

Featurization and Feature engineering

① Introduction :-

Text data :- Bow, tfidf, avg2v, tfidf w2v

also sequence data various featurization of text data.

Categorical data :-

- one hot encoding
- mean-response rate,
- domain Specific.

Time-Series :-

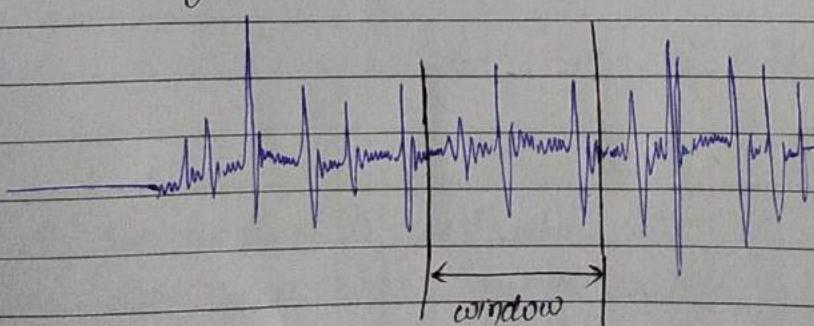
- heart rate
- ecommerce
- speech / audio
- stock market price

Image data :- face detection, face-recognition
X-rays, MRI scans, video

Digit + time series.

Graph data :- recommend a friend on FB

② Moving window for Time Series Data



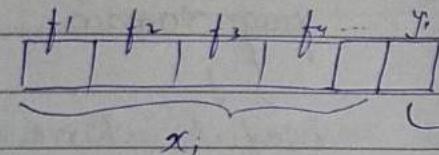
∴ window width depends on data.

- mean, std-dev
- median, Quantile
- max, min
- local maxima
- zero-crossing / mean crossing

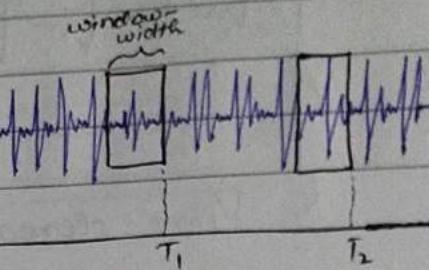
how many times wave curve crosses mean line

* Steps of featureization :-

- Decide window width
- Features that are useful

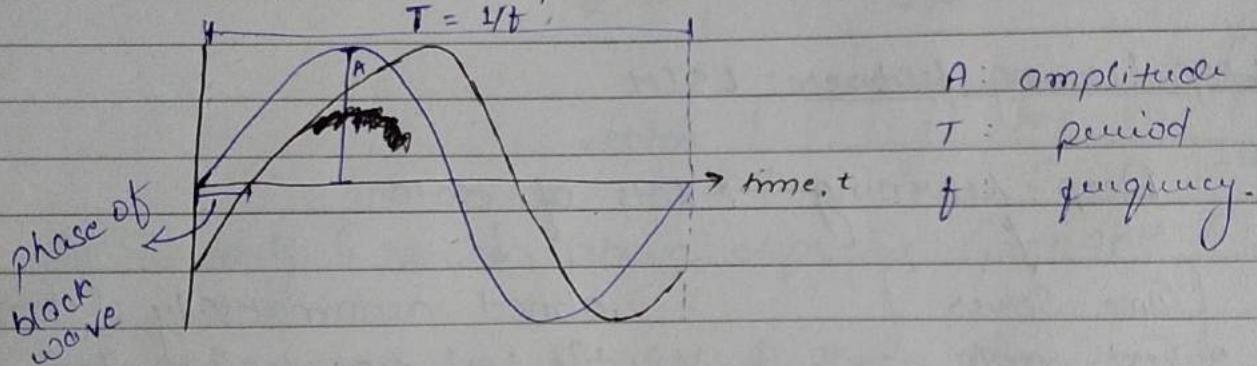


↳ boolean: heart attack in next 10 min or not.



* Fourier decomposition for Time Series Data :-

↳ Method to represent Time-Series.



e.g. if 1 osc. per second $\rightarrow 1 \text{ Hz}$
2 osc. per Second $\rightarrow 2 \text{ Hz}$.

Youtube: Programming Cradle

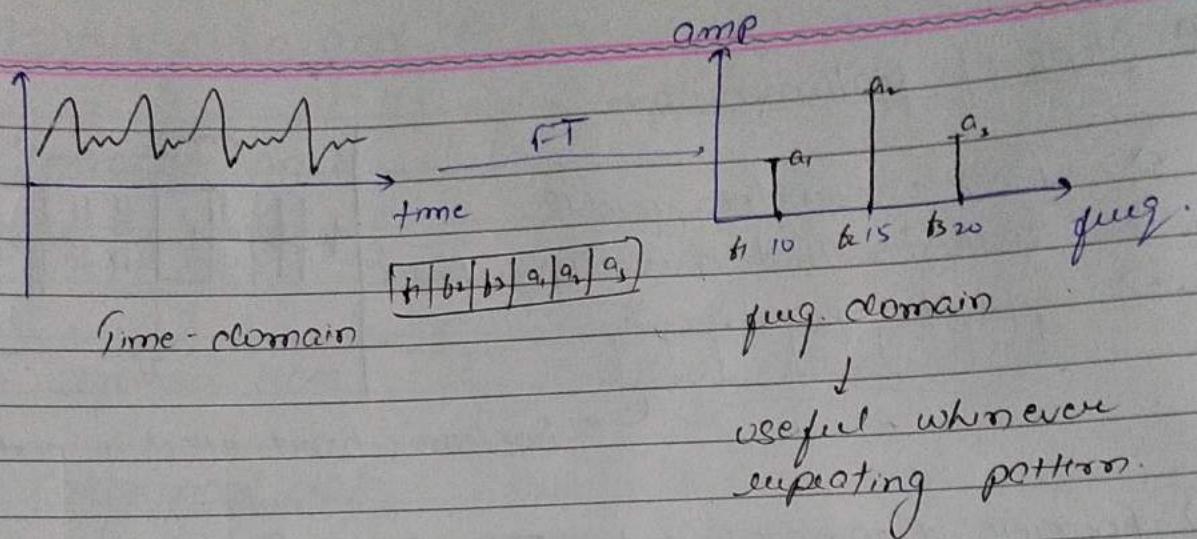
heart rate :- 60-100 beats per minute

|| ~~~~~
osc.

1-1.6 per Second. $\Rightarrow 1 \text{ Hz to } 1.6 \text{ Hz}$

* Fourier Decomposition / Transform:- A composite / complex wave can be represented as sum of various sine waves and convert it to freq vs Amp. plot.

Refer wikipedia FFT-Time-frequency-view.png



NOTE :- Whenever repeating pattern is there frequency is important.

④ Deep Learning features : LSTM

Deep - Learning → lots of data

↓

Time - Series almost automatically learn
Text - data the best featurization for
Image - data data

Deep Learnt features → best features to day.

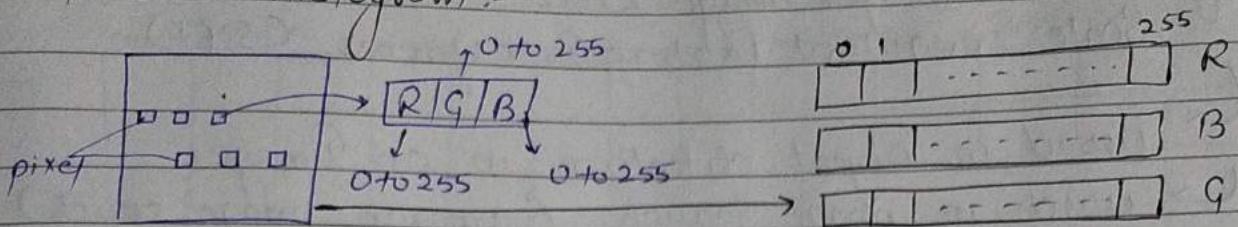
⑤ Image Histograms :- Youtube: Programming Cradle

Images :- faces, object, scans, x-rays, autonomous car

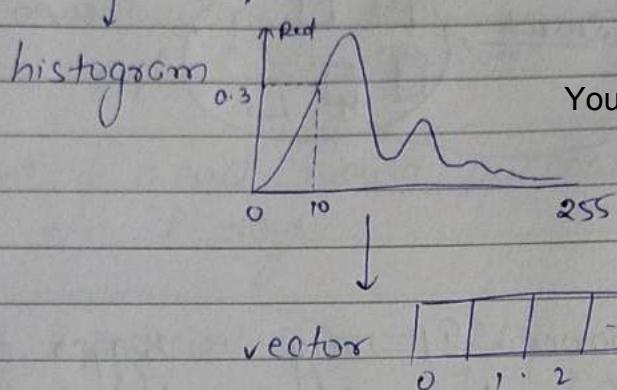
Two popular types of histograms:-

↳ Color histogram } basic and rudimentary
↳ Edge histogram. }

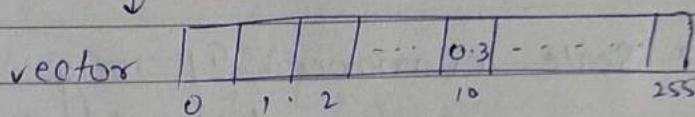
* Color histogram :-



① Collect Red values for each pixel

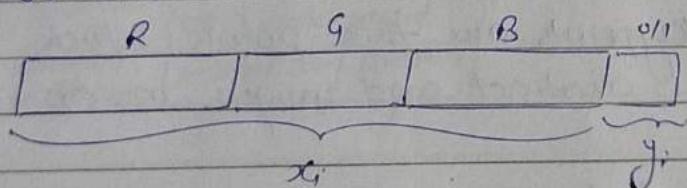


Youtube: Programming Cradle



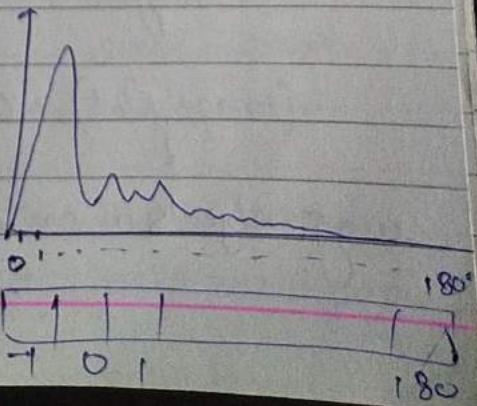
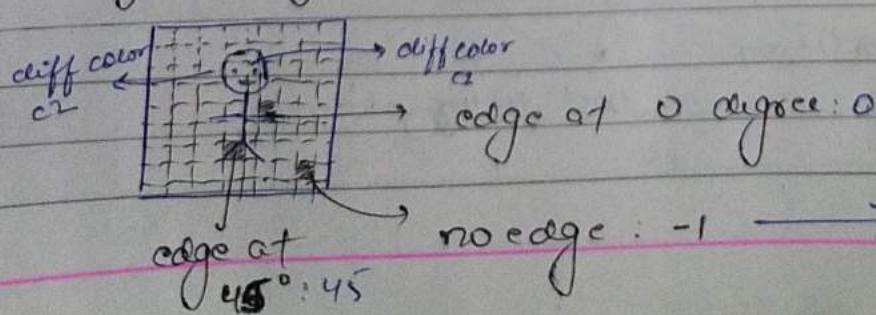
② Similarly I can do above steps for Green and Blue.

Example :- We want to detect if sky is there in the image or not.



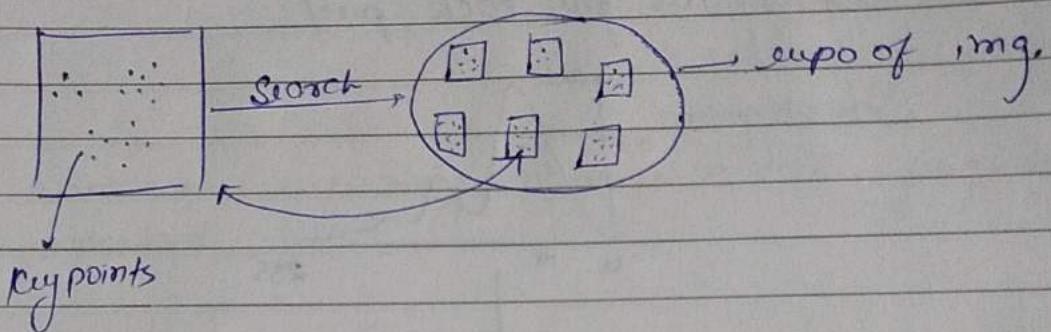
if sky present \rightarrow B will have more values $\rightarrow 1 = y_i$
 if sky not present \rightarrow B will have less values $\rightarrow 0 = y_i$

* Edge-histogram :-



* Keypoints :- SIFT for image Data
Scale Invariant Feature Transforms (SIFT)

- ↳ used in object detection in an Image.
- ↳ used in image search. (amazon image search)



- Scale invariance:- If image is bigger or smaller features do not change much.
- Rotation invariance:- If image is rotated slightly features do not change much (may change/ not work with too much rotation)
- Keypoint:- Keypoints are those points which are very distinct and unique in an image

* Deep learning features: CNN :-

Time-Series \rightarrow (LSTM) \rightarrow feature

Images \rightarrow (CNN) \rightarrow feature

(image Net) \rightarrow competition
! best featurization

many objects in an img

④ Relational Data :-

custid	cust pin code	custid	pid	time	custid	pid	time	pid	product type
eg.									

cust table

cust viewing /

visitation data

purchase

product data

Relational DBs :- MySQL, Oracle, SQL server

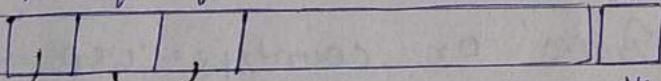
⑤ Task : predict if a cust would purchase a product in the next 7 days.

Given

cust id , product id → 1/0

1 ↗ do not buy
buy

t1 t2 t3



j1

no. of times the cust id visited the product page in last 24 hrs

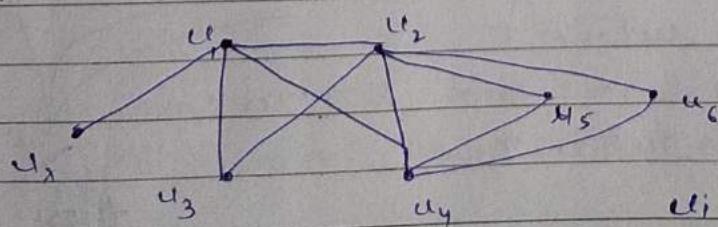
no. of times cust id visited any product type same as pid

Youtube: Programming Cradle

pin code (certain pin codes have richer people)

⑥ Graph Data :-

eg:-



Social - graph

ui → vertex (user)

edge → connection (friend ship)

* Task:- Recommend new friends for a user u_1 :

u_{10} Suggestions
 u_1 $\{u_3, u_7\}$

f_1 :- no. of mutual friends

f_2 :- no. of paths b/w 2 users

* Indicator variables :-

eg:-① 'height' as a feature

↳ pupil as is $h \in \mathbb{R} \rightarrow h \rightarrow$ real valued feature

↳ convert it to indicator variable

↳ if $h > 150 \rightarrow 1 \quad \} \text{ binary indicator}$
 $h \leq 150 \rightarrow 0 \quad \} \text{ gesture.}$

eg:-② Categorical feature:- Country.

Indicator variables
 $\left\{ \begin{array}{l} \text{if country == 'India' or country == 'USA'} \\ \quad \text{return 1} \\ \text{else} \\ \quad \text{return 0} \end{array} \right.$

* Feature binning :- Youtube: Programming Cradle

↳ Extension to indicator var.

eg- height $\in \mathbb{R}$

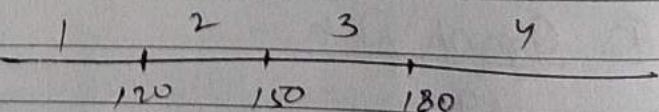
if $h < 120$

return 1

if $h < 150$ and $h \geq 120$ cm.

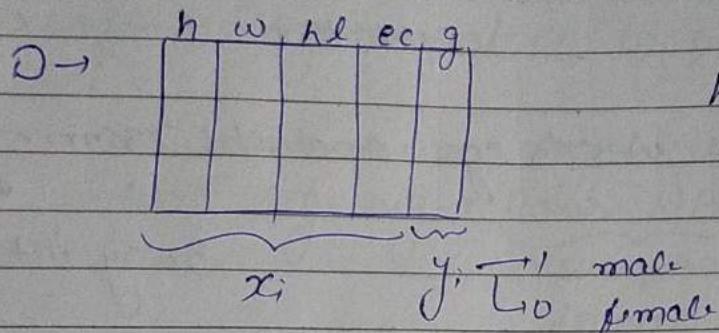
return 2

so on



these thresholds
are problem
specific.

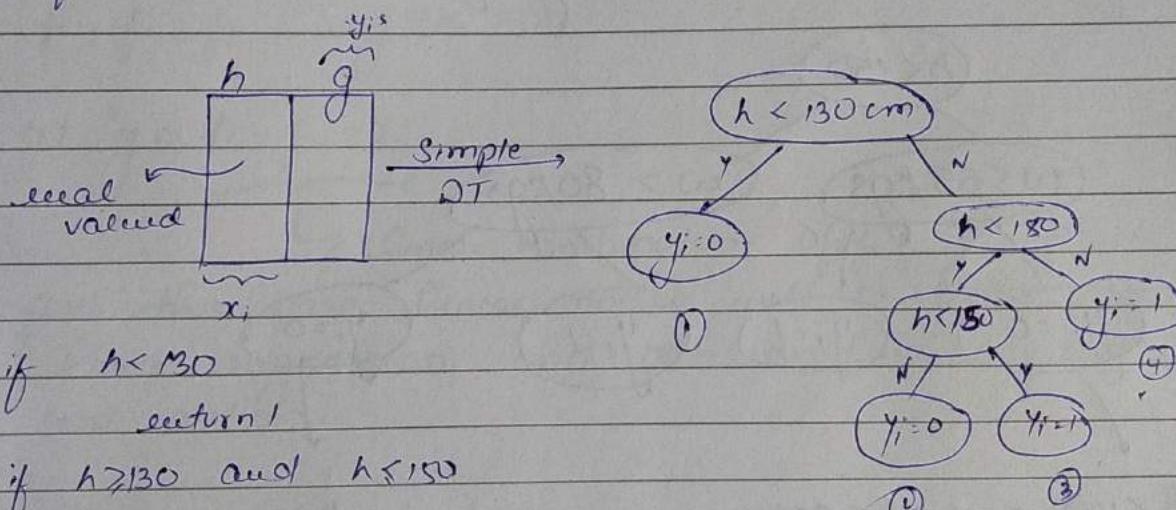
Task:- Predict gender given :- h, w, hair length, eye color



height → binning.

Challenge:- binning of 'h' using D (without domain knowledge)

one of the soln:- Youtube: Programming Cradle



if $h < 130$

return 1

if $h \geq 130$ and $h < 150$

return 2

if $h \geq 150$ and $h < 180$

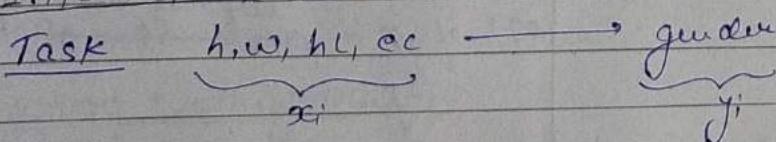
return 3

if $h \geq 180 \text{ cm}$

return 4

binned equal-valued features
using y_i 's and feature itself
using DT

Interaction variables:-



eg:- ① ($h < 150 \text{ cm}$) and ($w < 60 \text{ kgs}$) → logical 2 way

interaction feature

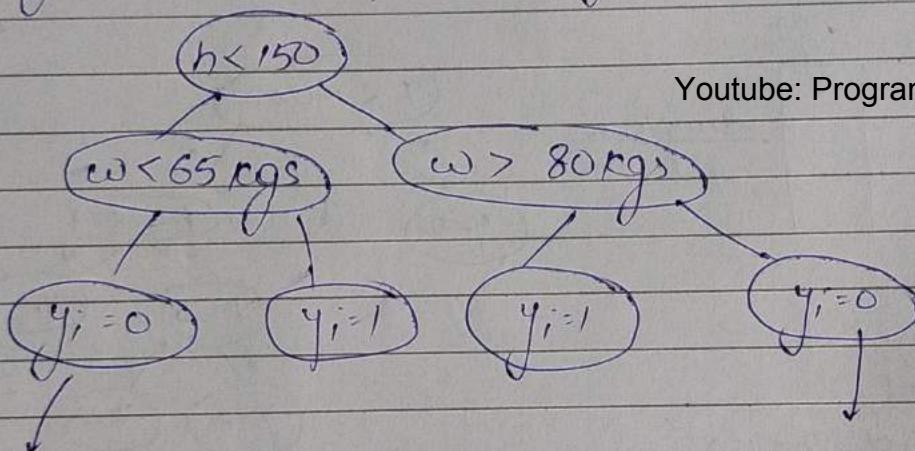
eg) ② $h * w \rightarrow f_2$

operation depends on
use case, can be any
operator +, -, *, / etc.

eg) ③ $h < 150 \text{ cm. and } w < 65 \text{ kgs and } hc > 5 \text{ cm} \rightarrow$

3 way interaction

Q) Given a task, how to find good interaction features?

Soln DTeg $x_i = h, w, hc, y_i = \text{gender}$  $h < 150 \text{ and } w < 65 \text{ kgs}$ $h > 150 \text{ cm and } w > 80 \text{ kgs.}$

★ Mathematical Transforms :-

 $x \rightarrow \text{Single feature}$

Q. what is the best transform?

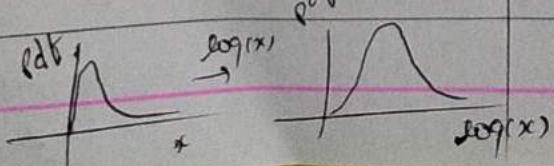
A) Problem specific.

eg) if $x \rightarrow$ ~~power law dist.~~ log normal dist. $\log(x) \rightarrow$ best transform
bcz it will make the dist as

normal dist.

Math transforms

$\log(x)$, e^x
\sqrt{x} , $\sqrt[3]{x}$
x^2, x^3 ... polynomial
$\sin(x), \cos(x), \tan(x)$



✳ Model Specific featureizations :-

e.g.: let $f_1 \rightarrow$ log normal. $\Rightarrow \text{log}(f_1) \rightarrow$ Gauss. dist.

We want to use logistic regression and we know that log-eug \rightarrow Gauss NB. \rightarrow becz of assumption we can assume: not f_1 directly. feature is hence we will transform using $\text{log}(f_1)$ can be used for logistic reg.

e.g. 2 $f_1, f_2, f_3 \rightarrow y \in \mathbb{R}$

at $y \approx f_1 - f_2 + 2f_3$

Youtube: Programming Cradle

↳ linear combination of f_i 's

\rightarrow for this case linear models will be better option like linear regression. and decision tree may not be best option.

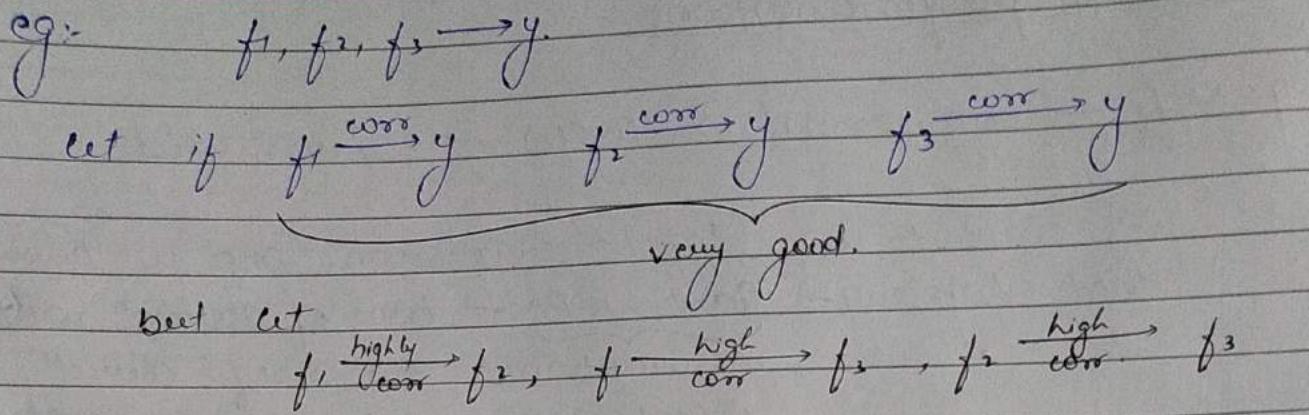
e.g. 3 $y \rightarrow$ interactions of f_1 and $f_2 \leftarrow$ domain knowl.
 $h.w \rightarrow$ gender \therefore DT / RF / GBDT best opt.

- text data : (BOW) \rightarrow linear models are better most of the times. as they have high dim. and linearly separation is easy in higher dim.

✳ Feature Orthogonality:-

*> the more ~~the~~ perp

*> the more diff / orthogonal the features are, the better would your model be.



→ because of high correlation among features effect will be less.

eg:- $f_1, f_2, f_3 \rightarrow y$. Youtube: Programming Cradle

let if $f_1 \xrightarrow{\text{corr}} y$, $f_2 \xrightarrow{\text{corr}} y$ $f_3 \xrightarrow{\text{corr}} y$

and $f_1 \xrightarrow[\text{corr}]{\text{not}} f_2$ $f_1 \xrightarrow[\text{corr}]{\text{not}} f_3$

→ in this case effect overall will be good.

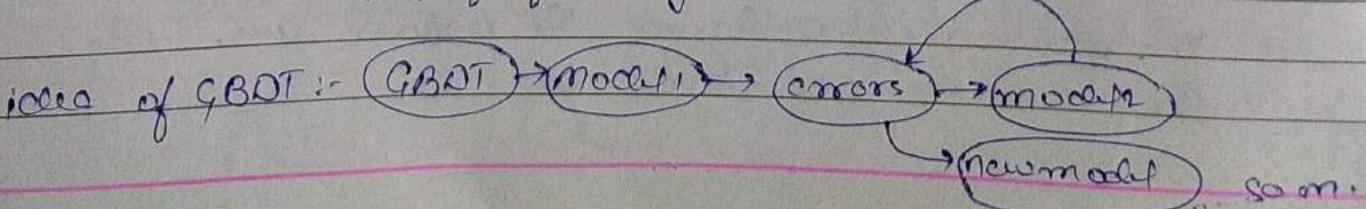
NOTE :- while creating new feature we should try (ideally) new feature is highly corr with dependent feature (y) and not corr with other independent features (f_i)

Q) How do I design a new feature f_4 s.t. $f_4 \xrightarrow{\text{corr}} y$; and less corr with f_1, f_2, f_3 ?

1 of the soln:-

or, idea:- Δx_i find error; = $y_i - \hat{y}_i$ $f_4 \rightarrow \text{error}$;

$f_1, f_2, f_3, f_4 \rightarrow y$.



* Domain Specific generalizations :-

e.g. heart attack \rightarrow ECG data

using this
we can create
new features.

important to research and study
existing generalizations doctors/
specialists.

* Feature slicing :- Youtube: Programming Cradle

ec \rightarrow eye color

e.g. h, w, hc, ec, country \rightarrow gender.

x_i

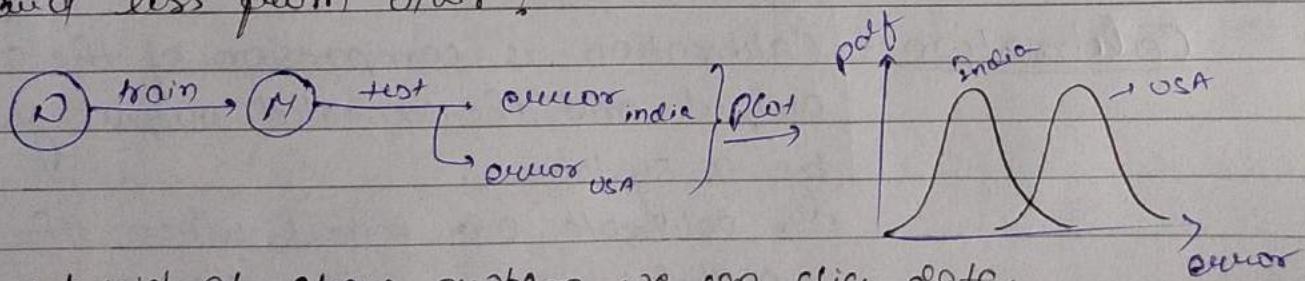
India USA

country	Ind	USA
Ind	1	0
USA	0	1
	80%	20%

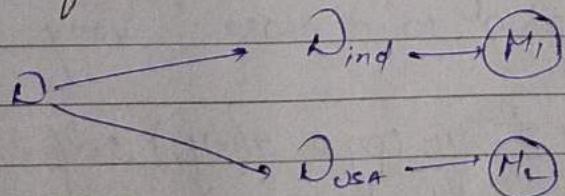
- In this case where we have more people / datapoint from India model will tend to perform better and for USA, model tend to perform poor.

Q) How to detect if more people are from specific country and less from other?

A)



∴ To get rid of above problem we can slice data.



∴ Have ~~diff~~ separate models for each D_{India} and D_{USA}

- Slicing the data based on feature's category. when
 - ① category1 and category2 are diff
 - ② Sufficient no. of pts for each slice of data.

* Calibration of Models: Need for calibration.

eg:- At 2 class classif? $y_i \in \{0, 1\}$, $D_{\text{train}} = \{x_i, y_i\}_{\text{train}}$

now given $x_q \Rightarrow f(x_q) = y_q$

\downarrow
1 or 0

[Model f(x)]

now I want to understand $p(y_q=1/x_q, f(x))$
 but the o/p of f may or not be a prob. score.

At talk about NB, we calculate $p(y_q=1/x_q)$ and $p(y_q=0/x_q)$ which is proportional to, prior * likelihood but not equal.

So, $f(x_q) = y_q$ and y_q may or maynot be the actual $p(y_q=1/x_q)$

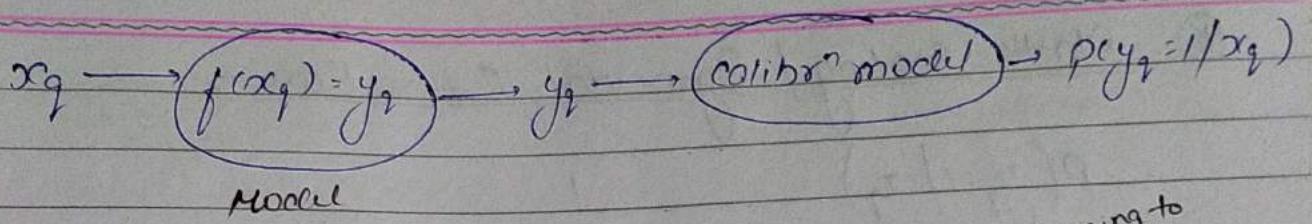
Calibration :- Calibration is comparison of the actual output and the expected output given by a system.

To calibrate our model when the probability estimate of a data point belonging to a class is very important.

eg:- At 2-class

$$\text{log-loss} = -\frac{1}{n} \sum_{i=1}^n y_i \log p_i + (1-y_i) \log(1-p_i)$$

here we use actual p_i , ~~but if p_i 's are wrong log-loss will be wrong. hence we need calibration~~



④ Calibration Plots :-

probability of given data point belonging to a specific class.

$$D_{train} \rightarrow f^n$$

probability

of

f(x)

↑

true class label

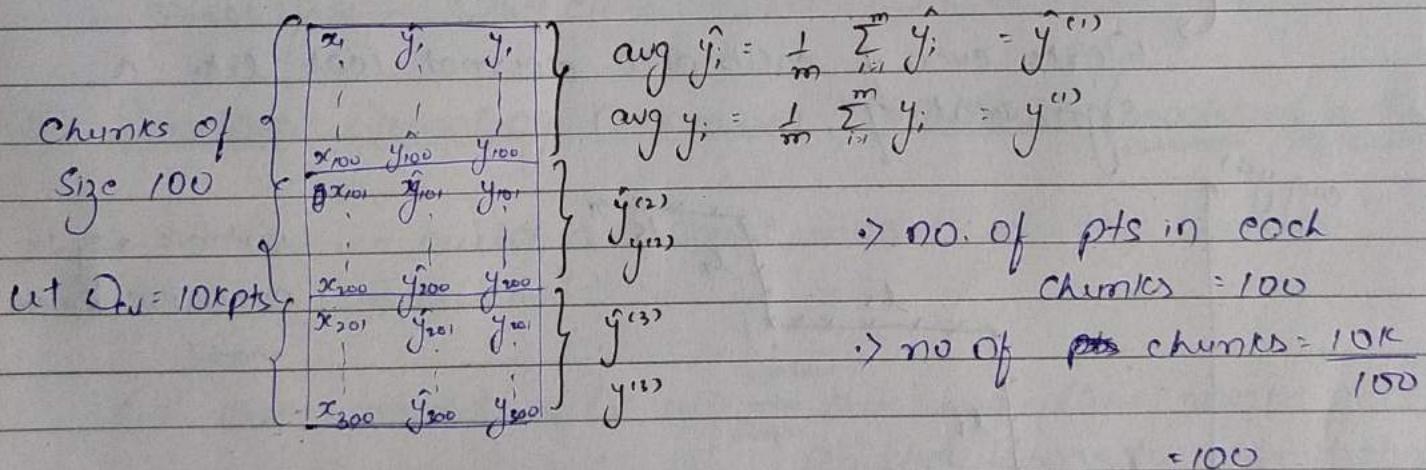
$$\textcircled{1} \quad D_{cal} \rightarrow f^n \rightarrow \hat{y}_i \quad \Rightarrow$$

x_1	\hat{y}_1	y_1
x_2	\hat{y}_2	y_2
:	:	:

$$\textcircled{2} \quad \text{Sort this table in } \uparrow \text{ order of } \hat{y}_i = f(x_i)$$

$$\hat{y}_1 \leq \hat{y}_2 \leq \hat{y}_3 \leq \dots$$

$$\textcircled{3} \quad \text{Break table obtained in second step in chunks.}$$



$$DataCalib = \{\hat{y}^{(c)}, y^{(c)}\}$$

Youtube: Programming Cradle

avg $f(x_i)$ → $p(y_i = 1/x_i)$

real world.

Calib plot :-

$p(y_i = 1/x_i)$, avg $\hat{y}^{(c)}$, $y^{(c)}$

ideal

$(y_{i00}, y_{i00}) \rightarrow$ set [chunk 2 → $\hat{y}_i \approx y_i$]

0.3

0.3

0.3

$\hat{y}^{(c)} = \text{avg } \hat{y}_i$

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

0.3

* Plot's Scoring / Sigmoidal Calibration :-

$$P(y_i=1/x_i) = \frac{1}{1 + \exp(Af(x_i) + B)}$$

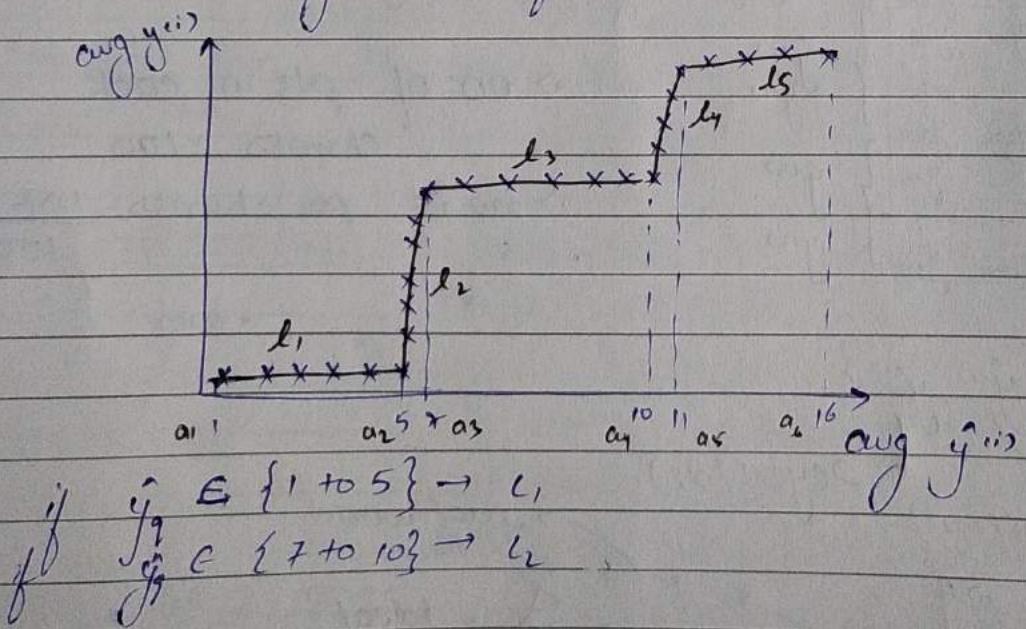
Q. How to find A and B?

A) Replace $f(x_i) \rightarrow \hat{y}^{(i)}$ and $P(y_i=1/x_i) \rightarrow y^{(i)}$ and then find best A and B by solving optimization problem. (there are various formulat of this opt^m problem)

NOTE :- PIAH's calibration works iff calib. plot looks like sigmoidal fn

* Isotonic Regression :-

↳ Works even if calib. plot does not look like a sigmoidal fn.



isotonic reg. :-
piece wise linear
model

$a_1, a_2, a_3, a_4, \dots$

$m \rightarrow \text{stop}$
 $b \rightarrow \text{intercept}$

$(m_1, b_1) (m_2, b_2) (m_3, b_3)$

- Optimization problem :-
- (find threshold a_i) \rightarrow find optimal m_i and b_i for each line
- Also there shouldn't be too many lines
(just enough no. of lines)

NOTE :- In isotonic reg; need many more points / data in Data than with Platt scaling (A,B) because in isotonic reg we have to find many more parameters $(a_i), (l_i)$

- Large Data \rightarrow Large CV \rightarrow Large Calib Data \rightarrow Isotonic reg.
- Small Data \rightarrow Small CV \rightarrow Small Calib \rightarrow Platt scaling.

* Random Sampling Consensus (RANSAC) [modeling in the presence of outliers]

Can we build a robust model in the presence of outliers?

Youtube: Programming Cradle

If our model is Lr. reg. \rightarrow Our hyperplane might get shifted towards dominating outliers present in the data.

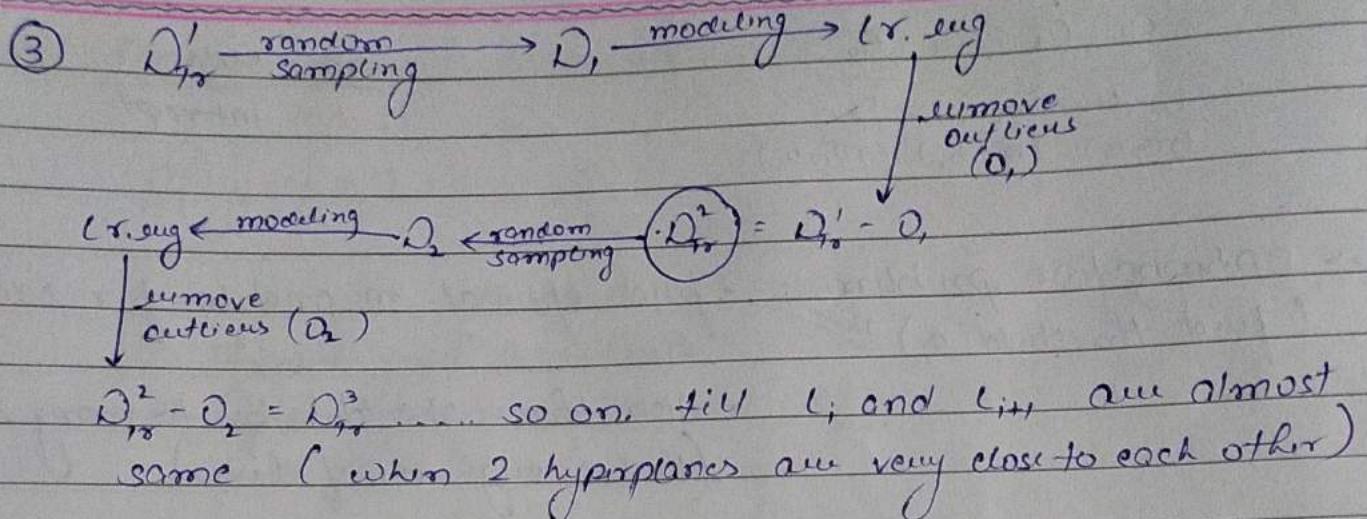
* RANSAC steps :-

① $D_o = \underbrace{\text{random sample from } D_{tr}}_{\text{without RANSAC}}$ \rightarrow This will also \downarrow chance of outliers by $\approx 50\%$.

② Compute outliers based on C_o if $y_i - \hat{y}_i$ is high \rightarrow outlier

$$O_o = \{\text{outlier pts}\}$$

$$\boxed{D'_{tr} = D_{tr} - O_o}$$

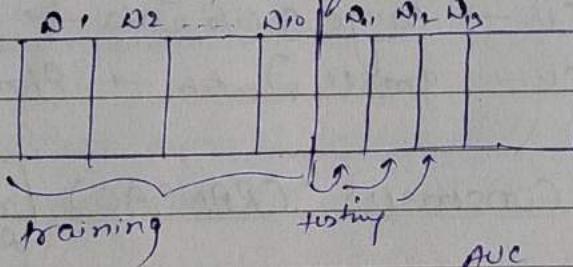


D_{t+1}^{tr} → outlier free data

l_{t+1} → outlier robust model

* Retraining models periodically. Youtube: Programming Cradle

e.g.: predict stock prices for the next day.



productionized model (M_{10})
on Day 11

1st Day 11 → $M_{10} \rightarrow 0.8$

Day 12 → $M_{10} \rightarrow 0.81$ → huge drop in AUC.

Day 13 → $M_{10} \rightarrow 0.6$

* How to determine if we have to retrain?

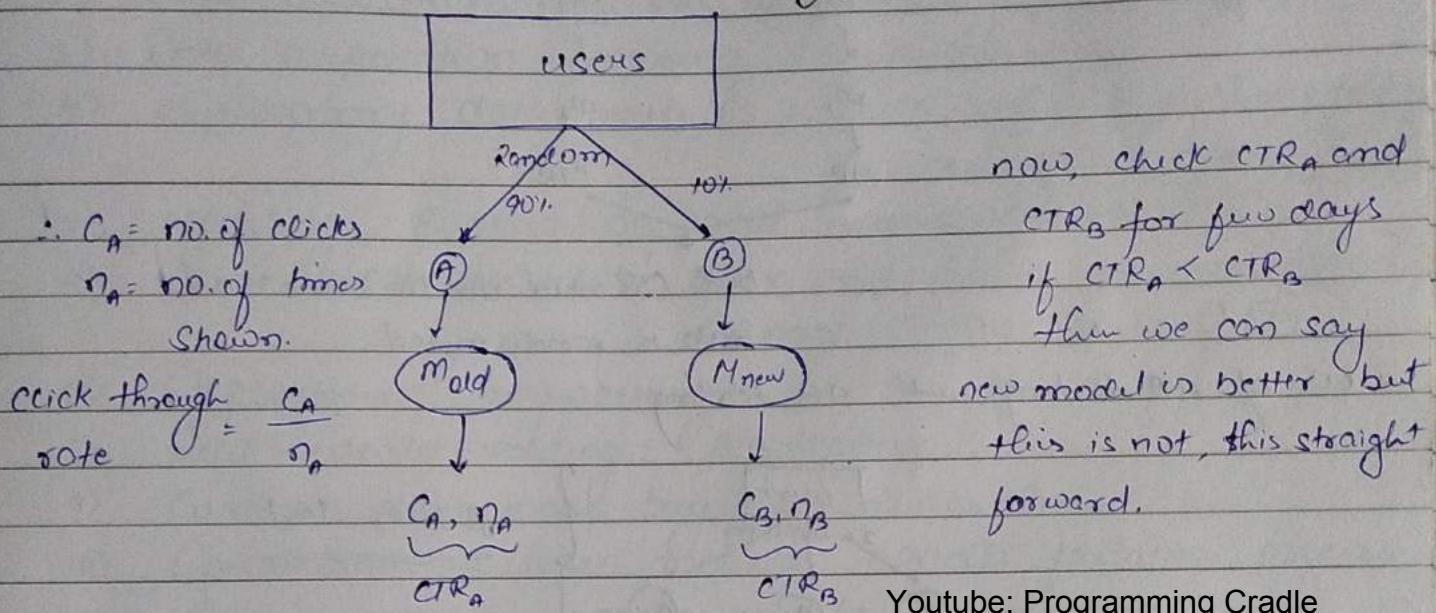
→ model performance dropped

→ dataset / data pts. changed. (non-stationary data)

if retraining cost is not too much retrain periodically and period can be small, as small as 1 min/hr.

★ A/B Testing :- Bucket testing or Split sum or controlled exp.

★ Eg:- Develop model for Google predict to which ad to show for a search query.



Youtube: Programming Cradle

$$\rightarrow \textcircled{A} \quad C_A \quad n_A \\ 100 \quad 100K$$

$$CTR = \frac{100}{100K} = 1\%$$

$$\textcircled{B} \quad C_A \quad n_A$$

$$2 \quad 100 \\ CTR = \frac{2}{100} = 2\%$$

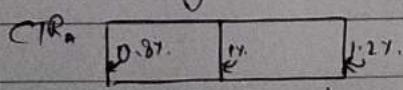
\therefore even though CTR of B is greater than A we can not say model B or new model is good becoz model B is overall showing less ads.

Hence we use confidence interval.

95% CI

$$CTR_A = [0.8\%, 1.2\%]$$

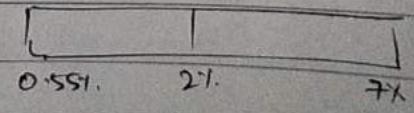
tighter



95% CI

$$CTR_B = [0.55\%, 7\%]$$

wide

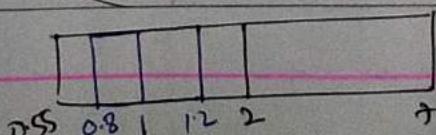


if

$CI(CTR_A)$

$CI(CTR_B)$

\therefore not overlapping hence can say CTR_B is better than A

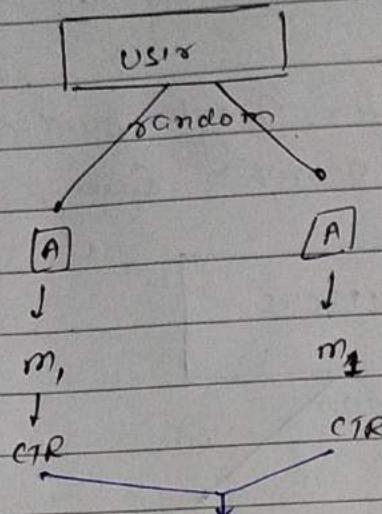


→ completely overlapping

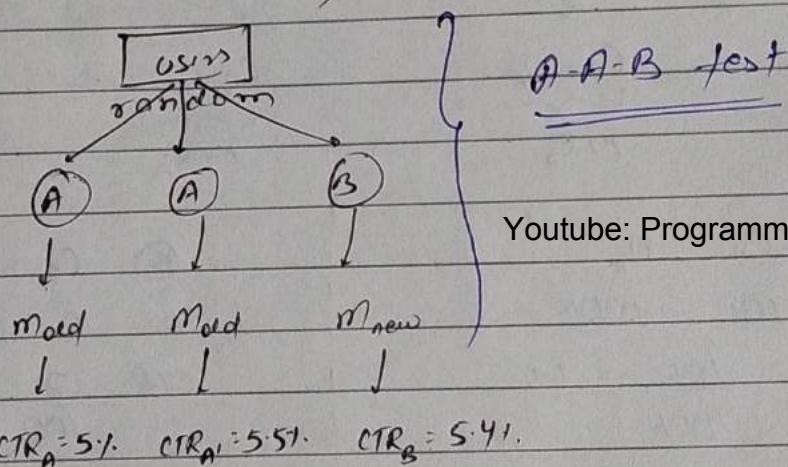
hence can't conclude

④ AA testing

tells the variance of same model on diff data



→ use of AA testing with AB testing



if we compare CTR_A and CTR_B we may conclude CTR_B is better but when we see all CTR_A , CTR_{A1} and CTR_B we can conclude this improvement might be just bcoz of randomization

* Data Science / ML project life-cycle

- 1> Understand business requirements :- define the problem, customer.
- 2> Data acquisition: ETL, SQL, DB, DW, log-files, Hadoop/spark]
- 3> Data preparation: cleaning, pre-processing.
- 4> Exploratory data analysis: plots, vizⁿ, hypothesis testing
slice and dice data
- 5> Modeling, Evaluation and Interpretation.
- 6> Communicate results: clear and simple to understand
(1 page - 6 page descriptⁿ)
- 7> Deployment: engineering.
- 8> Real world testing: A/B testing.
- 9> Customer / business buy-in, convince them.
- 10> Operations: retrain models, handle failures, process
- 11> Optimization: improve models, more data, more features, code optimization

* VC-dimension :- Vapnik-Chervonenkis

Q) How powerful is a class of models?

Youtube: Programming Cradle

- Linear-models
- boosting algo.
- NNs

VC-dimension:- Statistical ML, computation learning theory
It is a ^{theoretical} way to decide which model is better.
(not much used in applied/practical)

VC-dim :- max no. of pts. that can be shattered / separated by a model for all possible config.

VC-dim (lr-model) = 3 → 3 is the max no. of pts that can be separated by a lr model for all possible config.

VC-dim (RBF-SVM) = ∞

→ if $VC\text{-dim}(M_1) > VC\text{-dim}(M_2)$
then

Youtube: Programming Cradle

theoretically M_1 is powerful than M_2
but doesn't imply same in practical.