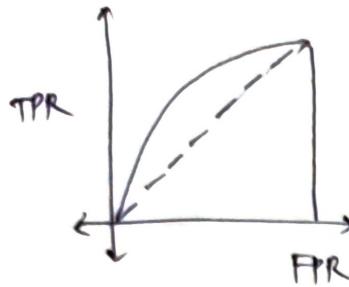


AUC-ROC

$\rightarrow \text{AUC} = 0.9$; what does it mean?



[ans] if we have 2 pts of different classes

and we don't know which is +ve and -ve,

$\text{AUC} = 0.9 \rightarrow 90\%$ it will classified correctly

$\text{AUC} = 0.5 \rightarrow \text{random model}$

[if $\text{AUC} < 0.5$] $\rightarrow \text{AUC} = 0.1$

swap -0 and +1

$\rightarrow \text{AUC} = 0.1$

$\rightarrow \boxed{\text{AUC} = 0.9}$

* In KNN why Manhattan > Euclidean in high dimensionality

Euclidean = $(x_1 - x_2)^2$ \rightarrow as sq \rightarrow it amplifies difference / distance

Manhattan = $|x_1 - x_2|$

etab. Manhattan \neq Euclidean

Θ 3D - 2D

to few z iterw.

high dimensional problems

short

dong web rebiese f1tab = 0.002

0.0005 - if hui elbrot f1tab =

program at kew sd f1tab = 0.0002

elbrot = 0.0002

\rightarrow short ab kewsd

Type I = sig level

$$\frac{d+s}{s} = (ds)$$

Show if norm standard

$$\frac{d+s}{s} = (ds)$$

more & big ben w: f1tab

more

f1tab \rightarrow (1000, 0.9) dim = 2000

$\frac{8AF}{8+A}$ si f1tab

more & siordaplo tar f1tab
of ben

more & siordaplo tar f1tab

$$\left(\frac{8AF}{8+A} \right)$$

(dis) norm standard:

$$1 - \left[\frac{\left(\frac{1+A}{A} \right)}{2} \right]$$

Hypothesis Testing

1) null hypothesis - ideal

2) Alternate hypothesis - opp of null

take sample from popn & calculate statistic

\rightarrow nabor ①

$$p\text{val} = P(\text{Obs} | H_0 \text{ is true})$$

\rightarrow p-value ②

p-value

H_0 : time spent = 20 min

H_a : μ time spent ≥ 20 min

$$\text{populn } \mu = \bar{\mu} = 25$$

take a sample

$$\text{of } N = 100$$

$$P\text{.val} = P(\bar{\mu} = 25 | H_0)$$

$$25.0 - 9.0$$

prob that our observation occurs given H_0 is true

$\alpha \rightarrow$ significance level

$$\underline{25.0 - 22.0} = P$$

if $P < \alpha$ H_0 rejected

$$P \geq \alpha : H_0 \text{ accepted}$$

generally $\alpha = 0.05$

Type I and Type II error

$$P(\text{Type I error}) = \alpha$$

	H_0 is true	H_0 is false
Reject H_0	TYPE I error	correct
Fail to Reject H_0	Correct	TYPE II error

reality

actual True \rightarrow Type I error

actual False \rightarrow Type II error

Type I : actual is True, you say False

Type II : actual is False, you predict True

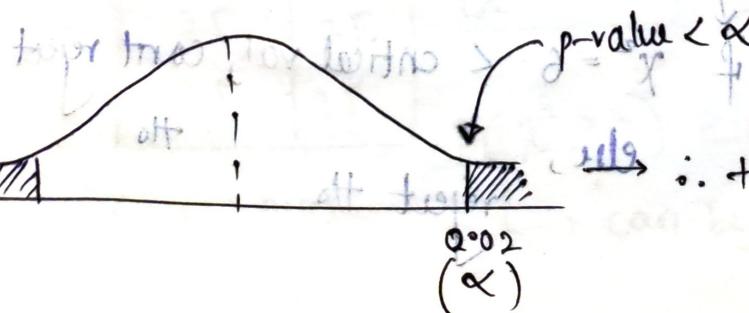
Type I : reject H_0 when it was True

Type I : reject H_0 when its true

\rightarrow assume H_0 is true and $\alpha = 2.1\%$

Type II : accept H_0 , when its false

$\epsilon = \text{nabor}$



\therefore there is 2.1% chance that you made on Type II error

* conditions for picking sample from popn

① Random

② Normal — at least 10 success and failure

③ Independence — size(sample) $\leq 10\%$ of population

z-table

$$x \rightarrow p(z \leq x)$$

L-value in z-table

problem

$$\text{spoke} \geq 1 \text{ lang} \quad \text{significance level} \rightarrow \alpha = 0.05 \\ H_0: p = 0.26$$

$$\text{sample} \rightarrow 120 \text{ ppn} = n \quad 10$$

signif. test abt proportion

with $H_1: p > 0.26$, left darg (40/120) speaks more than one language

alt: $z_i < 0$ or $z_{i+1} > 0$

$$(z_i - \mu) = \frac{\sigma}{\sqrt{n}}$$

$\bar{z} = \text{obtained prop} - \text{assumed proportions}$

$$\therefore 20.0 = \frac{p_0(1-p_0)}{n} \quad \text{std} \quad 0 \leq z \leq 9$$

p_0 — assumed popn proportion

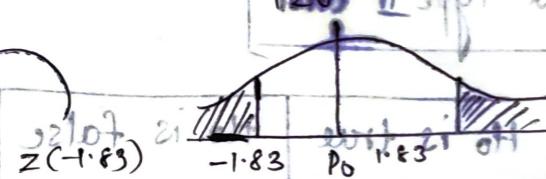
deviation from p_0

$$z = \frac{0.33 - 0.26}{\sqrt{0.26(1-0.26)}} = 1.83$$

$$\sigma = \sqrt{(p_0)(1-p_0)/n}$$

$$\therefore p\text{-value} \approx P(z \geq 1.83) = 0.0336$$

significant test abt mean



for evn. evn. \Rightarrow I sppt

CHI-S²

check off obv. the theory fits observations

→ theory: equal probs = 25, 25, 25, 25 (out of 100 students)

evn. evn. red obv. \Rightarrow I sppt

$$\chi^2 = 6 \quad \text{node} = 0.05 \text{ sppt} \quad \text{freedom} = 3$$

critical val for our χ^2 table

$\chi^2 = 6 >$ critical val, can't reject H_0

else,

reject H_0

H_0 : our theory is correct

H_1 : theory is not correct

Chi-square Test

χ^2 is distribution of distributions (picked from normal distribution)

degrees of freedom = $n - k$ (number of terms)

$X_i \sim \chi^2_{df=1}$ [if $x_i \sim N(0,1)$]

$\chi^2 = \sum X_i^2$ (for $n = 1, 2, \dots, n$)

$\chi^2_0 = \sum (x_i - \bar{x})^2 / s^2$

degrees of freedom (no of terms) = $n - 1$

chi-table	df	0.99	0.95	0.05	0.01
pdf \rightarrow critical	1	2.706	3.841	6.635	9.890
cdf \rightarrow p-value	2	5.991	7.879	10.828	13.819
df = $n - k$	3	7.815	10.828	12.834	16.812
df = $n - k$	n	10.828	12.834	15.812	19.812

Example: If you take a MCQ exam with 4 options

for χ^2 of degree = 3,

$P(X^2 \geq 7.81) \approx 0.05$

- (H₀) equal prob of choosing
- (H_a) : unequal

$$\chi^2 = \sum_{i=1}^n \frac{(obs - expected)^2}{expected}$$

sample size = 100

find probability of getting this if given H₀ is true

choice	expected	obs
A	25	20
B	25	20
C	25	25
D	25	35

chi-stat χ^2

$$\chi^2 = \frac{(20-25)^2}{25} + \frac{(20-25)^2}{25} + \frac{(25-25)^2}{25} + \frac{(35-25)^2}{25}$$

$$\chi^2 = 6$$

$$\text{degree of freedom} = (n-1) = 3$$

$$P(\chi^2 \geq 6) \approx 0.1$$

$$\alpha = 0.05$$

\hookrightarrow can be taken as p-value

then we get A

mean Est

Estimation

point estimation

Interval estimation

Q1 company sells pens over next month at 370.16 and $\sigma = 75$

Estimate mean demand for next month with 95% confidence

$$\rightarrow P\left(\bar{x} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

$\alpha = \text{significance level} = 0.05$

$$1 - \alpha = 95\% \Rightarrow \alpha = 0.05$$

$$Z_{0.025} = 1.96; 370.16 \pm 1.96 \cdot \frac{75}{\sqrt{100}} \Rightarrow 340.16 \leq \bar{x} \leq 399.96$$

chi-square independence test → check if 2 variables r independent

goodness of fit → check theory true

* 2 categorical variables

Ex:- a group of 120 people are given to choose their favorite social media app (i.e. social media app preference independent of gender)

H₀: gender & app are independent

	male	female	total
FB	15	20	35
INSTA	30	35	65
TIK	(15-25) + (30-25) + 20	(25-20) + (35-25) + 20	120
total	25	70	90

Female * Female

Total

Male * Female

Total

Male * Male

Total

Female * Male

Total

	male	female	total
FB	14.6	20.4	35
INSTA	27.1	32.9	60
TIK	8.3	11.7	20
total	50	70	120

OBSERVED

EXPECTED

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

$$\text{if independent, } \chi^2 \sim \chi^2_{(m-1)(n-1)}$$

$$\text{degree freedom} = (3-1)(2-1) = 2$$

$$20.0 = \frac{35 \times 50}{120}$$

frequencies analysis

Chebyshov's Theorem

* if data is gaussian distribution

$$\begin{aligned} \mu \pm \sigma &\rightarrow 68\% \\ \mu \pm 2\sigma &\rightarrow 95\% \\ \mu \pm 3\sigma &\rightarrow 99.7\% \end{aligned}$$

* if data is not normally distributed

At least $\left(1 - \frac{1}{K^2}\right)$ values will fall within $(\mu + K\sigma)$ and $(\mu - K\sigma)$

$(\mu \pm K\sigma)$ contains $\left(1 - \frac{1}{K^2}\right)$ proportion of values

sample μ and σ^2 → why $(n-1)$ in D^r ?

because when we take sample, it was found

$$\mu = \frac{\sum x_i}{n-1}$$

that, it was underestimated

$$\sigma^2 = \frac{\sum (x - \bar{x})^2}{n-1}$$

∴ reduce σ^2 by 1 to compensate

Z scores it represents no of SDs a value $\textcircled{2}$ is above or below the mean when dist^r is normal.

$$Z = \frac{x - \mu}{\sigma}$$

Ex: $Z = +2$ shows that value $\textcircled{2}$ is 2 SDs above mean

curve tail - x_{od} to abscissa of normal curve
curve right - x_{od} to abscissa of normal curve



* used to compare SD's with different means

Coef of variation

$$CV = \frac{\sigma}{\mu} (100)$$

The SD is $(CV)^{100\%}$ of the mean = $(80\%)^9$

$$L = (8)9 + (7)9 = (87.829) + (80.97)$$

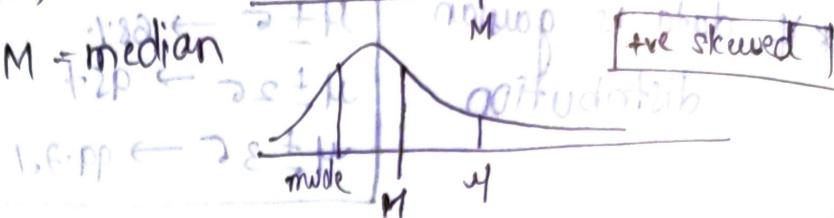
↳ relative component of SD to its mean

* Coef of skewness =

$$S_k = \frac{3(\mu - M)}{\sigma^3}$$

where,

M = median



Kurtosis

* Amount of peakedness of a distribution

Quartiles

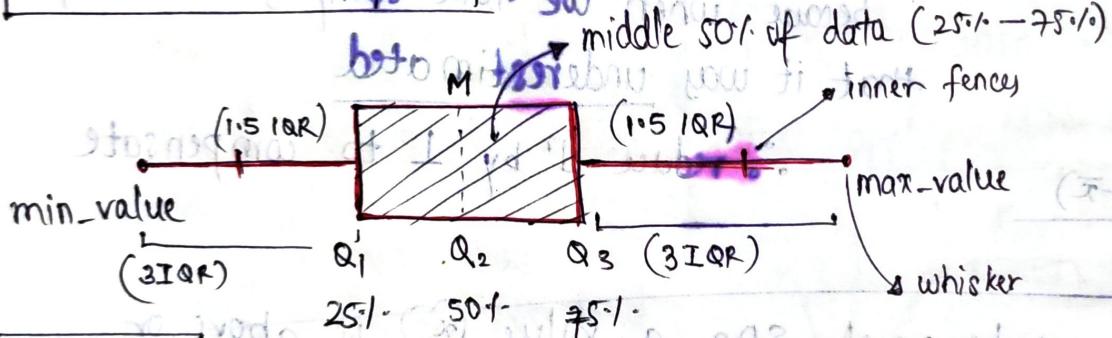
$Q_1 \rightarrow 25^{\text{th}}$ percentile

$Q_2 \rightarrow 50^{\text{th}}$ percentile / median

$Q_3 \rightarrow 75^{\text{th}}$ percentile

$$\text{IQR} = Q_3 - Q_1 \quad (\text{50% of data})$$

Box and Whisker plot



INTERPRETATION

* values after/before/outside inner fences may be treated as outliers

* values b/w (Q_1, Q_3) and inner fences \rightarrow mild outliers

* is distⁿ skewed if median is in right-side of box - left skewed
left-side of box - right skewed

Prob law of addition

$$P(A \cup B) = P(A) + P(B) + P(A \cap B)$$

$$P(A \cup B) + P(A \cap B) = P(A) + P(B) = 1$$

if A & B are mutually exclusive \rightarrow

$$P(A \cap B) = 0$$

Probability

* a priori → can be determined before the experiment

* mutually exclusive events — $P(A \cap B) = 0$

sampling from population

with replacement

$$(N)^n$$

N - population size

n - sample size

Probability Laws

addition

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

prob of either A or B happening

conditional prob

* multiplication

$$P(A \cap B) = P(A) \cdot P(B|A) = P(B) \cdot P(A|B)$$

prob of A and B both happening

$P(A) \cdot P(B)$ if A & B are independent

$$P(X|Y) = \frac{P(X \cap Y)}{P(Y)} = \frac{P(X) \cdot P(Y|X)}{P(Y)} = \frac{P(X) \cdot P(Y|X)}{P(X) \cdot P(Y|X) + P(\bar{X}) \cdot P(Y|\bar{X})}$$

* discrete distns

① Binomial

② Poissons

③ Hypergeometric

* continuous

① uniform

② normal

③ Exponential

distrbns

④ t-distribution

⑤ chi-square

⑥ F-distribution

Distributions

* Independent events

→ occurrence of one event does not affect the other outcome of any other event

$$P(X|Y) = P(X) \text{ and } P(Y|X) = P(Y)$$

WITHSTAND +
AMERICAN

! with replacement

$$N C_n = \frac{N!}{(N-n)! n!}$$

lorentz

gaussian distribution

choose error to meet

$P(A) \cdot P(B)$

$S_1 S_2 = 9$

if A & B are independent

Binomial Distribution

assumptions

- n identical trials
- two possible outcomes
- independent trials → with replacement → no of trials
- $p(\text{success})$ and $p(\text{fail})$ remain same throughout all trials

$$P(X) = \binom{n}{x} \cdot p^x \cdot q^{n-x} = \frac{n!}{(n-x)!x!} \cdot p^x \cdot q^{n-x}$$

$$\mu = np$$

$$\sigma = \sqrt{npq}$$

$p \rightarrow$ probability of success in 1 trial
 $q \rightarrow (1-p)$ for failure

Poisson Distribution

law of rare events

→ focus only on occurrences over a interval

$$P(X) = \frac{\lambda^x e^{-\lambda}}{x!}$$

$x = 0, 1, \dots, n$ (no of occurrences per interval)

λ long duration

$$\lambda = 2.7182$$

Hypergeometric distribution

* Binomial dist'n is for sampling with replacement

Sample without replacement
(trials r not independent)

Conditions

* size of population is known

* The no of success in the population must be known

$$P(X) = \frac{\binom{A}{x} \binom{N-A}{n-x} C_n}{\binom{N}{n}}$$

when to use hyper instead of Binomial?

① Sampling without replacement

② sample size $\geq 5\%$ of population

N - size of population

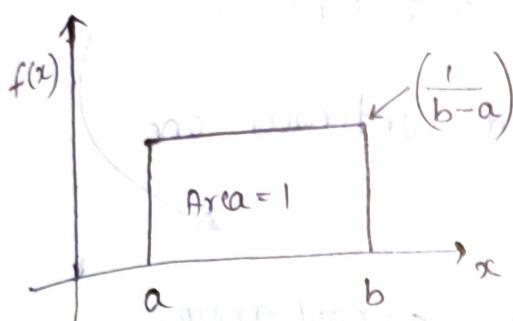
n - sample size

A - no of successes in pop

Sampling is done without replacement

x - no of successes in sample

Uniform Distribution



$$f(x) = \begin{cases} \frac{1}{b-a} & \text{for } a \leq x \leq b \\ 0, \text{ otherwise} \end{cases}$$

$$\mu = \frac{a+b}{2}$$

$$\sigma = \sqrt{\frac{b-a}{12}}$$

* Area enclosed = 1

$$(x-b)(x-a) = (b-a)(\text{height}) = 1$$

$$\text{height} = \frac{1}{(b-a)}$$

Probabilities in a uniform distr

$$P(X) = \frac{x_2 - x_1}{b-a}$$

where

$$a \leq x_1 \leq x_2 \leq b$$

Probability of val lies b/w x_1 and x_2

Normal Distribution

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

standardize
 $\mu=0$

$$z_i = \frac{x_i - \mu}{\sigma} \quad \forall x_i \in X$$

standard = ND with
distribution $\mu=0, \sigma=1$

* continuous and skewed to the right

* $0 \leq x \leq \infty$ and apex is at $x=0$

* as $x \uparrow$, curve steadily decrease

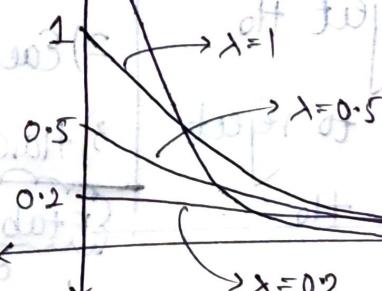
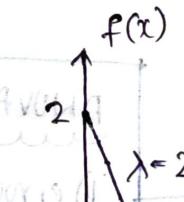
$$f(x) = \lambda e^{-\lambda x}$$

$$\mu = \frac{1}{\lambda}$$

$$P(X \geq x_0) = e^{-\lambda x_0}$$

Ex: prob of $X \geq 2 | \lambda = 1.2$

$$P(X \geq 2) = e^{-(1.2)(2)} = 0.0902$$



$H_{0.05} > H_{0.025} > \dots$

ANOVA

samples → can we say each grp is diff wrt continuous

ONE-WAY ANOVA

→ check if for a categorical variable given a continuous one means are significantly different

Ex:- f_1 : categorical with 3 groups ($n_1 + n_2 + n_3$) $c = \text{no of groups}$

$$SS_{\text{between}} = \sum (\bar{x}_i - \bar{x}_g)^2 = \sum n_i (\bar{x}_i - \bar{x}_g)^2$$

$$SS_{\text{within}} = \sum (x_i - \bar{x}_i)^2$$

$$F = \frac{MS_B}{MS_W} = \frac{\left[\frac{SS_B}{c-1} \right]}{\left[\frac{SS_W}{n-c} \right]}$$

F-score

Ex:- $g_1 \quad g_2 \quad g_3$
 $(1, 3, 5) \quad (5, 7, 9) \quad (4, 10, 5, 6)$

$$\bar{x}_1 = 3, \bar{x}_2 = 7, \bar{x}_3 = 5 \rightarrow \bar{x}_g = \frac{3+7+5}{3} = 5$$

$$\frac{3^2 + 7^2 + 5^2}{3-2} = \frac{15}{3} = 5$$

$$SS_{\text{W}} = [(-2)^2 + 0^2 + 2^2] + [-2^2 + 0^2 + 2^2] + [-1^2 + 0^2 + 1^2] = 18$$

$$SS_B = 3(3-5)^2 + 3(7-5)^2 + 3(5-5)^2 = 24$$

$$\therefore F_{\text{score}} = \frac{(24)}{(3-1)} / \frac{(18)}{(9-3)} = \frac{12}{3} = 4$$

use F-table find p-score

if p-score < α → reject H_0

p-score > α → fail to reject H_0

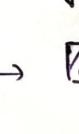
ANOVA assumptions

- 1) groups are conditionally independent
- 2) Each group is normally distributed
- 3) Have same variance

F-table

degree freedom $N_r - (c-1)$

degree of freedom
 $(N-c)$



critical val

if $F < \text{crit}$ → accept

H_0 : \bar{x}_j are same

H_1 : statistically different

PRE+RECALL

TIME SERIES

41

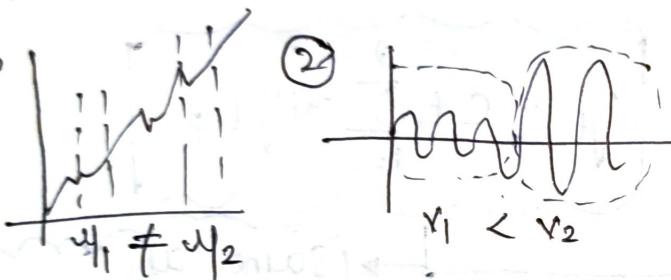
Stationary Time Series

① Mean constant over time

② variance constant

③ Co-variance is only a function \rightarrow take a gap/width of K
of gap/width

↓
Ex where non-stationary



	c_1	c_2
t	$t+K$	
$t-1$	$t+K-1$	
$t-2$	$t+K-2$	
⋮	⋮	⋮

$$\text{gap} = K$$

compute covariance
b/w c_1 and c_2

periodicity

* event occurs periodically
Ex: christmas season

trending

* overall trend
Ex: inflation

\hat{y}_t was predicted using
data $(y_1, y_2, \dots, y_{t-1})$

$c = \bar{y} = \text{mean of all values seen}$
till now

→ used in stationary-time series

ARIMA → moving average

auto regressive

Integrating

(18)

AutoRegression (P) model

→ assumes you can model \hat{y}_t using historic data

AR(P):

$$y_t = \mu + (\alpha_1 y_{t-1} + \alpha_2 y_{t-2} + \dots + \alpha_p y_{t-p}) + \epsilon_t$$

↓ constant ↓ $\alpha = [\alpha_1, \dots, \alpha_p]$ ↓ error term
hyper-param

* to generate o/p at y_t we use data from y_t to y_{t-p} .

$$y_t - \epsilon_t = \mu + \sum_{i=1}^p \alpha_i y_{t-i}$$

$$\hat{y}_t = c + \sum_{i=1}^p \alpha_i y_{t-i}$$

moving average models

MA(q)

$$y_t = \mu + \epsilon_t + \theta_1 \epsilon_{t-1} + \theta_2 \epsilon_{t-2} + \dots + \theta_q \epsilon_{t-q}$$

$$\epsilon_t = y_t - \hat{y}_t$$

↳ errors

Same as
Linear regression
on prev data

use linear
regressions

ARMA (P,q)

$$y_t = \mu + \sum_{i=1}^p \alpha_i y_{t-i} + \sum_{i=1}^q \theta_i \epsilon_{t-i}$$

AR - model on
prev values

MA - model on
prev errors

Non-stationary \rightarrow stationary

① Differencing

$$* y_t' = y_t - y_{t-1} \equiv \frac{\Delta y}{\Delta t=1} = \Delta y$$

$$* y_t'' = y_t' - y_{t-1}'$$

$$y_t = y_t' + y_{t-1}$$

$$\boxed{y_t' = y_t - y_{t-1}} \rightarrow 1^{\text{st}} \text{ order}$$

$$\boxed{y_t'' = y_t - 2y_{t-1} + y_{t-2}} \rightarrow 2^{\text{nd}} \text{ order}$$

$$y_t - 2y_{t-1} + y_{t-2}$$

original :- $y_0 \ y_1 \ y_2 \ y_3 \ y_4$

we can say $y_i' = y_i - y_{i-1}$

$$\begin{array}{cccccc} & y_0 & y_1 & y_2 & y_3 & y_4 \\ \downarrow & y_1' & y_2' & y_3' & y_4' & \\ \text{(Original)} & y_0 & y_1 & y_2 & y_3 & y_4 & y_5 \\ \downarrow & y_1' & y_2' & y_3' & y_4' & y_5' & \end{array} \quad \text{(1^{\text{st}} \text{ order}, } d=1 \text{) before}$$

$$\begin{array}{cccccc} & y_0 & y_1 & y_2 & y_3 & y_4 & y_5 \\ \downarrow & y_1' & y_2' & y_3' & y_4' & y_5' & \\ \downarrow & y_1'' & y_2'' & y_3'' & y_4'' & & \end{array} \quad \text{(2^{\text{nd}} \text{ order}, } d=2 \text{)}$$

Here, new series formed by $d=1, 2, \dots, 3$ may be stationary

Opposite of Differencing \rightarrow Integrating

ARIMA

\rightarrow hyper-params = P, Q, d

AR(p) + I + MA(q)

= ARIMA(p, q, d)

no of terms

no of errors

part p terms and part q errors

ARIMA is linear model applied on

time series post integrate/differentiation.

(4) Auto Correlation

ACF, it can be used for p,q initialization

why? to find best K

K is value by which if you shift
you will get perfect correlation ie $\rho = 1$.

→ How to calculate

① pick K

② construct 2 vectors v_1 and v_2

③ best K will have correlation $\rightarrow 1$ (it will be 1 for $k=0$,
but you can't pick 0)

* Best correlation value

expected (maybe 1 or -1)

to predict y_t
you have only

$y_{t-1}, y_{t-2} \dots$

ARIMA improvements

use features obtained
like $y_{t-1}, y_{t-2}, y_{t-3}$
 $\epsilon_{t-1}, \epsilon_{t-2}, \epsilon_{t-3}$

your own
categorical
features

better regression

algorithms (DF, EBP)

Ex: day,
holiday

if you ~~ARIM~~

ARIMA

Linear Regression = ARIMA

parameters | stop if no log | window

pair | sum pair

Precision Recall matrix

Step 1 :- Compute a confusion matrix

	0	1	2	3	4	...	9
Actual class	0						
predicted class	0	1	2	3	4	...	9
0	1	0	0	0	0	0	0
1	0	1	0	0	0	0	0
2	0	0	1	0	0	0	0
3	0	0	0	1	0	0	0
4	0	0	0	0	1	0	0
...	0	0	0	0	0	1	0
9	0	0	0	0	0	0	1

→ 9 values which belong to class 3 (5) are predicted to E2 class (3)

sum (row 3) = total vals which E2 class (5)

Precision matrix

CF + column normalized

sum of each col = 1

	pred	1	2	3	4	...	9
Actual	7						
pred	7	0.30	0.10	0.10	0.10	0.10	0.10
7	0.30	0.10	0.10	0.10	0.10	0.10	0.10

Out of all pts predicted to E2 class (4)
30% of them actually E2 class (7)

* precision → % of all predn are correct

↔ column-up

1. & out of all pred as P we have
val. P. e2 actual class K

Recall matrix

CF + row-sum = 0

pred	7	1	2	3	4	...	9
actual	3	0.02	0.02	0.02	0.02	0.02	0.02
7	0.02	0.02	0.02	0.02	0.02	0.02	0.02
1	0.02	0.02	0.02	0.02	0.02	0.02	0.02
2	0.02	0.02	0.02	0.02	0.02	0.02	0.02
3	0.02	0.02	0.02	0.02	0.02	0.02	0.02
4	0.02	0.02	0.02	0.02	0.02	0.02	0.02
...	0.02	0.02	0.02	0.02	0.02	0.02	0.02
9	0.02	0.02	0.02	0.02	0.02	0.02	0.02

* recall → how many % of actual class were predicted wrongly

row-down

30% of actual class K are predicted to E2 class P

* of all pts that actually E2 class (3), 21% are classified as E2 (7)

Ideal | of all pts E2 class (2),
diagonals = 1 100% are predicted E2 (2)

16

Compute PDF given CDF

Random variable is

discrete

$$\text{PDF}(x_i) = \text{CDF}(x_i) - \text{CDF}(x_{i-1})$$

① sort acc to x_i

discrete \rightarrow prob mass function

continuous \rightarrow PDF

(i/p) list of $x_i \rightarrow \text{CDF}(x_i)$ (CDF = cumulative sum)

(o/p) find $\text{PDF}(x_i)$

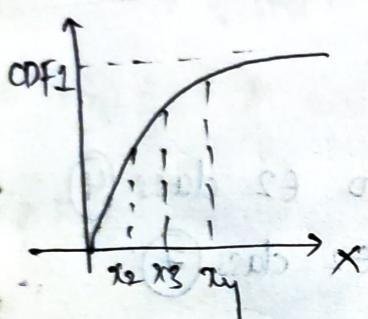
RV is continuous

\rightarrow sort x_i

\rightarrow PDF = derivative of CDF

$$\lim_{(x_2 - x_1) \rightarrow 0} \frac{F(x_2) - F(x_1)}{x_2 - x_1}$$

$$\begin{aligned} F &= \text{CDF} \\ f &= \text{PDF} \end{aligned}$$



so $\rightarrow + x_i \text{ CDF}(x_i)$ in arr:

$$\text{PDF}(x_i) = \frac{\text{CDF}(x_i) - \text{CDF}(x_{i-1})}{x_i - x_{i-1}}$$

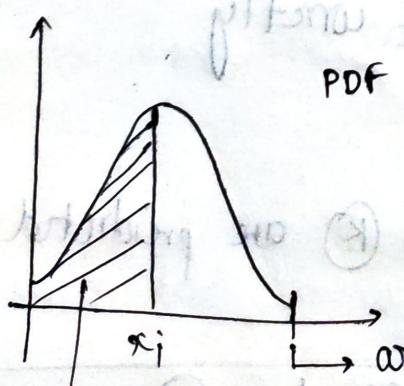
nearest
(x_i, x_{i-1})

How is CDF calculated?

Integrate PDF \therefore PDF values may be > 1

\therefore CDF last val is ① but PMF values are always b/w 0 to 1

area under whole PDF



$$\text{area} = P(x \leq x_i) = \text{CDF}(x_i)$$

0 to 1

Permutation sampling get p-value

Q: 2 samples of spending of men and women

$$S_m = [\dots 100]$$

$$S_w = [\dots n=100]$$

Question: compare S_m and S_w

$$\underline{\text{metric}} : \Delta = |\bar{w}_{\text{men}} - \bar{w}_{\text{women}}|$$

$$\rightarrow \text{observed metric} = |\bar{w}_{\text{men}} - \bar{w}_{\text{women}}| = 0.21 \leftarrow \text{given}$$

Permute + resample

① Null Hypothesis = men and women's spending habits are same

② mix $S_m + S_w = [\dots]_{n=100}$ (as H_0 : same)

why?

if \bar{w}_m and \bar{w}_w are same (H_0)
then samples r also similar

randomly sample S_m and S_w with 100 samples each] repeat K
and calculate metric times

③ let $K=1000$ (u get 1000 test-metrics $-(ts)$)

Now,

$$p\text{-value} = P(X \geq \text{obs}) = P(ts > 0.21) = 0.02$$

prob of atleast seeing value equal to or

greater than observed

Result

p-val: prob of getting val atleast
as extreme as observed

→ if $p\text{-val} = P(ts > 0.21) < 5\%$ → reject H_0 bcoz in our testing

few than 5% of test samples

had metric atleast (\geq)
observed

i.e. prob of observing metric is
very low

(*) **KSTest** check if ~~this~~ both r from same distribution
 → parametric

test-statistic

$$D_{n,m} = \sup_x |F_{1,n}(x) - F_{2,m}(x)|$$

$F_{1,n}$ → CDF value of x_1 with n samples

\sup = supremum fn

+ Null hypothesis is rejected at $\alpha = 0.05$

if

$$D_{n,m} > c(\alpha) \sqrt{\frac{n+m}{nm}}$$

$$c(\alpha) = \sqrt{-\frac{1}{2} \ln(\frac{\alpha}{2})}$$

Note:

if m, n are large,

CDF's are almost similar (very less gap)

Exact p-value

$$D \sqrt{\frac{nm}{n+m}} > \sqrt{-\frac{1}{2} \ln(\frac{\alpha}{2})}$$

$$\Rightarrow \exp\left[-2D^2 \left(\frac{nm}{n+m}\right)\right] < \alpha$$

↓ squaring and mult by -1

$$\therefore p\text{-val} = \exp\left[-2D^2 \left(\frac{nm}{n+m}\right)\right]$$

$$-D^2 \left(\frac{nm}{n+m}\right) < \frac{1}{2} \ln\left(\frac{\alpha}{2}\right)$$

further if $p\text{-val} < \alpha$, reject

VIF variable inflation factor (49)

step 1 : Try to predict f_i using other independent features

step 2 : If they can predict f_i well

→ multi-collinearity

$$\boxed{VIF(f_i) = \frac{1}{1-R^2}}$$

$\rightarrow \left(1 - \frac{SS_{res}}{SS_{tot}}\right)$ of model which fits f_i

as $R^2 \rightarrow 1$, better fit

as $R^2 \rightarrow 1$, $VIF \uparrow$

$\rightarrow VIF = 0$ (no collinearity, $R^2=0$) doesn't fit

$VIF > 10, 15 \rightarrow$ multi-collinearity

$$est - (P \cdot f_i + b) \div b$$

$$\begin{array}{|c|} \hline 9 \cdot \bar{w} = 6 \\ \hline 11 \cdot \bar{w} \\ \hline \end{array}$$

$$\begin{array}{|c|} \hline 9 \cdot \bar{w} = 6 \\ \hline 11 \cdot \bar{w} \\ \hline \end{array}$$

$b = \frac{9 \cdot \bar{w} - 6}{11 \cdot \bar{w}}$ = with no $(9, 6)$ to calculate b

→ is it really following a linear trend?

with b

$est = P \cdot w + b$

against dependent value

$10 + 8 \cdot \bar{w} / 11 = 6$

$26 + 8 \cdot \bar{w} / 11 = 8$

$$(6 + p_w)b = (7 + q_w)b + (8 + r_w)b$$

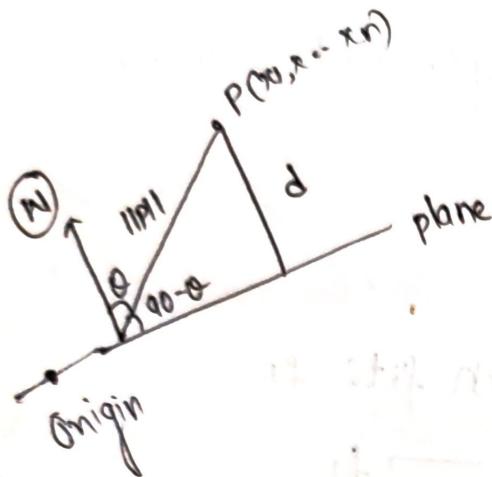
the same

$$\frac{(6 + p_w)b}{(7 + q_w)b + (8 + r_w)b} = \frac{6}{7 + 8}$$

$\rightarrow \text{if } p_w = q_w = r_w$

Q5

Distance of point from plane



$$\sin(90 - \theta) = \frac{d}{\|P\|} \Rightarrow d = \|P\| \cos \theta$$

$$W \cdot P = \|W\| \|P\| \cos \theta$$

$$\frac{W \cdot P}{\|W\|} = \|P\| \cos \theta = d$$

$$d = \frac{W \cdot P}{\|W\|}$$

Ex: plane in 2d = line

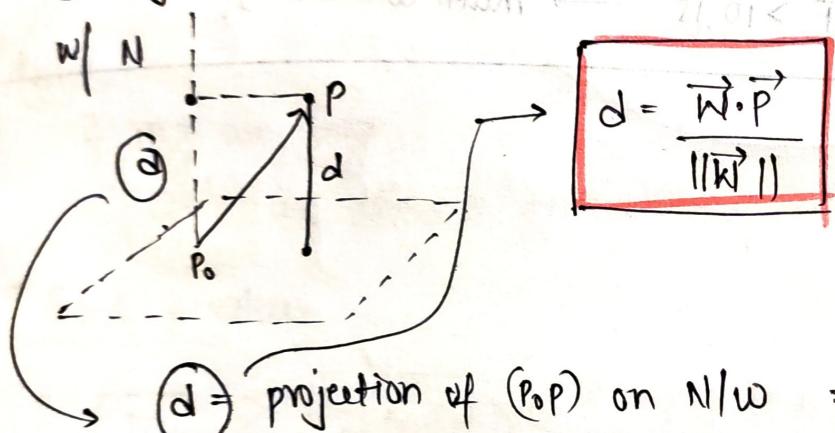
$$P(7, 9)$$

$$Line = 7x + 4y - 120 = 0$$

$$W = (7, 4) - 120$$

$$d = \frac{(7 \cdot 7) + (4 \cdot 9) - 120}{49 + 16}$$

(ii) Projection view



$$d = \text{projection of } (P_0P) \text{ on } N/W = \text{Proj}_{\vec{W}} \vec{P} = d$$

Shortest distance b/w 2 parallel planes π_1, π_2

if $\pi_1 \parallel \pi_2$,

then $\nabla W_1 = W_2$,

only intercept/b changes

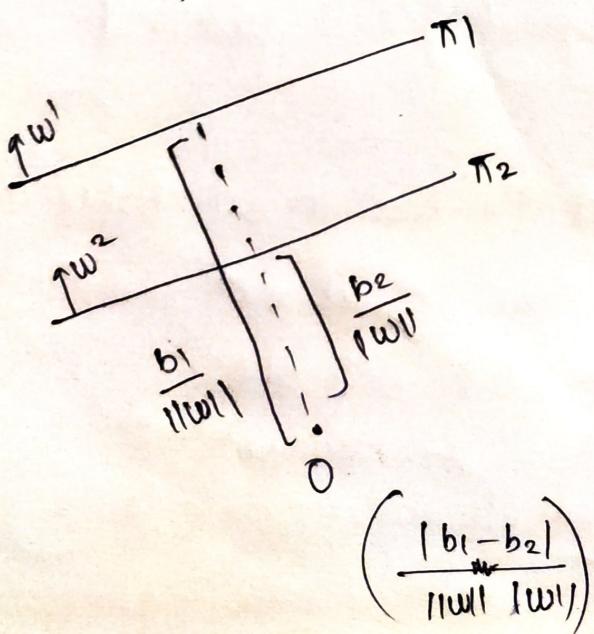
$$\pi_1: W_1 x + b_1$$

$$\pi_2: W_2 x + b_2$$

$$\therefore d(\pi_1, \pi_2) = d(\text{origin}, \pi_1) - d(\text{origin}, \pi_2)$$

$$= \frac{b_1 - b_2}{\|W_1\| \|W_2\|}$$

$$= \frac{\|b_1 - b_2\|}{\|W_1\| \|W_2\|}$$

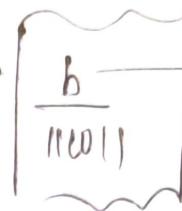


dist of origin from plane

(51)
term

Eqn of π passing through origin = $w^T x = 0$

not passing $O(0,0)$ = $w^T x + b = 0$



$\vec{x} = \vec{w}$

$$d = \frac{w \cdot p + b}{||w||}$$

dist of point from plane
which doesn't pass origin

$p = O(00\ldots 0)$

$$d = \frac{w \cdot [0\ldots 0] + b}{||w||} = \frac{b}{||w||}$$

\therefore shortest dist b/w

parallel planes π_1 and π_2 =

$$\frac{|b_1 - b_2|}{||w||}$$

R² and adjusted R²

$$R^2 = 1 \leftarrow \text{Best} \quad (\text{SSres} = 0)$$

$$R^2 = 1 - \frac{\text{SSres}}{\text{SStot}} = 1 - \frac{\sum(x_i - \bar{x})^2}{\sum(x_i - \bar{x})^2}$$

↓
when more features are added → as nfeat↑
to regression (multi-linear regression) SSres ↓ → R²↑
the added feature is not correlated with y

penalize more features

$$R_{\text{adj}}^2 = 1 - \frac{(1-R^2)(N-1)}{N-p-1}$$

← penalize features when newly added features are not correlated with y.

(N) = no of sample

(P) = no of predictors(features)

↳ no of independent features