

$$\text{Sale price} = Q \cdot y * \text{Sale and } (1 - P), \text{ Cost per unit} = Q \cdot y + \text{Cost out}, \text{ Total cost} = \text{Sale A} + \text{Cost A}$$

$M = \text{Sale A} + \text{Cost A}$

## Bhavuk Gandhi

1) Regression project: Not time series data,  $y$  is cont.

$$a^0 = 1$$

$$0^b = 0$$

$$0^0 = ?? = 1$$

$$0.9^{0.9} = 0.909 \quad 0.6^{0.6} = 0.736 \quad (0.2)^{0.2} = 0.72$$

$$0.8^{0.8} = 0.836 \quad 0.5^{0.5} = 0.707 \quad (0.1)^{0.1} = 0.79$$

$$0.7^{0.7} = 0.779 \quad (0.4)^{0.4} = 0.6963 \quad (0.01)^{0.01} = 0.95$$

$$(0.3)^{0.3} = 0.693 \quad (0.001)^{0.001} = 0.993$$

4) We cannot compute this, we have to use limit

$\lim_{x \rightarrow 0} n^x = 1$ , as  $x$  tends to 0, it approaches to 1.

$$\lim_{a \rightarrow 0} \frac{1}{a} = \infty$$

Derivatives

5) Derivatives

Rate of change

$\frac{d}{dx} f(x)$  infinitesimal change in  $f(x)$  per infinitesimal change in  $x$ .

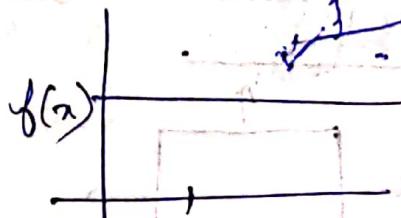
What is the change in  $f(x)$  when  $x$  changes by  $\frac{1}{\infty}$  ( $\rightarrow 0$ )?

$f(x)$  when  $x$  changes by  $\frac{1}{\infty}$  ( $\rightarrow 0$ )

or infinitesimal amount

6)

$$\frac{d}{dx} c = 0$$

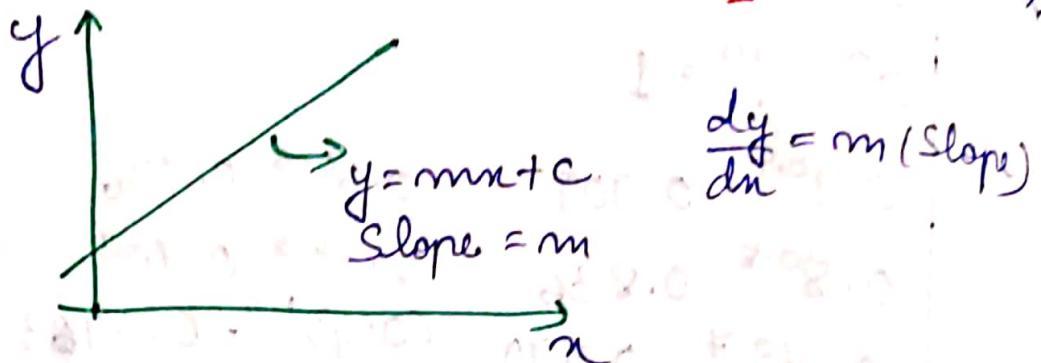
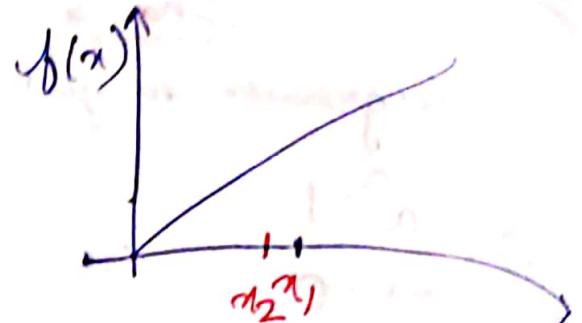


$f(x)$  doesn't change even if we change  $x$  very small amount.

$$1) \frac{d}{dx} f(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}$$

$$2) \frac{d}{dx} (x) = 1$$

$$\frac{x_2 - x_1}{x_2 - x_1} = 1$$



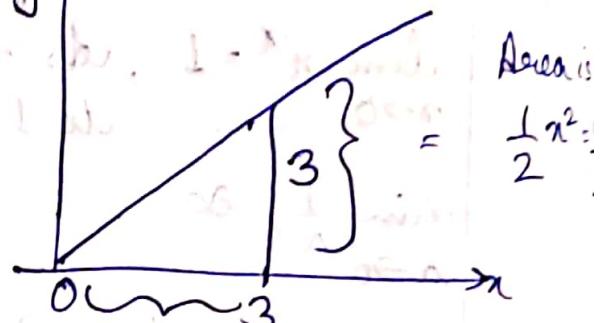
Slope of curve is the rate of change

## INTEGRALS:

$$3) \int x dx = \frac{x^2}{2}$$

$$4) \int_0^3 x dx = \frac{9}{2}$$

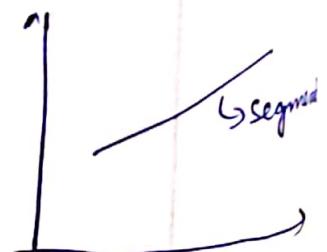
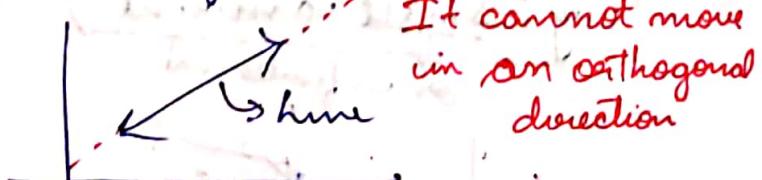
Insp 5)



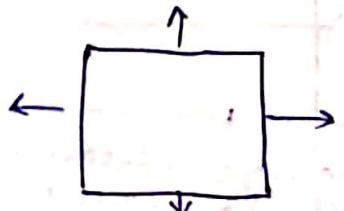
Integrals tells us the area under the curve

## LINES & PLANES

A line can expand either of the two direction



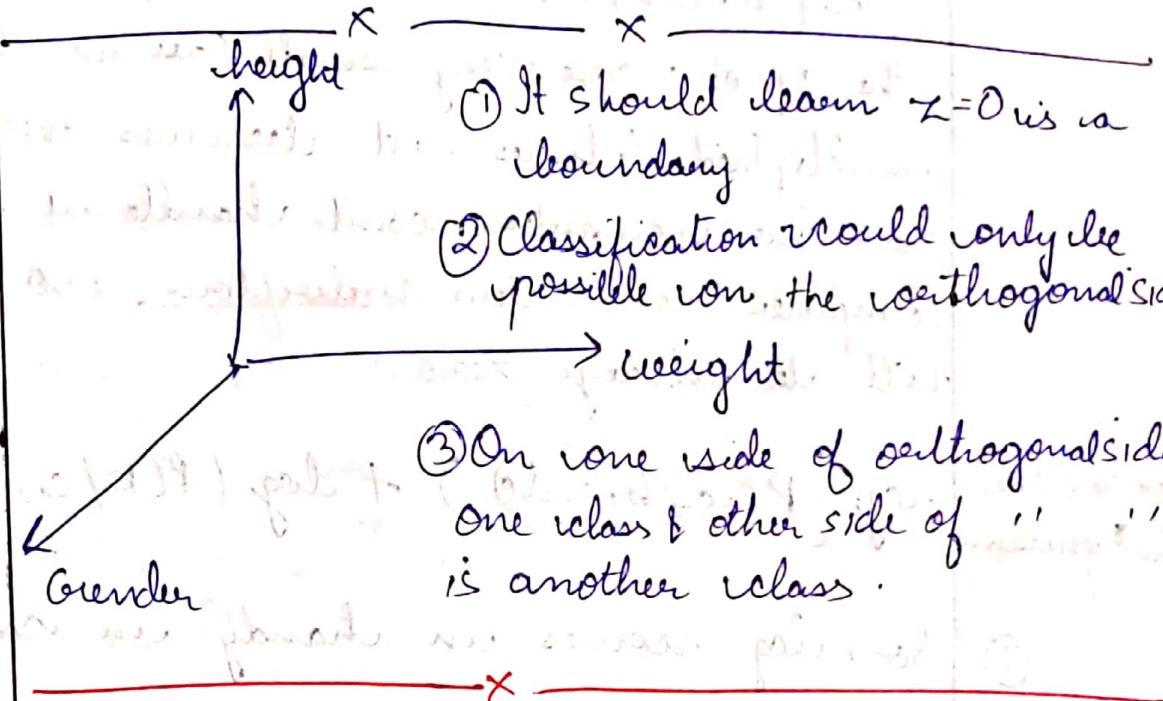
## PLANE:



It cannot move in an orthogonal direction

- ① An  $n$ -dimensional plane can move only in those  $n$ -dimensions, not in other direction.  
 E.g. 5-d can move only in 5 directions.

(2)



## LOGARIT- - HMS

$$\log_{10} 100 = 2 \quad (\text{Representing } 100 \text{ in decimal it takes } 3 \text{ digits } (2 + \frac{1}{5} \text{ always}))$$

$$\log_2 1000 = 9.9 \dots \quad (\text{representing } 1000 \text{ in binary})$$

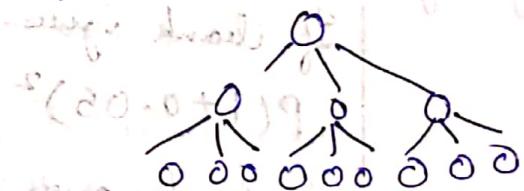
length  
defn:

How much does it takes to represent certain amount of information

$$\log_3 9 = 2$$

$$\log_3 81 = 4$$

number of leaf nodes



(Becoz this is a ternary tree, that is it splits into 3)

1) A very large number, if we represent in log is very small, & very small nos. represented by log makes it large.  
 either way is true that is

$$\log_{10} 0.0000000001 = -100$$

0 followed by 100 zeroes

Chain Rule  
of probability

$$P(a, b, c, d, e) = P(a/b, c, d, e).$$

$$P(b/c, d, e) \cdot P(c/d, e) \cdot P(d/e) \cdot P(e)$$

$P(a/b, c, d, e)$  = is not possible  
the prob. are very small (assume  $10^{-2}$ )  
multiplied 4 times, it becomes very small.  
So, our computer can't handle it. So,  
computer encounters underflow, so the  
will be always zero.

So to solve  
this problem:

①

So, log comes in handy in case of  
NLP task.

Euler's  
No.

e  
(Transcendental  
Number)

Origin of e

~~Step 1~~ ~~dx~~ ~~dx~~ ~~dx~~ ~~dx~~

$a^x = e^x$

$P(1+x)^m$

If bank gives 5% rate of interest for 2 years

$$P(1+0.05)^2$$

If bank gives once a quarter in these 2 yrs

$$P\left(1+\frac{0.05}{4}\right)^8$$

If we say give every month in these 2 yrs

$$P\left(1+\frac{0.05}{12}\right)^{24}$$

Or we inc. this like every hour,  
minute & every second. so, we keep  
dividing the base & multiply the power,

## vectors and Matrices

①

At a certain point it will not increase & stops at a particular number that is 'e'

\_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_

① Vectors : which has magnitude & direction

②  $x, y(x)$

①

Matrices are functions <sup>for</sup> of vectors, which will transform these vectors.

②

They will take a vector & then transform into a vector <sup>& then, it gives the importance { of the features eigen value}</sup>

\_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_

## linear

Eq<sup>m</sup> of line:

$$y = mx + c$$



①

Linear is not necessary to have everything to do with the line.

Eq<sup>m</sup> of plane:

Here, 3-dimensional

$$y = a_1x_1 + a_2x_2 + a_3x_3 + a_0$$

Linear doesn't do with line, it basically deals with plane. (that is always 1)

X X X

Linear will not create a curve, we can only have line & plane

a+b

2 Harmonic Mean  $\frac{2ab}{a+b}$

2 Geometric Mean  $\sqrt{ab}$

When to use what

dist (Const.)  $\frac{d}{t}$  = Speed  
time (A variable)

## Properties of H.M

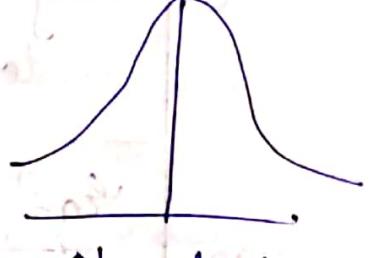
If we have bunch of positive numbers, harmonic mean is smaller than A.M. It tends to smaller values. It doesn't get effected by outliers as much as Arithmetic mean does. (In case of outliers H.M is used)  $\rightarrow$  only in ~~in~~ positive outliers.

Here, Median is expensive so we will use H.M.

## Central Limit Theorem



Mean  $M$   
Mean  $M$   
 $M$



## Std. Error

$\frac{\sigma}{\sqrt{m}}$ , Suppose,  $m = 100$

$$\text{so, } \frac{\sigma}{\sqrt{m}} = 0.1 \sigma$$

if  $\sigma \downarrow$ , then  $0.1\sigma \downarrow$

So, we want  $\sigma \uparrow$ , then  $0.1\sigma$  is significant

original dist is very spread out.

miss out some region of data  
dont use divided by  $\sqrt{m}$ .

## Statistics

### Z-test

### t-test

t  $\rightarrow$  slightly more spread out (student t-dist)

Z  $\rightarrow$  less spread out.

If small sample size is there, we want to have slightly more spread out test, so as to capture all things in that small size sample data. (Moreover when we don't know Std. Dev. of population)

30 sample  
size ??



tails were broaden

tails were almost touch  
to the axis.

Answer: →  
so here  
cutoff was  
arbitrary

They could fit 30 values in a paper;  
They were fitting & creating t-table & z-table  
They didn't use another paper

So, anywhere between 20 to 40 can have z-test  
But in case of very less value z-test  
will perform worst.

## Hypothesis

### Testing

Null

Alternative

p-value: given the results how likely is it that  
the null hypothesis is not true

If p-value ( $< 5\% - \text{Arbitrary std. number used}$ )  
not rejected

It depends some company uses 1%,  
10%.

A/B  
Testing  
(hypothesis  
is used  
most)

1) For instance, we will say we want to  
change the election of buy movie in Amazon.  
Some pop. see it in below

Then see the order to  
see the concession

One scenario conversion was 0.2.

Another " " " 0.2001

Then we use H-Testing (A/B Testing)

What is the prob. of value 0.2001 coming

Here sample size should be sufficient.

P Demography is not considered

~~x~~

~~x~~

Assumption: No Autocorrelation

But in Time series data we must have  
Autocorrelation, so we ~~cannot~~ <sup>use lin. Reg.</sup> use ~~on~~ on  
Time series data

Homoscedasticity: Variance of a variable should be constant

\* The residuals need to be normally dist.

Why? What model would learn from training data  
it would ~~not~~ perform well on testing data.



Hypothesis  $h(x) = \theta^T x$

Cost function

Minimise Error

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (y_i - \theta^T x_i)^2$$

→ L: will give much more minima, so we  
will use only b 2. than square

Why  
Square

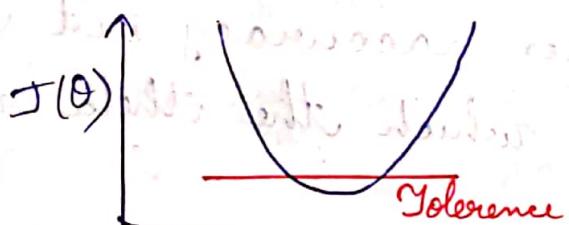
## Gradient Descent

$$\sqrt{25} = \sqrt{-5} \quad i\sqrt{-5} = i\sqrt{5} \times i\sqrt{5} = -1 \times 5 = -5$$

$$\theta_i := \theta_i - \alpha \frac{\partial J(\theta)}{\partial \theta_i}$$

↳ learning Rate  $\rightarrow$  Gradient / Slope = Direction & Magnitude

also  $\frac{d}{dx} = \text{(derivative)}$   $\frac{\partial}{\partial x} = \text{(Partial Derivative)}$   
when we have multiple values



In curve, it takes many close datapoints which is classically a line & we calculate the slope.

## Why learning Rate?

① We can control of how much change we want.

② Without  $\alpha$  value it will never work.

## Where to Stop?

There are 2 variables

If  $\alpha$  is low  
high value  
of tolerance

If  $\alpha$  is high  
tolerance is  
also high

① When  $\theta_i := \theta_i$  are approximately close to each other we will stop. (That is Tolerance)

② Learning Rate

③ Tolerance = If change betw. theta is less than one particular value, i.e. 0.001 then we will stop.

i) Accuracy will be low.

ii) Convergence time will be more.

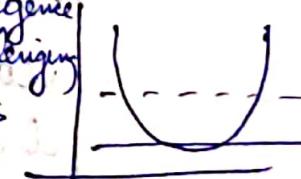
$\alpha$  is high

Tolerance is low

$\alpha$  is low

Convergence is challenging

Acc. is fine



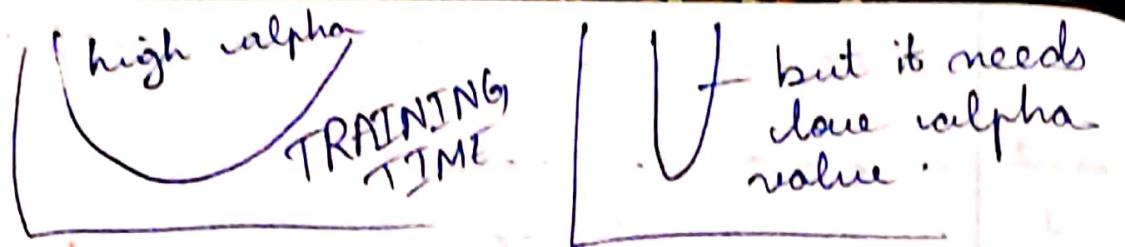
i) I will converge very quickly

ii) Accuracy will be low

$\alpha$  is low  
Tolerance is low  
Extremely convergence

more time to converge  
Acc. is reasonable

## Scaling



- ① Same alpha value cannot work for both the features of different scales
- ② So, this is the reason we rescale the features.
- ③ ~~it describes accuracy, it describes the speed at which the theta will converge.~~ X

## Overfitting

Noise: Is basically the data that is unexplainable by given variable. the variable which is unable to explain the data.

### Source of Noise:

- ① Data input / tracking error
- ② Calibration Error [Ex: Temp. calibration]
- ③ Survey error [Ex: Smoking]
- ④ Unaccounted variables [Ex: lung Cancer patient then a person who hasn't taken cigar, smoking also had shown lung Cancer, due to the asbestos factory to which he was living. That is for our model, it is noisy.]  
Since, we are telling that we are giving all variable, it tries to fit everything. That's why model has noise. We said that our objective is to fit everything so that training error is less.

What stops  
the model  
from  
overfitting

- ① No. of features  $\downarrow$
- ② Degree of polynomial  $\downarrow$
- ③ Complexity of the Model  $\downarrow$

Complexity  $\Downarrow$   
The amt. of the data model  
needs its store to save itself

2 variables, Degree = 2

$$6 \text{ Theta values } \{ \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2 + \theta_6 \}$$

Avoiding

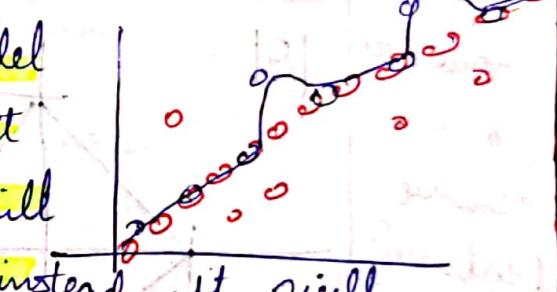
- ① Increasing the size of the training data.

Justify of Above Statement:

Character. of Noise: Random to the model  
It will look  
The model is unable to spot these noise  
patterns, that's  
why noise is random here.

(Ex: Generating pattern for 3 random nos. is  
easy than " " " 30 .. numbers)

On inc. the data pts, model  
will find difficulty to fit  
the noise, so model will  
not go to fit the noise instead it will  
fit the signal those rare many in nos,  
model cannot ignore that signal



- i) ✓ Reduce Complexity
  - ii) " Features
  - iii) " Degrees
  - iv) " Hidden Nodes
  - v) " Tree Depth
- ii) Inc. training data
  - iii) update objective function

Regularization: don't just  
reduce error

instead reduce error plus something.

## Regularization

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i)^2 + \gamma \sum |\theta_j|$$

to differentiate this it have to do conjugate derivative since its not a continuous function

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i)^2 + \gamma \sum \theta_j^2$$

to differentiate this it is easier since its a continuous function

## Lasso Vs Ridge

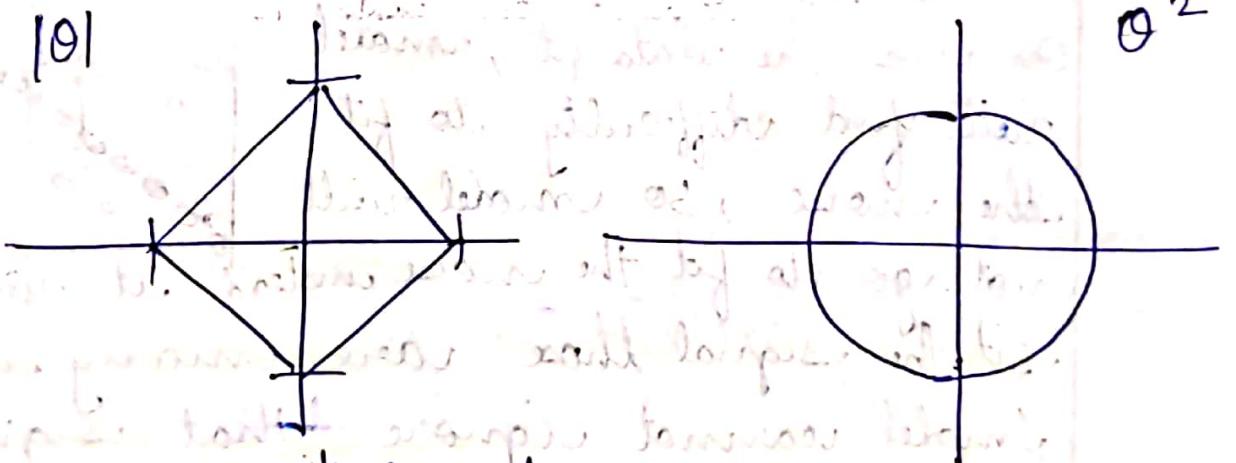
### ② Speeder

Lasso: ① we can use lasso for feature selection

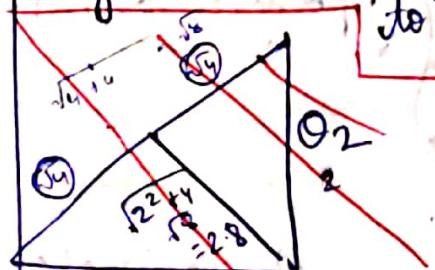
② Nos. of features > Nos. of obs.

③ Sparse: When there is data points that is widely spread in multidimensional with lot of zero element.

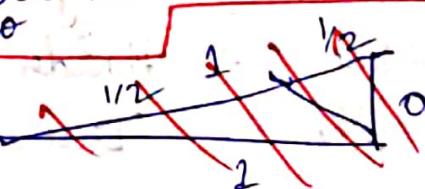
## How does Lasso Remove Features?



Here minimum has only 4<sup>maxima</sup> point. It can pick any value out of this 4 points to make variable zero



Here there is no minima



Ridge  
vs  
Elastic Net

Ridge

- ① less tuning
- ② Speed

Elastic

- ① Better in case of sparse dataset.

- You reduce the dimension its more dense.  $\Rightarrow$  less order you inc. the " " " " sparser the dim.  
So that we make more diff. for model to fit it.
- So internal func. inc. the dim.  $\rightarrow$  so easier to separate data

vectoriza-  
tion

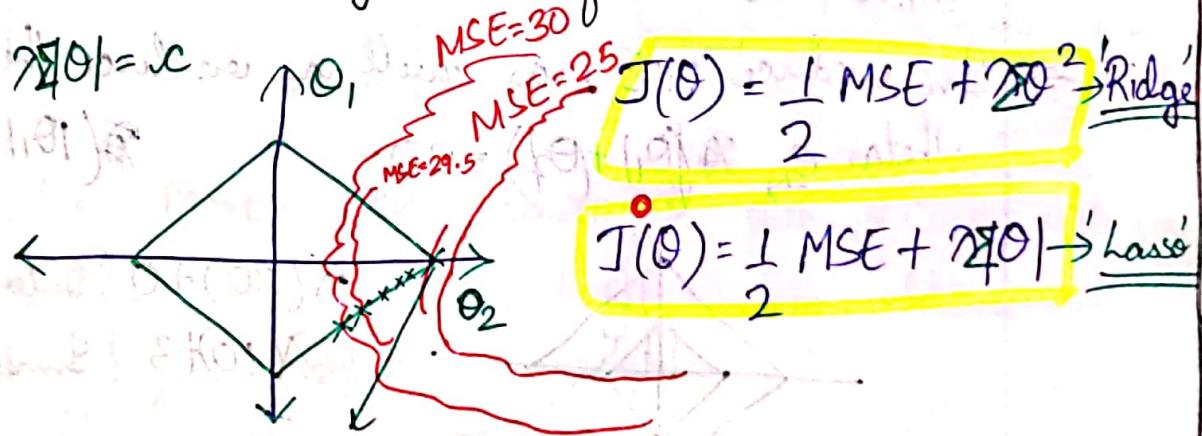
- ① I have to compute space of all used vars.

New Class:

lasso, Ridge, AIC, bic, onehot vs labelEnc, Elastic net, logistic DTree

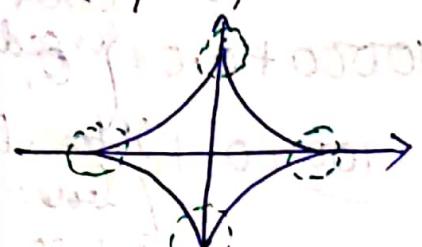
12/08/19

## lasso & Ridge Classification

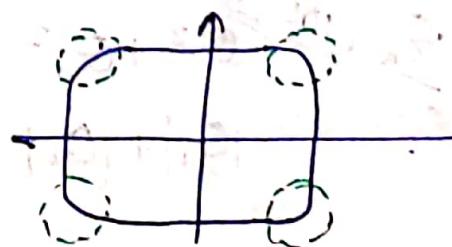


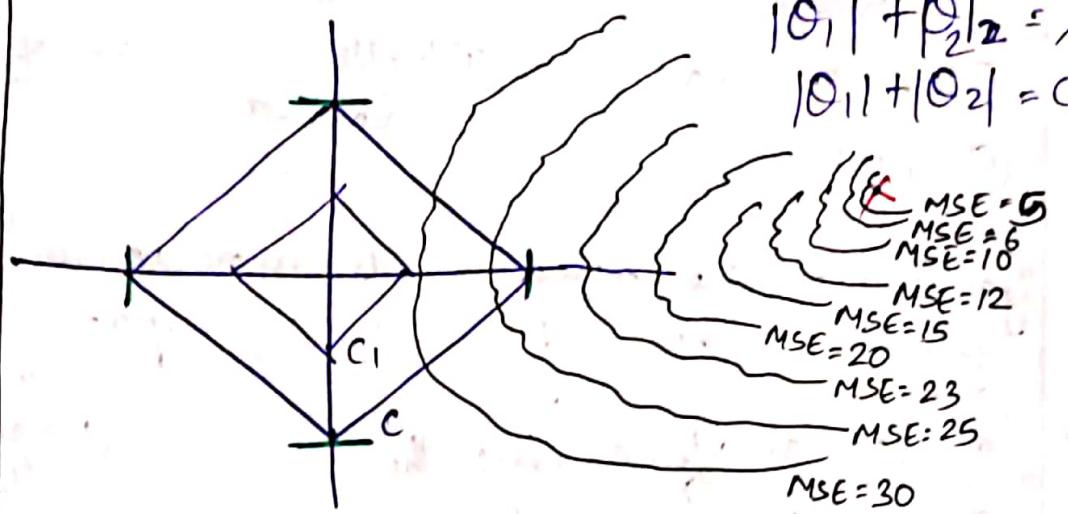
Finally we will stop at this point.

$$C_{0.5} \quad \lambda \sum |\theta_i|$$



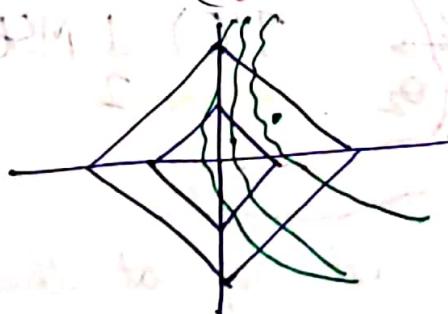
$$\lambda \sum \theta_i^2$$





$\times$ : If I say reach this point then we will end up increasing  $\theta_1$  and  $\theta_2$

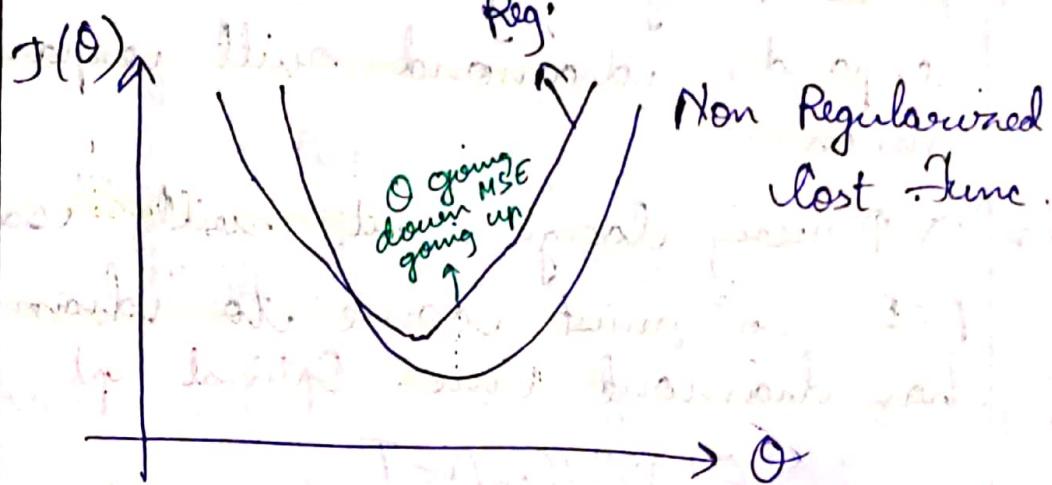
- Some try to keep the contour within the diamond so that it is at the minimum end of the diamond
- Now, there is reg. so we can inc. MSE so that it intersect the diamond
- To reduce MSE, as well as reduce the sum of  $\theta_1$  and  $\theta_2$ ,  $\|\theta_1\| + \|\theta_2\| = 20$



$$\begin{aligned} &\|\theta_1\| + \|\theta_2\| = 20 \quad \{\text{Outer diamond}\} \\ &\|\theta_1\| + \|\theta_2\| = 10 \quad \{\text{Inner diamond}\} \end{aligned}$$

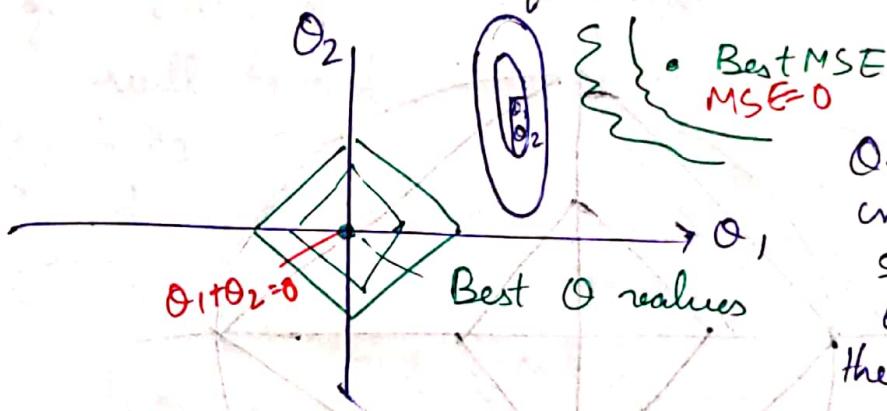
- $L_{0.5}$  = It is more strict, so it will stream  $\theta$  to zero more often than lasso.
- For ridge, it suggests to reduce  $10000 + 0.01\|\theta\|^2$
- For lasso, it does not matter to reduce  $\theta_1$  or  $\theta_2$ .

$$\|\theta_1\| + \|\theta_2\|$$

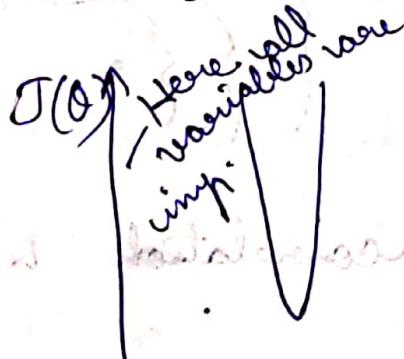


So, when MSE reduces,  $\theta$  goes up  
when  $MSE \uparrow$ ,  $\theta \downarrow$

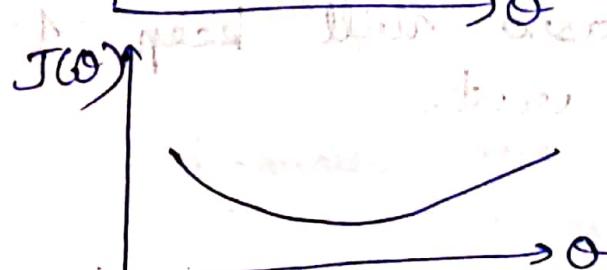
Best case scenario for Theta & MSE



So MSE &  $\theta$  cannot be zero, we  
are inc.  $\theta$  & MSE simultaneously.



We can do regularization here, since  
on red.  $\theta$ , MSE  
shoots up very much  
counterf. so reg. cannot do  
much effect on such  
type of curve.



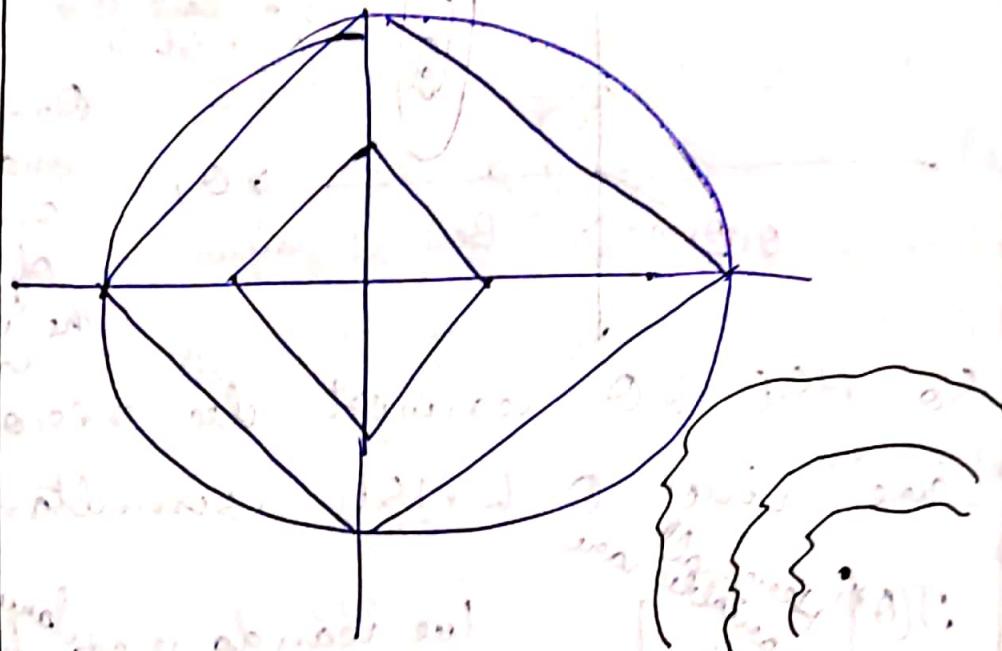
→ Here reg. will  
play more role  
 $\theta \downarrow$ , MSE is hardly  
changing

- Scaling doesn't mean every curve will be same
- If  $\gamma_1$  is very tiny, contours won't expand, diamond will expand more much.
- If  $\gamma_1$  is very large, it will say that MSE must come to diamond, rather than diamond line. Optimal opt. is  $(0,0)$

### ELASTIC NET

$$\frac{1}{2} \text{MSE} + \gamma_1 \sum | \theta_i | + \gamma_2 \sum \theta_i^2$$

↓  
doesn't care Impact large  $\theta$  value



- When features are correlated & there are multiple features

- o Lasso will keep 1, remove the rest.

En: I have Room Size  
 and  $5 \times \text{RoomSize}$   
 Let's say I remove  $5 \times \text{RoomSize}$ .  
 Carpet Area - <sup>Inside</sup> Room carpet area.  
 Built-up = Takes area into the walls.  
 Super Build-up = Net Area  
 Total Apartments.

→ Hasso will remove any 2 variables.

$$1000 \times C.A + 50 \times S.B.U + 100 \times B.U \\ \times 0.2 \times C.A \quad \times 0.2 \times C.A$$

→ The red. in MSE will decrease, but Hasso will turn 2 of them to zero.

→ Ridge will keep these variables

→ E.N. will depend on  $\lambda_1$  &  $\lambda_2$ . For val. of  $\lambda_1$  it will remove all the cluster & keep all the real. of  $\lambda_2$  it will keep the cluster. Here, tuning  $\lambda_1$  &  $\lambda_2$  is difficult to tune.

Akaike Info. Criteria { It is used to compare Bayesian " " models }

BIC assume that one of them is a true model & try to rank in that.

AIC doesn't assume that anything is a true model { Not a single model can be perfect }

" All models are wrong, that is each model will have some relevanceback "

BIC  $\rightarrow$  picks more ~~true~~ model present in the cluster of models.

AIC  $\rightarrow$  doesn't pick.

BIC

Difference

Stricter on nos. of parameters.  
 $\rightarrow$  picks model that is simpler.

AIC

Not stricter.  
Might  $\rightarrow$  picks complex model.

$\rightarrow$  Both checks  $R^2$  & MSE.

$$\begin{aligned} \rightarrow AIC &\rightarrow 2K + \boxed{\ln RSS} \\ BIC &\rightarrow K \ln(n) + \boxed{\ln RSS} \end{aligned}$$

$K = \text{nos. of parameters}$   
 $n = \text{nos. of data pts.}$

If  $m > e^2$ , then  $k \ln(n) > 2K$   
 $\downarrow$  will be always greater,  $e^2 = 9$ ,  $e = (2.737)^2 \approx 9$

So, BIC is always practically greater than AIC.

$\downarrow$  has higher weightage on  $k$ , so it always choose simpler model.

BIC  $\rightarrow$  better in case of overfitting.

AIC  $\rightarrow$  " " underfitting.

So, both are better but people tends more towards AIC since no model is good.

If  $m \downarrow$ , then AIC ends up picking the better model.

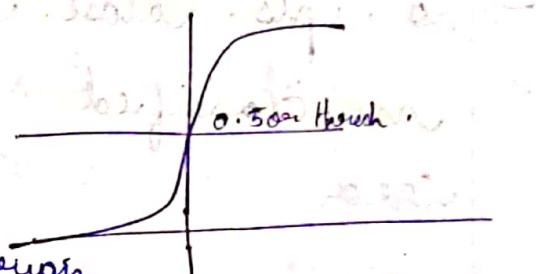
When  $m < k^2$ , always go for AIC.

↓  
Sent data, TF-IDF, features don't represent much.

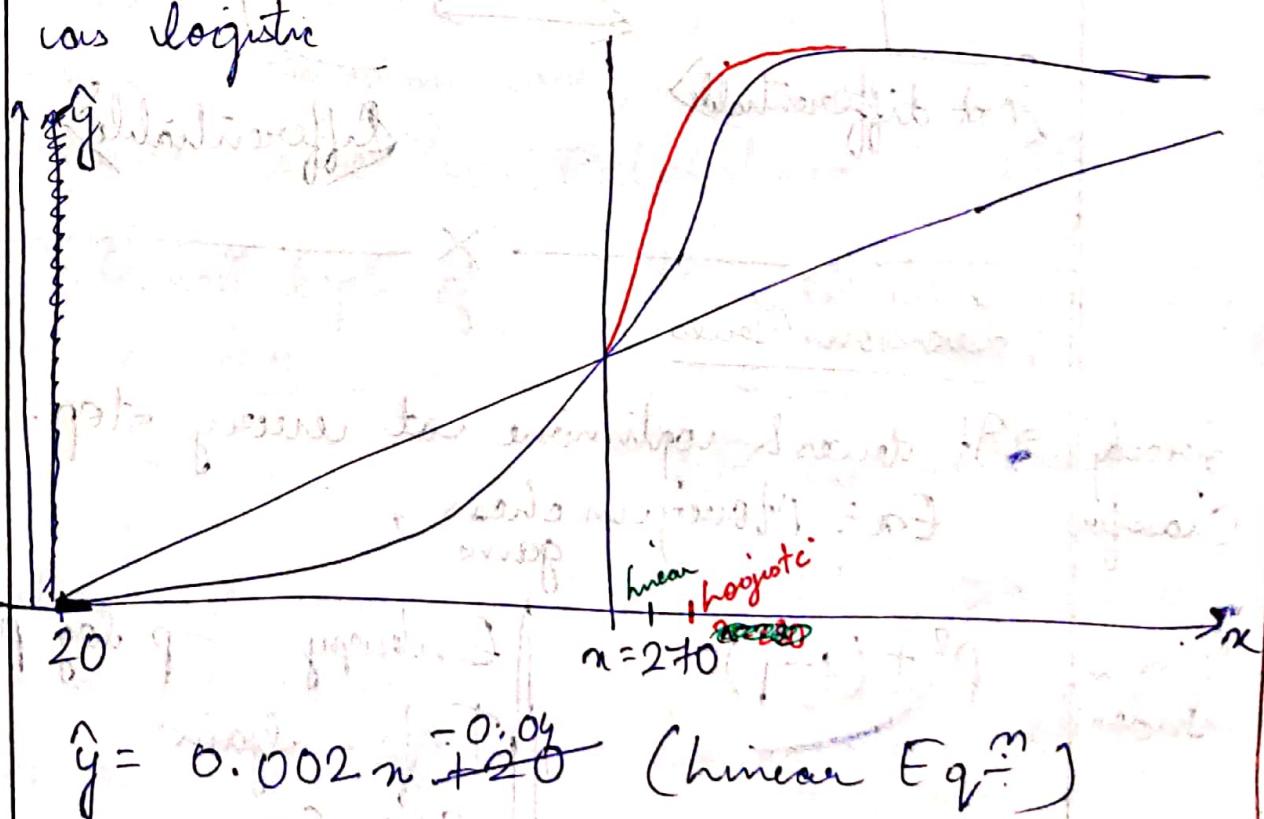
## LOGISTIC REGRESSION

Sigmoid Function:

$$0 \leq \hat{y} = \frac{1}{1 + e^{-\theta^T n}} \leq 1, \text{ formula 1}$$



- if a pt. is close to origin
- In log., even if a point is slightly mis-classified, it will penalize heavily.
  - But linear will not penalize as much.



$$\hat{y} = 0.002n - 0.04 \quad (\text{linear Eq. 1})$$

when  $n = 270$

$$\hat{y} = 0.5$$

When  $x = 280$

$$\hat{y} = 0.52 \quad | \quad y = 0.6 \text{ (Logistic)}$$

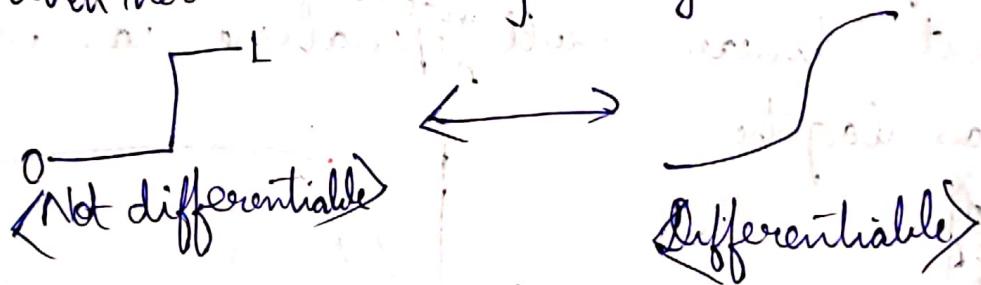
$$\text{When } \hat{y} = 0.6 \quad | \quad \hat{y} = 0.85$$

→ So it allows it to learn much faster.  
Since it induce much higher error in logistic than linear.

→ So, pts. close to origin & if slightly misclassified. Sigmoid induces much high error.

(check)  
Cont. function

When there is arbitrary change in slope.



### Decision Trees

Greedy  
Classific.

→ It iterates & optimizes at every step.  
Ex: Moving in chess game.

Gini  
Index

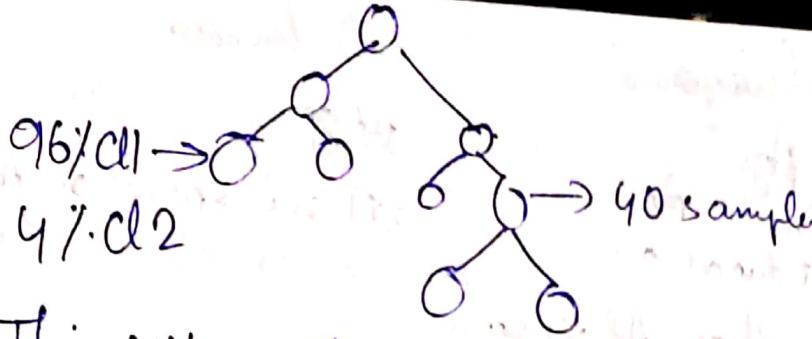
$$p^2 + (1-p)^2$$

$$\text{Entropy} = -p \log p$$

Info. Gain

Chi Sq.

Variance Red<sup>2</sup>



- This 4% could be noise, { Plz don't split when seeing few samples }  
Becoz: it try to reach homogeneity.
- 2 data pt. are identical then acc. is less than 100%.  
Or else in every other case accuracy is 100%!
- So, its imp. to put restrictions on it to remove this overfitting.
- Some restrictions are:
  - ① Node size?
  - ② What is the I.G or Gini
  - ③ Overall depth of the tree.

### DT      vs      kR Logistic Regression

- ① Can fit very complex func. less likely to overfit.
- ② Robust to outliers we have to handle outliers
- ③ OK with missing values
- ④ Missing value will always go into the 'NO' part
- ⑤ Better w/ class skewed

Categorical

- 1) Age
- 2) State Resi.
- 3) Education

Continuous

- 1) Age
- 2) Time spent on Ad.

4) Weather Ad has been clicked or not.

(categorical variable)

continuous variable

Age

<25

24

>25

26, 55

→ In cont. variable do put a boundary is a good idea.  
since, 24 years old person is similar in behaviour to 26 year old person.

Categorical

$$M=0$$

$$Fem = 1$$

Either 0.6 \* 1 (0.6 \* 20 - UP)

or 0.6 \* 0 (0.6 \* 1 - Andhra)

Doesn't Make sense

THETA VALUE  $0.5 \times 20$  No Sense  $0.5 \times 1$  Sense

Cont.

$$0.37 \times \text{Age}$$

$$0.5 \times \text{Age}$$

Age  
12      13

$$39 \times 0.5$$

$$41 \times 0.5$$

$$59$$

Age  
28      29.5

Makes sense  
Similar

→ Logistic suitable to cont. variable.

Logit would work with Cat. but it would be difficult for it to optimise the cost func; whereas in case of DT it is very easy for model to optimise those variables.

D.T can build many type of non-linear boundary very very quickly.

Naive Bayes is very very very very less likely to overfit.

Class skewed is handle at lower levels

Suppose Class 1 - 90% Class 0 = 10%

For above case, first Class 0 will split at, & then so Depth (↑) needed to use the depth class 1 will split starting

Robust to outliers, outliers usually remain split at last Depth (↓)

at end so we need to reduce the depth

D.T	logis R
Better with Categorical var.	Better with Continuous Variables

After Break

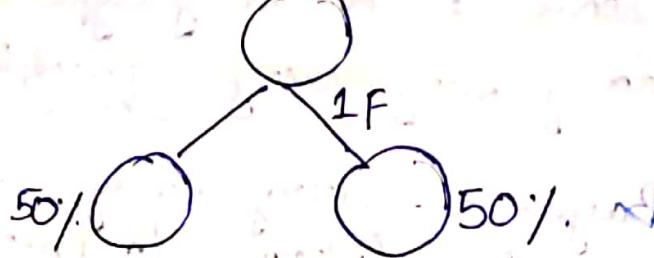
When output variable is continuous

$y = \text{Weight}$

$x$  is Height, Gender or any categorical var

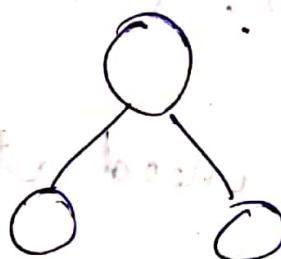
→ Gini Index will not work

Gini Index will not work



Reduction in variance: feature  $\rightarrow$  If feature is reducing the variance of other features  
 If variance is reduced by large amount  
 Ultimate gain: There should be zero variance

### Classification ( $y = \text{late.}$ )



→ Prob. how to decide  
 Statement from which values to split

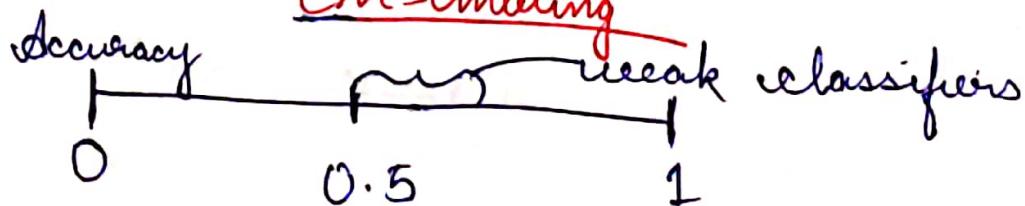
where to create boundary when we want to analyze the time spent on social media  
 values  $0, \dots, 10^6$  go from:

Gini vs ChiSq vs Entropy

Faster log is slower, complex expensive computation

Result for all three, will mostly be same.

### Ensembling



weak classifier: better than random guess.

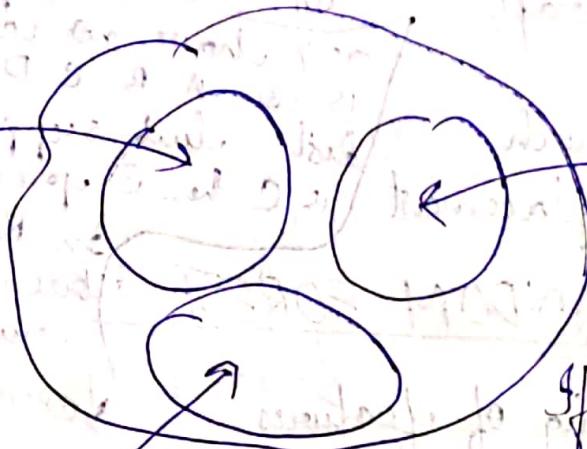
Ex: If tossing a coin 5 times then  $p \geq 0.5$

$C_1$  is

wrong

$C_2$  is

wrong



If combining all  
 $C_1, C_2, C_3$  we will  
get 100% accuracy.

50% i.e.,

$C_1$  is

wrong

but  $C_2 \cap C_3$

is correct

$C_1 \& C_2$  is wrong

$C_2$  is

wrong

Correct 40% of time

$C_2 \& C_3$  is wrong.

$C_3$  is wrong

When intersection takes 50% of area,

then the classifier is not so good.

75% correct

$C_1$  is

wrong

correct

$C_2$  is

correct

Correct 15% of time

At least 2 needs to  
be correct which  
is shaded area for  
the entire classifier  
to be correct.

$C_3$  is correct

Cond: At least each classifier should  
get correct value 50%

In  
ENSEMBLE  
COND

# BAGGING

pt 3) Explanation

Bootstrap

Aggregation

Sampling with  
Replacement

RANDOM FOREST

In KBC, suppose audience poll 10% of pop. knows ans.

90% have no clue what the is so A B C D has each 22.5% dist' but suppose correct ans is C then C prop. is 32.5% resp.

In RF also, suppose tree have rule & 50% have its tree, so it goes to tree taking less overfit

① Sampling of features

② Sampling of records / data points

③ Majority vote / depth, mean ↓ overfit

Reducing features have much impact on overfitting than reducing sample of data.

→ Features Every tree is exposed to something from training data.

Each tree will overfit to data that it has seen.

→ If 20% of data has been used by 1-D-Tree 1 data pt. will seen in the 20% of the data.

On an average each point is exposed to 20% of tree. So it has seen 20% of data & 80% of data it has not seen so it will not overfit.

→ If we are exposing 80% of data then we are also reducing the effect of 2nd point that is we are not exposing 20% of data.

So, in this case also it will manage overfitting.  
Bagging loose out on speed, interpretability, reduces overfit

### Rand. Forest

### Decision Trees

- |                                                                                        |                                        |
|----------------------------------------------------------------------------------------|----------------------------------------|
| ① Less likely to overfit                                                               | Interpretable                          |
| ② Handles Noisy labels                                                                 | Speed<br>(Especially prediction Speed) |
| ③ Tells Feature Importance<br>(It checks what high feature is present in each D. Tree) | Less likely to underfit                |

Noisy labels: Fraud detection problems

have noisy labels.

Transaction Fraud - Small lag of time (Anomaly Detection)  
Human Fraud - Large lag of time (Sup. learning)

Let's say we took last 6 months of data (May-Oct) we give to model.

Here we make which transaction is fraud or not.

On 25 root. fraud happens, but in data -  
these 25 rows has ~~not~~ fraud, so it will show <sup>not</sup> fraud on 25 Oct in model.

→ None Deep learning do this.

Auto Encoders we don't use labels.

→ Because of leaving 20% data we are exposing, so it is ~~not~~ handles noisy labels.

→ RF is worst since one rare reducing feature  
RF will underfit most.

En, 10000 data pts. & 3 features.  
RF will take 2 features & rebuild D.T

n nos. of times but none of the time,  
D.T is able to predict the labels.

WHEREAS, D.T will take all 3 features  
10,000 data pts & will be able to  
predict the labels.

## BOOSTING

Suppose we have bunch of three  
classifiers:  $h_1, h_2, h_3$ ,

$$h_1, h_2, h_3 \xrightarrow{+1} \xrightarrow{-1}$$

$$h(x) = \text{sign}(h_1(x) + h_2(x) + h_3(x) + \dots)$$

What if I boost all misclassified datapoints  
in our next classifier.

$\sum w_i = 1$ , lets scale the weights so that its  
easier to calculate

Now boost stronger classifiers

$$h(x) = \text{sign}(\alpha_1 h_1(x) + \alpha_2 h_2(x) + \alpha_3 h_3(x))$$

## Steps

① Starts with

$$w_i = \frac{1}{N} \rightarrow \text{pick } \alpha \text{ that minimizes } E_{\alpha}$$

$w_{ii} = \frac{1}{N} \rightarrow$  Pick  $h_t$  that minimizes the  $E_t$

Ada Boost

Calculate  
 $w_{it}$

Pick  $\alpha_t$

$$w_{ii}^{t+1} = \frac{w_{ii}^t}{Z} e^{-\alpha_t h_t(x_i) y_i}$$

$$\alpha_t = \frac{1}{2} \ln \frac{1 - E_t}{E_t}$$

$$e^{-\alpha_t (+1)}$$

when pred. is  
Same as model

when  $h^t(x_i) = y_i$   
when  $h^t(x) = 1, y_i = 1$   
or  $1 = -1, y_i = -1$

if correct

when opp. pred.

$e^{-\alpha_t (-1)}$   
if wrong

$$w_{ii}^{t+1} = \frac{w_{ii}^t}{Z} e^{\ln \sqrt{\frac{E_t}{1 - E_t}}}$$

$$= \frac{w_{ii}^t}{Z} e^{\ln \sqrt{\frac{1 - E_t}{E_t}}}$$

$$\Rightarrow \frac{w_{ii}^t}{Z} \left( \sqrt{\frac{E_t}{1 - E_t}} + \sqrt{\frac{1 - E_t}{E_t}} \right) = 1$$

$$\Rightarrow \frac{1}{Z} \left( \sqrt{\frac{E_t}{1 - E_t}} \sum w_{ii}^t + \sqrt{\frac{1 - E_t}{E_t}} \sum w_{ii}^t \right) = 1$$

$$\Rightarrow Z = 2 \sqrt{E_t(1 - E_t)}$$

~~$$\sum w_{ii}^t = \frac{1}{Z} \sqrt{\frac{E_t}{1 - E_t}}$$~~

$$\sum w_i^{it} = (1 - \varepsilon)$$

correct

$$\sum w_i^{-it} = \varepsilon$$

wrong

$$\text{So } z = \frac{\sqrt{\varepsilon t}}{\sqrt{1-\varepsilon t}} \times (1-\varepsilon t) + \frac{\sqrt{1-\varepsilon t}}{\sqrt{\varepsilon t}} (\varepsilon t)$$

$$\text{So we get } z = 2\sqrt{\varepsilon t(1-\varepsilon t)}$$

$$w_i^{it+1} = \frac{w_i^{it}}{z} \sqrt{\frac{\varepsilon t}{1-\varepsilon t}} \quad \text{if correct}$$

$$= \frac{w_i^{it}}{2\sqrt{\varepsilon t(1-\varepsilon t)} \sqrt{\frac{\varepsilon t}{1-\varepsilon t}}}$$

if wrong

$$= \frac{w_i^{it}}{2\sqrt{\varepsilon t(1-\varepsilon t)} \sqrt{\frac{1-\varepsilon t}{\varepsilon t}}}$$

$$\frac{w_i^{it}}{2\sqrt{\varepsilon t(1-\varepsilon t)} \sqrt{\frac{1-\varepsilon t}{\varepsilon t}}}$$

$$= \frac{w_i^{it}}{2(1-\varepsilon t)}$$

$$\frac{w_i^{it}}{2\varepsilon t}$$

$$\sum w_i^{it+1} = \frac{1}{2}, \quad w_i \text{ when correct is } \frac{1-\varepsilon i}{1-\varepsilon i}$$

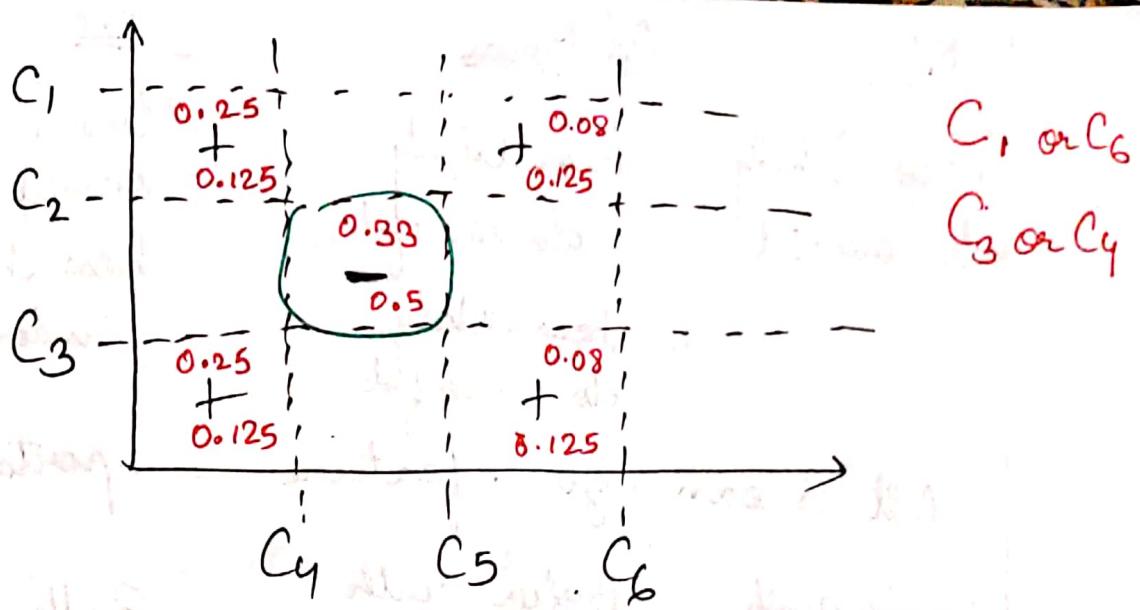
$$\text{So } \frac{1-\varepsilon i}{2(1-\varepsilon i)} = \frac{1}{2}$$

$$\& \sum w_i^{it+1} = \frac{1}{2}, \quad w_i \text{ when incorrect is } \varepsilon i$$

$$\text{So } \frac{\varepsilon i}{2\varepsilon t} = \frac{1}{2}$$

This works if we have more 50% correct data so that it weights correct as  $\frac{1}{2}$  & incorrect as  $\frac{1}{2}$

→ Suppose 80% of correct & 20% of incorrect data is there, so it gives 80% of desired & correct 20%



Let's say, I have some classifier  $h(x_i)$  GBM's

$$y_i^t = h^t(x_i) + \text{err}^t$$

$$\text{err}^t = h^{t+1}(x_i) + \text{err}^{t+1}$$

$$\text{err}^{t+1} = h^{t+2}(x_i) + \text{err}^{t+2}$$

At some step, error becomes very small.

$$y_i = h^t(x_i) + h^{t+1}(x_i) + h^{t+2}(x_i)$$

→ So, in GBM's it tries to correct the error from the previous one.

## RF

Less likely  
to overfit

works with  
noisy labels  
Since we use D-Tree, so it can  
handle  
Missing Value,  
Outliers

## AdaBoost

Less likely  
to overfit  
Less likely  
to underfit

All 3 can give feature importance

Better with  
class skew.  
Handle

Better with  
class skew.

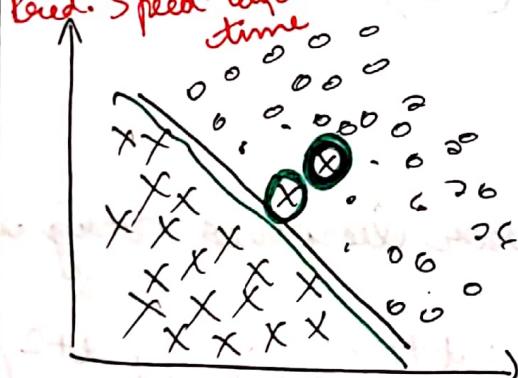
Depends on what type of  
underlying classifier we are  
using.

Training Speed  
When parallel  
Pred. Speed takes time

Prediction  
Speed

Prediction Speed.

< Logistic is faster than all >



→ Ada boost

Overfitting is when it doesn't  
perform well on unseen  
data. But since we

do separate model, Ada  
Boost does well on Test or unseen data.

## XGBoost

GBM + Regularization + Speed.

Reduces overfitting, & its fast.

## Problem with Boosting:

We need lot of hyperparameter tuning.

Day 3: 21-Dec-19

Till now we have studied Discriminative Generative model, today we will learn v/s Generative model.

⇒ In Generative Model we will learn Joint probability  $P(c, d)$

Joint prob of c and d is

$$P(c, d) = P(d/c) \cdot P(c) \Rightarrow P(c/d) \cdot P(d)$$

$$P(c/d) = \frac{P(d/c) \cdot P(c)}{P(d)}$$

$$\text{Max}(P(c/d)) = \text{Max}(P(d/c) \cdot P(c))$$

$P(d/c)$  = Let's say I have 5 features

$$P(x_1, \dots, x_5 | c)$$

Joint prob of all the features given the class

Now expanding this thing:  $P(x_5 | c)$

$$= P(x_1, x_2, x_3, x_4 | x_5, c) \cdot P(x_1, x_2, x_3 | x_4, x_5, c)$$

$$P(x_4/x_5, c) P(x_1, x_2 | x_3, x_4, x_5, c) \cdot P(x_1/x_2, x_3, x_4, x_5, c) \cdot P(x_3/x_4, x_5, c)$$

This is very difficult to compute since the probability gets smaller and smaller.

$$P(x_2/x_3, x_4, x_5, c)$$

So, we assume that all the features are independent of each other given the class C.

Assumption

Complexity  
of Naive Bayes  
if there's  
no assumption

$O(1 \times 1^m \cdot |C|)$  parameters.

So, after assumption we can write as,

$$P(x_1, \dots, x_m | C) = P(x_1 | C) \cdot P(x_2 | C) \cdots \frac{P(x_m | C)}{P(C)}$$

$$C_{MAP} = \underset{C \in C}{\operatorname{argmax}} P(x_1, x_2, \dots, x_m | C) P(C)$$

$$\operatorname{argmax}_{C \in C} P(C_j) \prod_{x_i \in X} P(x_i | C)$$

Word positions in that doc.

$$C_{MB} = \underset{C_j \in C}{\operatorname{argmax}} \prod_{w \in V} P(w_i | c_j)$$

1st Attempt: Maximum likelihood estimates  
Simply use the frequencies in the data.

$$\hat{P}(C_j) = \frac{\text{doc-count}(C = C_j)}{N_{\text{doc}}}$$

$$\hat{P}(w_i | C_j) = \frac{\text{count}(w_i, C_j)}{\sum_{w \in V} \text{count}(w, C_j)}$$

→ Zero prob. cannot go away, it's a big challenge, so people use smoothing.

What if we have seen no training doc. with the word "fantastic" & classified in the topic "positive".

$$\hat{P}(\text{"fantastic"} / \text{positive}) = \frac{\text{count}(\text{"fantastic"}, \text{positive})}{\sum_{w \in V} \text{count}(w, \text{positive})}$$

① ~~replace smoothing~~: Just adding one in the numerator & the denominator.

- vantages
1. Very Very Fast.
  2. Low Storage Requirements.
  3. Extremely robust to irrelevant features.  
(Least likely to overfit than e.g., D.T, Boosting, less than Random Forest, its a benchmark)

Ex: Whether somebody going to cancer or not  
and if it will or not, so on that particular day

Sol: Both Cond<sup>1</sup> has no correlation

$$\text{So, } P(R=Y/C=\text{Cancer}) = 0.3 \quad P(R=1)$$

$$\text{So, } P(C=Y/R=\text{Not Cancer}) = 0.3$$

So, reweightage for both the classes will be equal, so no overfitting.

4. Both of the features are given equal importance but D.T suffers from fragmentation in such cases - especially if little data.

Ex: Whether an email is spam or not-spam.

They will just repeat the subject. So, if a particular feature is imp. just duplicate it & give 6 features if 5 features is there - whether or not. There is 100% correlation, it will work.

Indeed by giving certain importance to  
 - that particular features whereas all  
 so ML model will struggle like correctly  
 → It works really well in very high dimen.  
 datasets, so high dimen. datasets like text  
 like 1 millions docs, if I take 10K for  
 bag of words & do n-grams & n-grams  
 So that I have 1 lakh features. It's  
 very difficult to figure out right points  
 when there are lot of features.

~~Discriminative~~  
 It learns better & lot faster, less likely to underfit

It generalizes very very well that is it will never overfit

→ Predicting stock price - Difficult epochs. (use discriminative classifier)

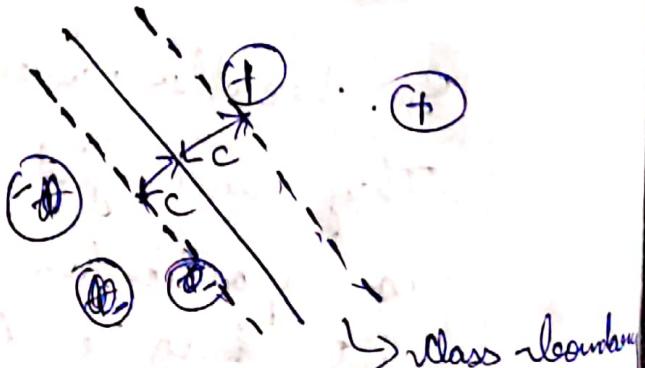
→ When it's classify task - Simple epochs. (use generative classifier)

## SVM

- ① The distance between points & boundary is as high as possible as in logistic func. is no such rule.

$x_i$  = data pt.

$w$  = the vector that allows me to determine the dis. b/w boundary vector  $w$  & decision boundary.



$$\vec{w} \cdot \vec{u} + b > 0$$

$$\vec{w} \cdot \vec{u} + b > 0, \text{ here } b \text{ is constant}$$

For logistic, the objective is to just create a boundary or just separate the data.

Now separating tree and -ve data points

$$(\vec{w} \cdot \vec{u} + b) \geq 1 \quad \text{for +ve}$$

$$(\vec{w} \cdot \vec{u} + b) \leq -1 \quad \text{for -ve}$$

Combining both the eq:

$$y_i(\vec{w} \cdot \vec{u}_i + b) \geq 1$$

Outcome either +1 or -1

$$\vec{x}_+ = (-b) \frac{\vec{w}}{\|\vec{w}\|}$$

$$\vec{x}_+ - \vec{x}_- = \frac{2\vec{w}}{\|\vec{w}\|}$$

$$\vec{x}_- = (+b) \frac{\vec{w}}{\|\vec{w}\|}$$

So, we want to maximize the Margin,

by maximizing  $\frac{1}{2} \|\vec{w}\|^2$  with

$y_i(\vec{w} \cdot \vec{u}_i + b) - 1 \geq 0$  as my constraint

→ having say I can combine these 2 eq & minimize this

$$\frac{1}{2} \|\vec{w}\|^2 + C (\sum y_i (\vec{w} \cdot \vec{u}_i + b) - 1)$$

Margin

Error point

If violated

$$1 \|\vec{w}\|^2 - C (\quad " \quad )$$

I want to min  $\frac{1}{2} \|w\|^2$

I say that the error is not really that is if its not violated then it will always be zero.

If error is -ve then net is -ve.

$\alpha$  is a weightage which give weightage to some part.

If it is violated then we will multiply by a weightage

If  $\alpha$  is very very large, then will never get violated

$$\frac{\partial L}{\partial w} = 0$$

$$\vec{w} - \sum \alpha y \vec{x} = 0$$

$$\vec{w} = \sum \alpha y \vec{x}$$

$$\frac{\partial L}{\partial b} = \sum \alpha y = 0$$

$$L = \frac{1}{2} \sum \alpha_i y_i \vec{u}_i \cdot \sum \alpha_j y_j \vec{u}_j$$

$$= \sum \alpha_i y_i \vec{u}_i \cdot \sum \alpha_j y_j \vec{u}_j > 0$$

I did several linear or for linear class.

$$L = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j (\vec{u}_i \cdot \vec{u}_j)$$

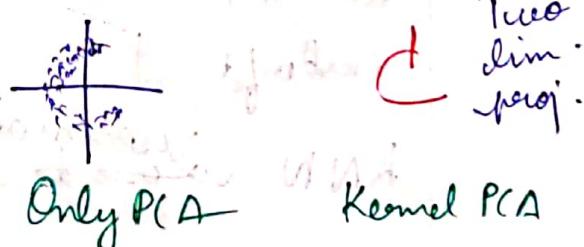
Now, I can build a kernel func<sup>n</sup>, so that I can transform its new dim, so I need to transform the dot products to the transform space.

→ new dot prod. of every comb. of our data

$$K \Rightarrow (\vec{u} \cdot \vec{v} + 1)$$

→ It is better in high dim. data (like text data). So, kernel func. does this. On nuc. dimen. it will not overfit since its method is to minimize the margin

→ PCA always do linear comb. of features that only we use. kernels in PCA



disadv-

- ① Slow, need to figure out right kernel, need to do tuning.

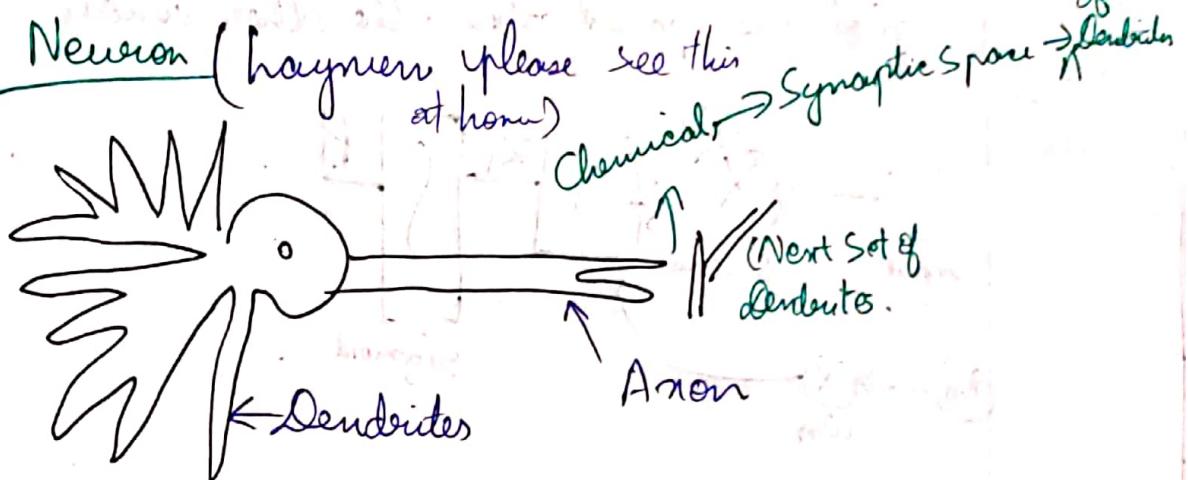
## Neighbourhood based Prediction (KNN)

- 1) It is computationally expensive when 1 Million data pts are there.
- 2) So, there are bunch of tree stems. <sup>vantage point</sup> KDT<sub>Tree</sub>, VP Tree. These trees allow to calc. dist. ~~very quickly~~.
- 3) So after mapping faces to see whether a face is matched cor. or not.  
Ex: In Shadi.com, criminal faces are not allowed. If face matches to a person, it goes to CRO [he says give id & proof.]
- 4) It is still slower but with these VP trees & all we can make it faster than RF & SVM.
- 5) It is very easy quick to underfit. Depends on If to check the temp. condition in Mumbai. If I keep  $K=1$  Million, then it practically gives the avg. of the entire dataset (Underfit)  
If  $K=1$ , it will take immediate neighbourhood point. (Overfit)

$K=3, 4, 5$  are good beyond that it will underfit

KNN ~~Voronoi~~ diagram ??

# Deep learning

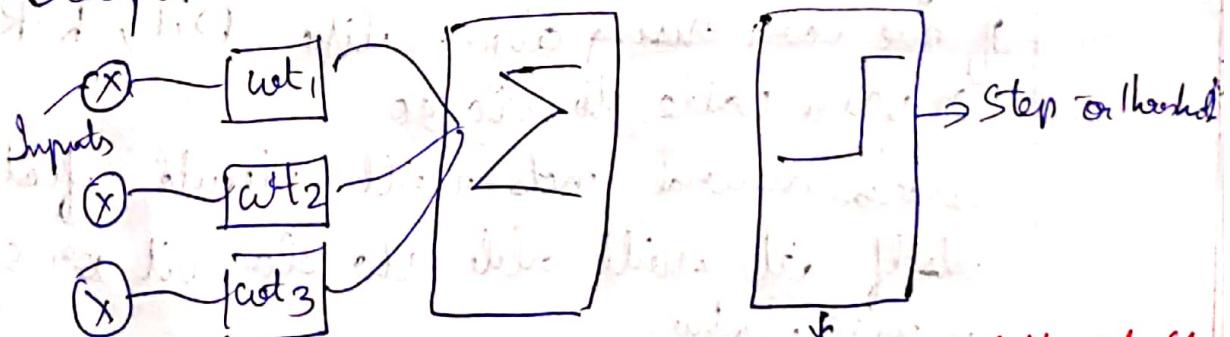


2) If no threshold is not crossed, no signal will pass  
"All or None"

3) So, there are dendrites which are closer to a particular Axon

Cumulative Influence: Multiple neuron triggering  
So multiple neurons cumulatively influence a reaction.

4) Synaptic Weight: A particular dendrite will have certain amt. of weight for a particular receptor.



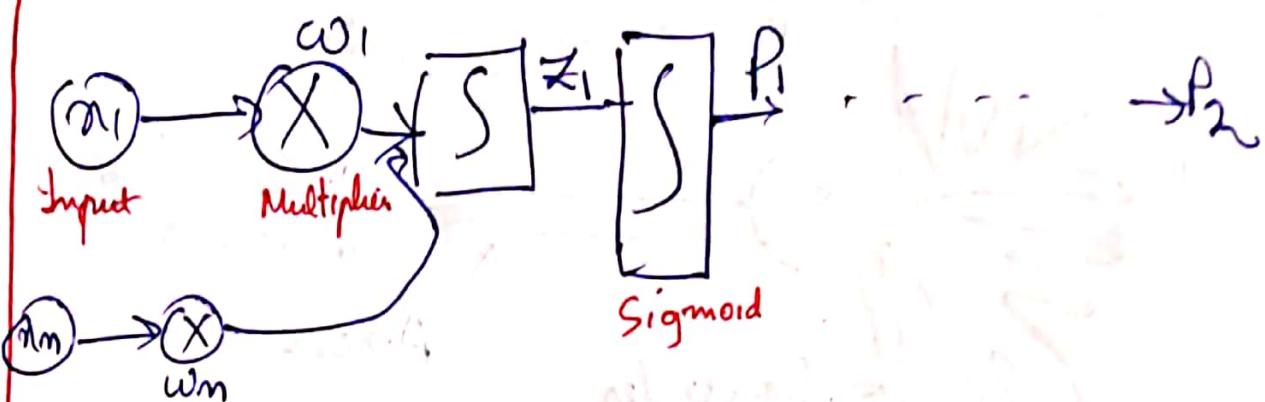
This was not differentiable

so they replaced with sigmoid.

Refractory Period: After one signal is passed it needs sometime to for other signal to pass

In neural nets no such concept is there

Axonal Bifurcation: That is if one axon more electrical activity than in one signal is passed then the other axon.



$$\frac{\partial P_2}{\partial w_2} = \frac{\partial P_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2}$$

$$\frac{\partial P_2}{\partial w_1} = \frac{\partial P_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial P_1} \cdot \frac{\partial P_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

Linear in Width Depth

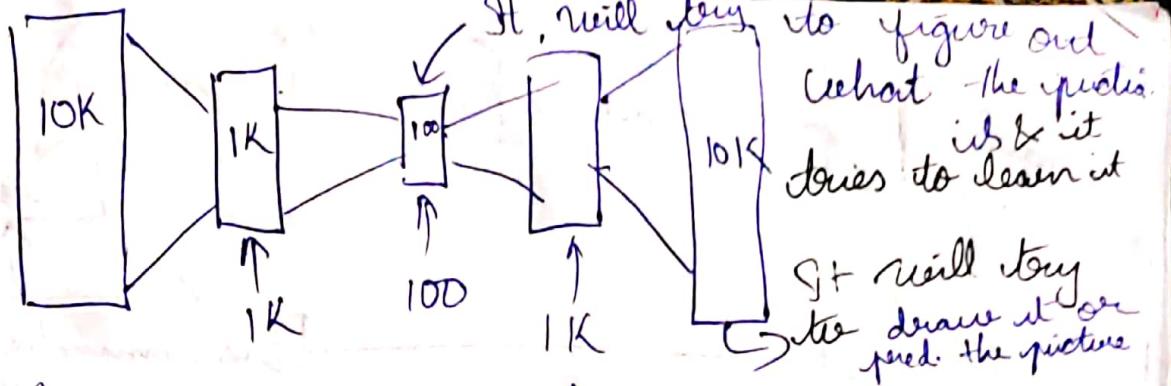
Quadratic in Depth Width

Rm. 1) In Shaddai.com suppose a person is from Punjab & looks懦弱 he will have a tough time if he leaves in Chennai.

If we were using algo. like DT, LR, SVM it can't do so.

Whereas neural nets will build features by itself it will able to do it or solve the above feature.

## Auto Encoders



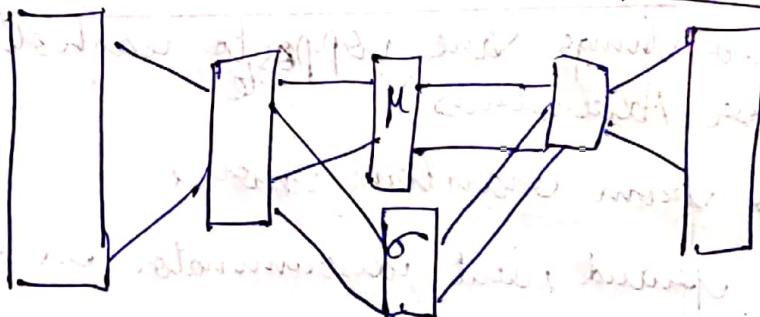
If I give 10K values, then give me 10K value at output

The problem here is it cannot remember all the values.

### Denoising Autoencoder (DAE)

- 1) It will reconstruct from partial data
- 2) Make model robust to <sup>noisy</sup> data

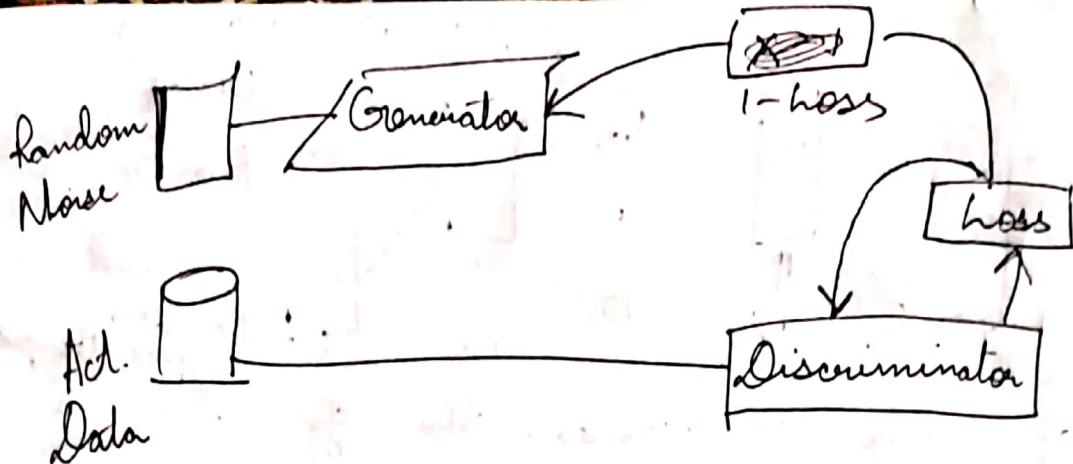
### VAE (Variational Autoencoder)



### Used in Anomaly Detection

- How? If S.D. falls beyond a particular S.D. like 3...5, we will reject it.
- It checks mean & SD & then learns from the mean & SD generated.

### GANs (Generated Adversarial Networks)



Let's say I am generating faces ( $100 \times 100$ )  $\rightarrow$  give this to discriminator, & I tell him that if you predict that this is a real then this is wrong.

<sup>Prop</sup>  
Discriminator will be better in segregating which is real & non-real.

Generator will be better if it produce data which is very very close to real.

Adversarial: Two things are opposite to each other (two competitions or Adversaries).

Generator learns from random noise

Generator make fraud, but discriminator is able to discriminate b/w fake & real.

Challenge: If one model is better than other (like Dis. is very very powerful than Generator) then the system stops learning.

Powerful discriminator able to  
generator will learn since everytime loss will

be large

Powerful generator

Discriminator could be able to distinguish b/w fake & real & moreover loss is same everytime

## CNN



Conv.      Pooling      FCN (Fully connected net)

Here there is func. that is applied to every part of our data.

like detection edges, we get certain output

Then we apply pooling (typically max pooling)

→ what will happen if pooling layer is removed.

It keeps what is up.

Percent Overfit.

There is significant impact on computational time.

→ Conv. will detect edges.

→ why Max pooling is used?

Pooling lay itself is useless.

Justify what's happening in Conv., becoz conv. is detecting the prominent parts of image.

detecting the prominent features is detected

$$\begin{matrix} 8 & 2 \\ 3 & 4 \end{matrix}$$

By itself it's useless, so it's not the 1st layer everytime.

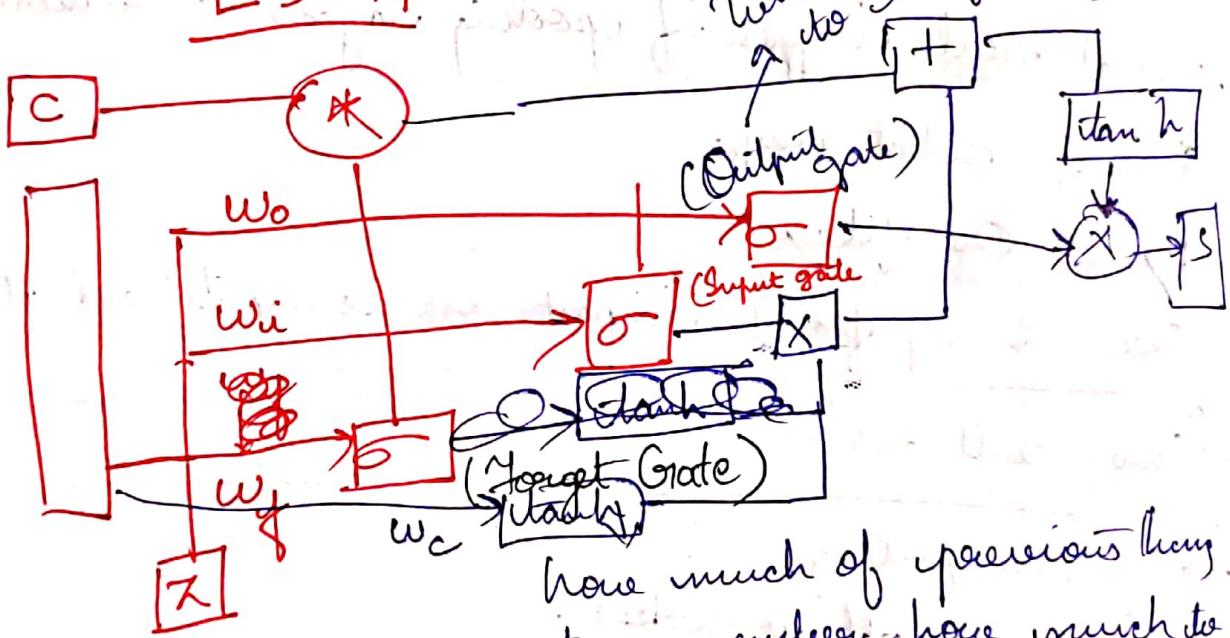
→ So pooling pick most prominent pieces from the conv. layer.

- different kernel is used by the model to detect the edge.
  - Kernel is in motion
  - " will give the size of the biggest edge
  - Once it is done moving pooling decides whether to carry forward the edge or not etc.
- Dropout layers: (to reduce overfit)

Neural Network

- Dropout layers try to remove local optima
- Enable to achieve global optima

## LSTM



how much of previous things to remember, how much to forget  
 The food is good but (forget things)  
; (semi colon forget)

I/P gate: What from current world to be remembered.

Q. When Target Crate B is removed?  
I provided the ability but this doesn't mean that it will forget. It's just a architecture.  
Whether a model uses all of that depends on the type of training data.

### Activation Function

Sigmoid

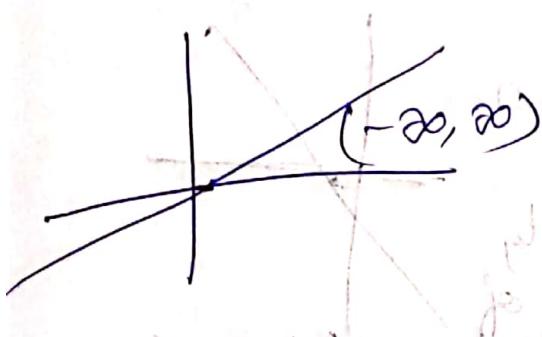
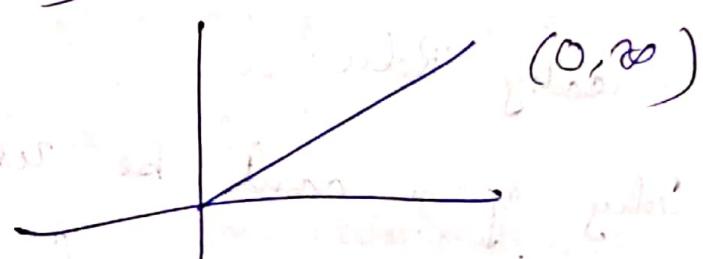
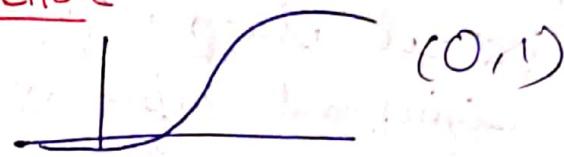
Tanh

Relu

LRelu

PRelu

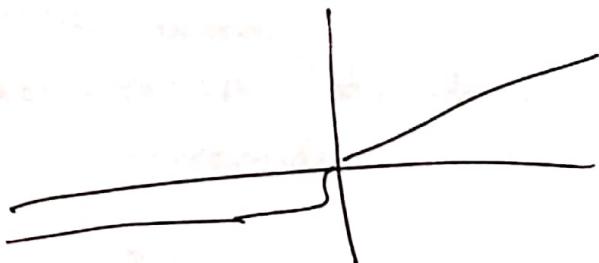
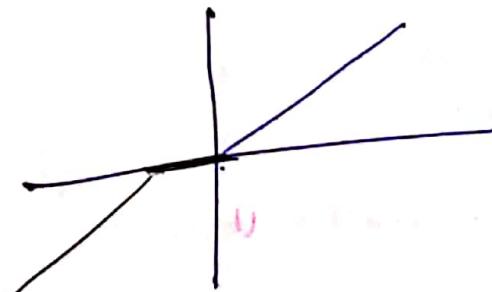
Selu



If Relu if +ve then its n  
-ve then its 0. In other words

PRelu (Probabilistic Relu)

Selu

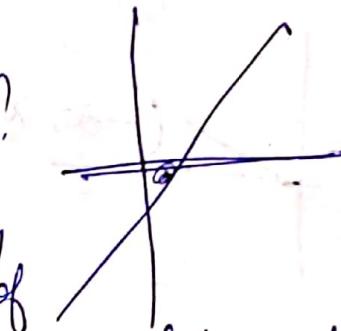


## Vanishing Gradient

Sigmoid } Squashing } When squashing happens  
Tanh } function } when we do backpropagation for value  $0.5 \text{ k } 1$ .  
There is no significance difference.  
Because it's ~~it's~~ deep squashing, so there is no significant difference when it kept squashing.

Prob. in Relu is Dead Neuron, so we use Leaky Relu.

- ① Why you can't be used?
- ② There is no non-linearity in the data. The purpose of deep learning is to include non-linearity in the data.



1. They are need to be personalized.
2. It needs to provide choice to the users.
3. They are inherently ranking system, it is important. Ex: It will present result in a ranking fashion (e.g., Google, YouTube, Netflix) [Top to bottom OR left to right]

How  
Sorting is  
done

- 1) Sorting Phase: Sorting is done based on the recent post in Facebook.
- 2) Formula Phase: It gives weightage to a topic, friendship. It then sorts by this weightage.
- 3) Model: Developing formula for every person is difficult. There comes a D. Scientist to build a Model.

The model will predict how likely a people will comment or like the post.

→ There was a problem in the model like birthday, anniversary but it does not mean that we want to see that.

3) Creating Multiple Model & combining the models.

Google Ad: One model for how likely a user will click the ad.

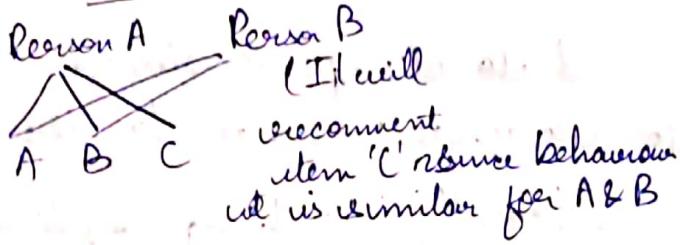
Other model for a user to stay at least for 60 seconds.

#### 4) Micro-Personalization Phase

- 
- 1) Popularity based
  - 2) Collaborative
  - 3) Content
  - 4) Hybrid

Types  
of Reco-  
mmender  
system

## Collaborative



## Content based

If user recommends item similar to his item user likes  
 Ex: If one likes Star Trek movies, will recommend Star Wars.

→ A/B testing is used to check current system with the new R.S., where current is exposed to 50% & rest new is exposed to 50% where we will see where ~~we are getting more revenue~~

## Types of Data to be used:

### 1. Explicit Data: (User has given directly)

- Customer Rating
- Feedback/reviews
- Demographics
- Ephemeral Needs (e.g. watched Chinese movies)

### 2. Implicit Data: (Infer something based on the customer activity)

- Purchase history
- Click or Browse history

Condition: Many people need to participate.

2. Must have efficient algo.
3. There must be a simple way to express.

- Time spent viewing the product
- Nos. of times the product is viewed
- Products in wish-list
- Products bought

Collaborative

Implicit Feedback

## Explicit Feedback

- o Rating given to a product {This is}
- o Comments left for a product better

Explicit data is very less. Since people don't give rating everytime. It is inaccurate.

Mostly Implicit data is used, it depends on company based on the count. of implicit or explicit data we have. We will give weightage to implicit as well as explicit.

## Strong Generalization vs Weak Generalization

→ If all test set is present in train set, it is weak generalization.  
→ Here uniqueness is user-item pair.

→ But if none of the users are present in the train set then it is a strong generalization

→ Its important to see your results on both kind of test sets.

## -ve Sampling

- 1) When we assume negatives we have to use this Ex: If I show you 10 item from my system, I clicked on item 3 & item 5. I will say there is +ve for 3, -ve for 5 & I assume that there is -ve for 1, 2, 4, 6, 7, 8, 9, 10. So here I will sample below these nos. since I don't know whether it is actually -ve or not
- 2) So all +ves taken & sample of -ves are taken

## Techniques in use: Similarity Index

### Approach:

- 1) Represent a customer as an N-dimensional vector of items.
- 2) Vector is +ve for purchased or positively rated items and negative for negatively rated items

3. Based on cosine similarity : find similar customers/users, bet. A,B or A,C or A,- or

$$\text{similarity } (\bar{A}, \bar{B}) = \cos(\bar{A}, \bar{B}) = \frac{\bar{A} \cdot \bar{B}}{|\bar{A}| \cdot |\bar{B}|}$$

4. Generates recommendations based on a few customers who were most similar to the user.

5. Rank each item according to how many similar customers purchased it. (Taking top 10 or 100 customers)

### Problems:

1) Computationally expensive:  $O(M \times N)$  is the most case  
 $M = \text{nos. of customers}$  if  $M$  is in millions  
 $N = \text{nos. of items}$   $N$  is in billions  
Not feasible

2) Dimensionality reduction can increase the performance. But, also reduce the quality of the recommendation.

3) For very large data sets, such as 10 million customers and 1 million items, the algorithm encounters severe performance and scaling issues.

So people started using Clustering Methods.

### Approach

- 1) Divide the customers into segments & treat as classification problem
- 2) Assign user to the segment containing the most similar customers.
- 3) Use the purchases and ratings of the customers in the segment to generate recommendations
- 4) This method have better online scalability and performance. Implicitly clusters happened based on usage

### Problems:

- 1) Quality of recommendation is low.
- 2) To improve quality, it needs online segmentation, which is almost as expensive as finding similar customers using collaborative filtering.

DBSCAN - fast (anomaly detection) | K-Means - fast.  
Hierarchy - Accurate and slow

### User Factorization

It assumes that:

- 1) Each user can be described as  $k$  attributes or features: E.g; Feature 1 can be user (that says how close it is to user), "2" is user liking movie, "3" is user's age etc.
- 2) Each item can be described by an analogous set of  $k$  attributes or features. Example, feature 1 would say or describe to some other feature.
- 3) If we multiply each feature of the user by the corresponding feature of the movie & add them all, this will be a good approximation for the rating of the user.

### Matrix Factorization

• User matrix  $\times$  Item matrix

• User matrix  $\times$  User matrix  $\times$  Item matrix

• User matrix  $\times$  User matrix  $\times$  User matrix  $\times$  Item matrix

• User matrix  $\times$  User matrix  $\times$  User matrix  $\times$  User matrix  $\times$  Item matrix

• User matrix  $\times$  Item matrix

• User matrix  $\times$  Item matrix

• User matrix  $\times$  Item matrix

• User matrix  $\times$  Item matrix

• User matrix  $\times$  Item matrix

• User matrix  $\times$  Item matrix

• User matrix  $\times$  Item matrix

## User Based R.Sys.

- 1) Each user is represented as a vector of items
- 2) Cosine

UCF → Not great for inactive users & gave rating for 1 item then & if not watched then the vector is empty.

IBCF → Not great for inactive items [When more inactive items up] Great for old items [Since we can recommend similar item based on current items]

t-SNE: Best for data visualization, but takes time.

## Evaluation for Recommender Systems

- RMSE
- MAP : When we have binary feedback, it is best. It takes account crowding or ordering.
- NDCG : Normalized Discounted Cumulative Gain  
When feedback is numeric, it account position  
Ex: Enum or Ordinal
- Explicit Feedback:  
MPR & MRR

MAP (Mean average precision) : If must have the results at the top.

Implicit  
Feedback



NDCG<sub>1</sub>

The relevance is discounted by  
 $\gamma_{ui} = \frac{1}{\log_2(i+1)}$  and the sum @k is normalized

by its upper bound - the IDCG<sub>1</sub>

$\forall k=3$

$$\begin{aligned} & \text{DCG}_1(\text{Discounted cumulative gain}) @ k \\ &= \frac{1}{\log_2(1+1)} + 0 + \frac{5}{\log_2(3+1)} \text{ Rating for Result 1} \\ & \quad \text{Rating for the result which is ranked 3rd} \\ & \quad \text{pos. no. } \log_2(1+1) \text{ pos. no. } \log_2(3+1) \\ & \text{Ideal (for ideal ranking) Rating of 1 Rating of 2} \\ & \text{IDCG}_1 @ k = \frac{5}{\log_2(1+1)} + \frac{3}{\log_2(2+1)} + \frac{1}{\log_2(3+1)} = 7.39 \\ & \quad \text{pos. nos } \log_2(1+1) \text{ pos. nos } \log_2(2+1) \text{ pos. nos } \log_2(3+1) \end{aligned}$$

$$\text{NDCG}_1 @ k = \frac{3.5}{7.39} = 0.47$$

If compared to position nos. 3, I am half way the ideal

Ground Truth

True Result  
1: Rating 5

True Result  
2: Rating 3

True Result  
3: Rating 1

Rank Induced  
by algo

True Result  
3

- True Result  
2

True Result  
1

True Result  
2