether to Reject $H_0$ or Not !

**Ex:** Want to Test whether your data comes from a certain Gaussian distrib?

→ Kolmogorov - Smirnov (KS) Test ; From CDF (empirical)

Normalized Histogram ($\approx$ pdf) → CDF.

empirical pdf = Normalized Histogram.



$\rightarrow x$

empirical CDF

$F_x(x)$ : true Gaussian CDF

$\hat{F}_x(x)$

$H_0: X \sim N(0,1)$

A possible (KS test)

$$D_n = \max_x \left| F_x(x) - \hat{F}_x(x) \right|$$

If $D_n$ is small

→ Do Not Reject $H_0$.

CDF

1

$\frac{1}{2}$

0

$-4 \qquad -1 \quad 0 \qquad\qquad 4$

↑ mean

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

KS test is frequently used $D_n$ has a known
calculated distrib. $\longrightarrow$ tabulated   prob values of $D_n$.

$$\xi = \frac{1.36}{\sqrt{n}} \qquad (n : \#\,data\ points.)\ w/\ 5\%\ \text{rejection prob.}$$

If $D_n \geqslant \xi \longrightarrow$ Reject $H_0$.

_____

THE   END.

$\rightarrow$ Now, you have learned the basics of statistical
methods, any statistically-literate engineer/scientist should
know about.