

Machine Learning 601

Lecture ① :

- a) get data
 - b) Problem objective
 - c) decide algo
 - d) run
 - e) validate result.
-

estimation.

generalization.

•) Supervised learning:

$$D_m = \{(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)\}$$

Input Output.

$x^i \in \mathbb{R}^d$ → dimension.

$y^i \in \{+1, -1\}$ → classification
 ↓
 positive example Negative example:

$\psi(x)$: feature representation $\in \mathbb{R}^d$

•) Regression → Supervised wity $\{y \in \mathbb{R}^d\}$

•) unsupervised learning:

Density estimation:

Given: $x^1, x^2, \dots, x^m \in \mathbb{R}^D$ } drawn from
B. $\Pr(x)$
Distribut.

Predict

then $\Pr(x^{m+1})$ of an elem drawn from
same distribution.

Setzen

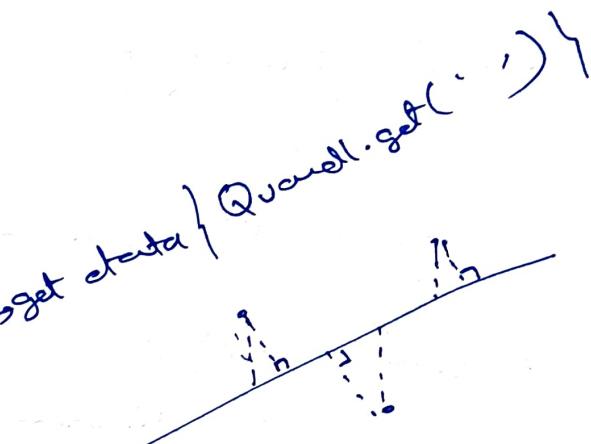
~~1)~~ ML intro

2) Regression:

pip install sklearn
pip install Quantile
pip install pandas

(feature, labels) \rightarrow supervised..

Input Output



Scaling \rightarrow Feature $\rightarrow (-1, +1)$

Cross-Validation \rightarrow shuffles the data & breaks it into training & testing.

Sklearn \rightarrow preprocessing.

Sklearn . Cross-validation

Sklearn . LinearRegression

Sklearn . SVMs

used to thread the regressor process
parallelly -

$X_{\text{train}}, X_{\text{test}}, Y_{\text{train}}, Y_{\text{test}} = \text{cross_validator}.\text{train_split}(X, Y, \text{test_size} = 0.2)$

Classifier = LinearRegression();

clf . fit(X_{train} , Y_{train}); \rightarrow training.

clf . score(X_{test} , Y_{test});

accuracy

kernel = 'polymer' / 'linear'

④ Classifier = SVM . SVR()

Support Vector Regression

L-5

Regression for future prediction

(pickle file)

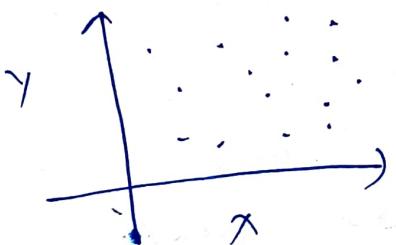
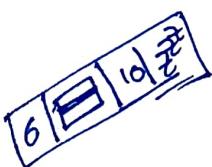
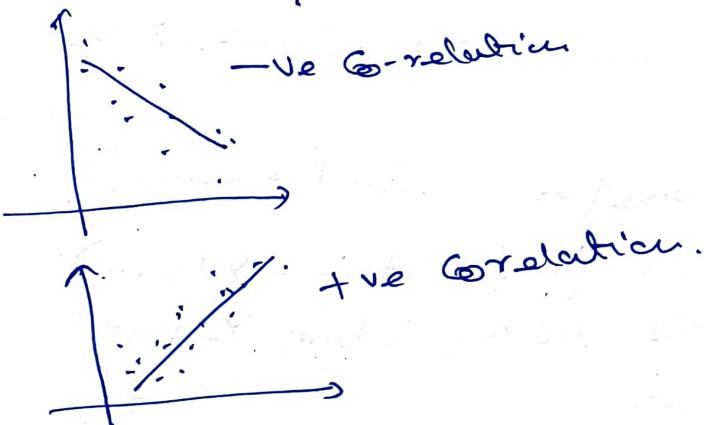
Pickling : ~~saves~~ lets us to not to

~~by saving~~ run the classifier for the whole data again and again.

if the new data is entered.

L6:

~~linear reg~~
Regression

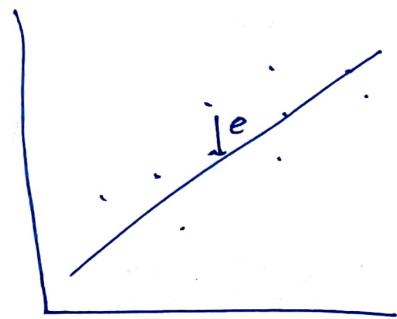
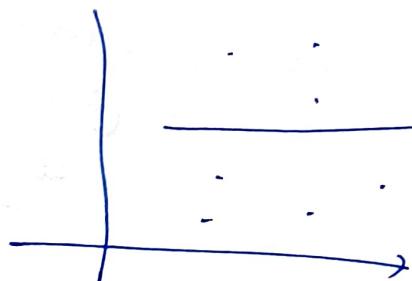


→ x is not very well correlated with y , therefore doing something like LRF will not be a good idea.

$$Y = mx + b \rightarrow \text{best fit}$$

$$\left. \begin{array}{l} m = \frac{\bar{x} \cdot \bar{y} - \bar{xy}}{(\bar{x})^2 - (\bar{x}^2)} \\ b = \bar{y} - m\bar{x} \end{array} \right\} \text{L.R.}$$

L10 : $\frac{r\text{-squared}}{r^2}$ (Coefficient of Determination)
accuracy



(r^2)

$$r^2 = 1 - \frac{SE(\hat{y})}{SE(\bar{y})}$$

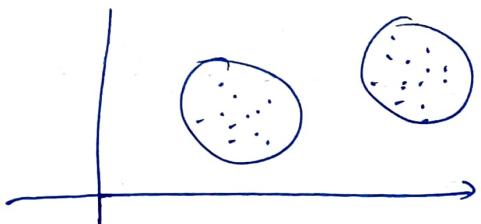
(0-1)

r^2 should be high

LII : Testing Assumptions :

Creating your own dataset

LIII⁶ : Classification :



$k = \text{odd}$
To break tie
 $k > \text{no. of groups}$

Clustering

$\leftarrow (- - +) \Rightarrow 66\% \text{ confidence}$
confidence
not accuracy

Euclid dist

Classification = k Neighbors Classifer)

sklearn.neighbors

$$\text{LIV: Euclid dist} = \sqrt{\sum_{i=1}^m (q_i - p_i)^2}$$

self neighbor.

L17

To inc. accuracy \Rightarrow we can shuffle the data

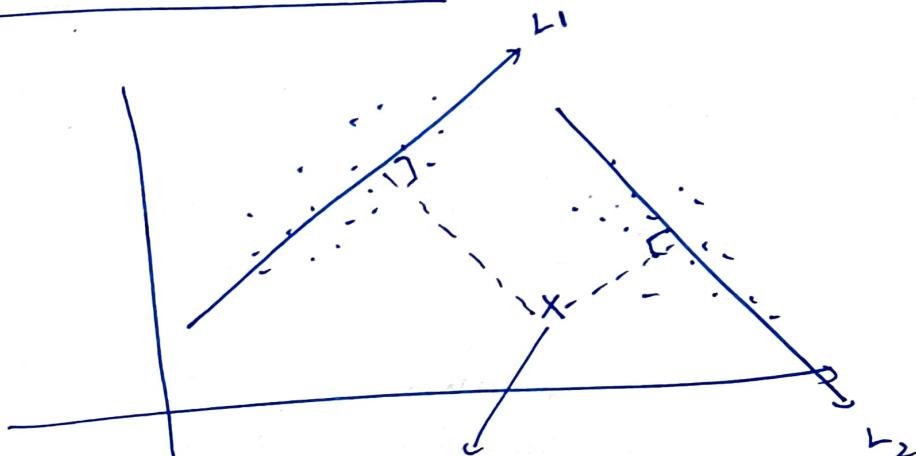
L18

k, accuracy & Predictions

There is two diffⁿ b/w Confidence & Accuracy.

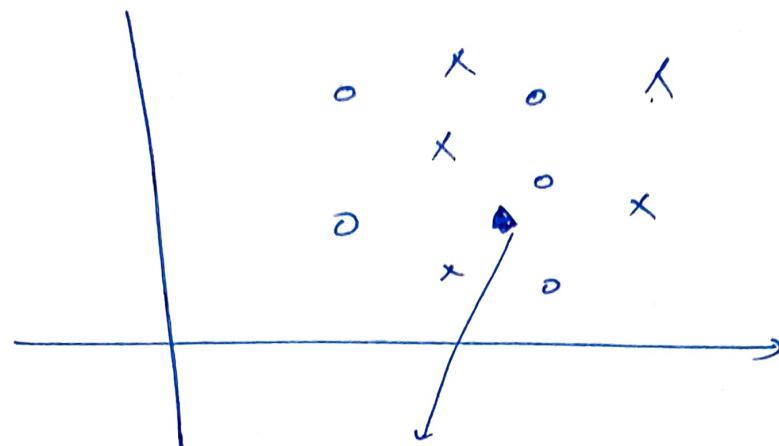
We \rightarrow can thread, to see ~~the~~ + make it fast.
 \rightarrow we can use Radius Concept.

KNN \longrightarrow Linear regression



We can measure
B/SF of the point from
the 2 regression line and
put it B there, as it less
error.

But for



to be predicted

here we can't use L_P loss

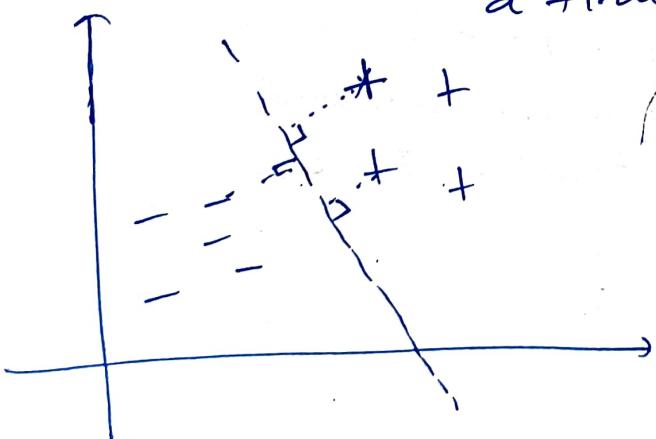
Means we can use but the accuracy will be very bad due to very bad correlation.

But we can use KNN here.

L20 SVM : Supervised Machine Learning Classifier

Support Vector Machine (Is a Binary Classifier)

means it can only separate two groups at a time.

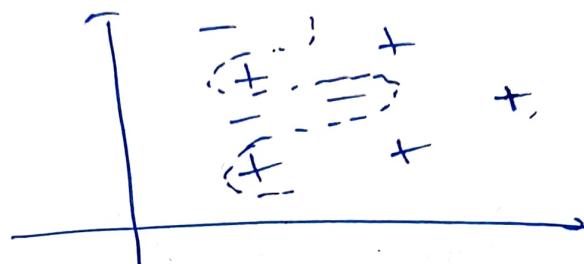


SVM finds the best hyperplane (best separating boundaries).

Maximal

Hyperplane \rightarrow has the best possible separating distance b/w the points on the two sides w.r.t. it.

As it was hyperplane \therefore we can do it on a linear data & Not some data like



\Rightarrow hyperplane not exists

SVM

\rightarrow Sum.SVC()

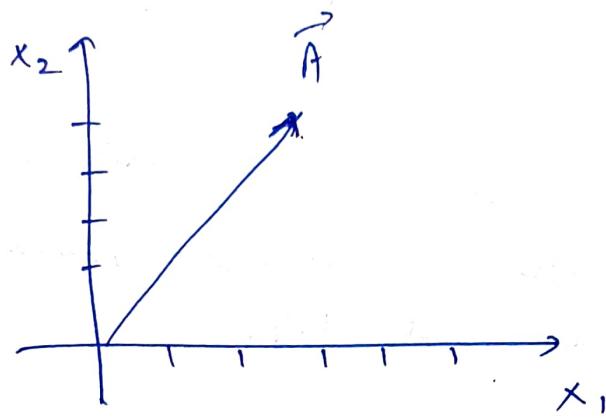
Support Vector Classifier

{ Including id Gram only reduced the accuracy of our classifier. }

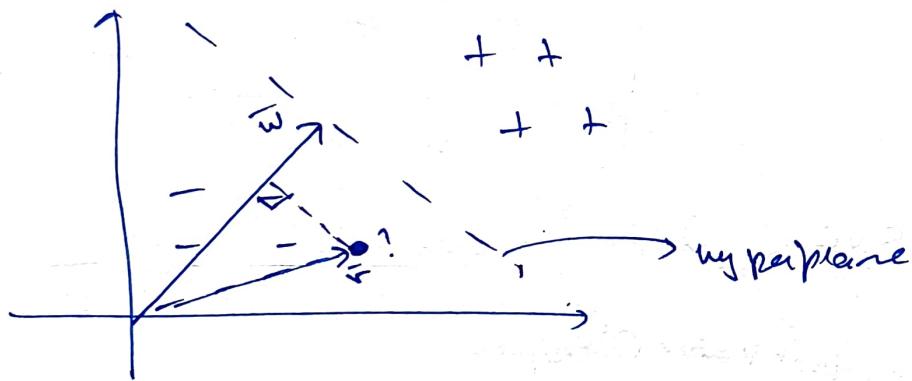
L21
SUM

$$\vec{A} = \{3, 4\}$$

$$\|\vec{A}\| \Rightarrow \text{Norm/Mag.}$$



L22:



$$\vec{v} \cdot \vec{w} + b \geq 0 \Rightarrow '+' \text{ class.}$$
$$\leq 0 \Rightarrow '-' \text{ class.}$$

i) $\vec{v} \cdot \vec{w} + b = 0 \Rightarrow \text{on the Decision boundary.}$

$$\begin{aligned} & \vec{x}_{-sv} \cdot \vec{w} + b = -1 \\ & \vec{x}_{+sv} \cdot \vec{w} + b = +1 \end{aligned} \quad \left. \begin{array}{l} \text{---} \\ \text{---} \end{array} \right\} \begin{array}{l} \text{+ve} = +1 \\ \text{-ve} = -1 \\ \text{(Class)} \end{array}$$

$\begin{array}{l} \text{+ class : } (\bar{x}_i \bar{w} + b = 1) \times y_i \\ (\bar{y}_i = 1) \end{array}$

$\begin{array}{l} \text{- class : } (\bar{x}_i \bar{w} + b = -1) \times y_i \\ (\bar{y}_i = -1) \end{array}$

$y_i(\bar{x}_i \bar{w} + b) = 1 \cdot y_i$

$y_i(\bar{x}_i \bar{w} + b) - 1 = 0$

$(\bar{x}_i \bar{w} + b) y_i = -1 \times -1$

$\Rightarrow g(\bar{x}_i \bar{w} + b) - 1 = 0$

Same Condition:

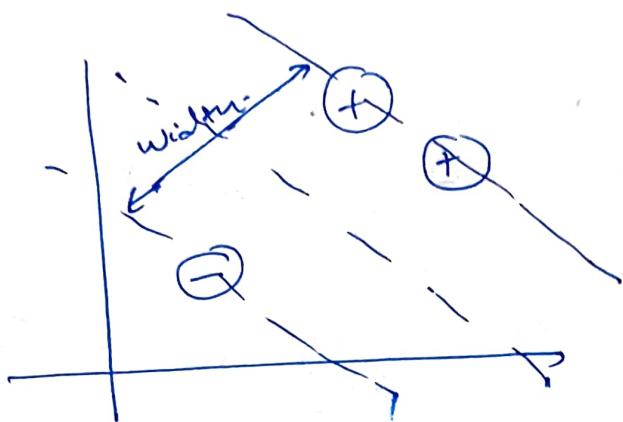
$y_i(\bar{x}_i \bar{w} + b) = 1 = 0$

y_i
 1
 -1

Moving w with b/w the support vectors.

23°

$$y_i(\bar{x}_i \bar{w} + b) - 1 = 0$$



$$\arg \max_{\text{width}} (w \cdot \text{width}) = (\bar{x}_+ - \bar{x}_-) \cdot \frac{\bar{w}}{\|\bar{w}\|}$$

subject to

$$\left\{ \begin{array}{l} y_i(\bar{x}_i \bar{w} + b) - 1 = 0 \\ \bar{x}_+ = 1 - b \\ \bar{x}_- = 1 + b \end{array} \right.$$

ansolute we
find

maximum

$$\text{width} = \frac{2}{\|\bar{w}\|}$$

we minimize $\|\bar{w}\| \Rightarrow$
 $\therefore \text{minimize} = \boxed{\frac{1}{2} \|\bar{w}\|^2}$ for
 mathematical convenience

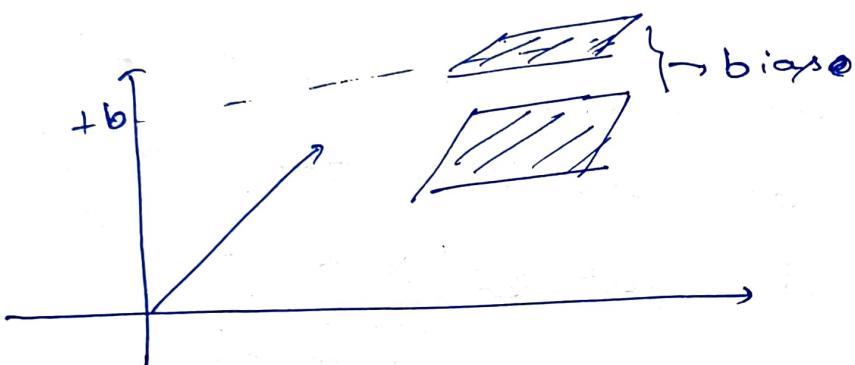
$$\left(\frac{1}{2} \|w\|^2 \right)$$

\hookrightarrow Maximize with constraint $y_i(\bar{x}_i^T \bar{w} + b) - 1 = 0$

Lagrange,

$$L(w, b) = \frac{1}{2} \|w\|^2 - \sum \alpha_i [y_i(\bar{x}_i^T \bar{w} + b) - 1]$$

we want to minimize ' w ' and maximize ' b '
bias



$$\frac{\partial L}{\partial w} = 0, \quad \frac{\partial L}{\partial b} = 0$$

$$\frac{\partial L}{\partial w} = \bar{w} - \sum \alpha_i y_i \bar{x}_i = 0$$

$$\therefore \bar{w} = \sum \alpha_i y_i \bar{x}_i \quad \text{--- (1)}$$

$$\frac{\partial L}{\partial b} = -\sum \alpha_i y_i = 0 \quad \text{--- (2)}$$

$$L = \sum_j \alpha_j - \frac{1}{2} \sum_{ij} (\alpha_i \alpha_j) y_i y_j \cdot (\bar{x}_i \cdot \bar{x}_j)$$

quadratic

downside

SUM \Rightarrow has same problem of speed or in KNN, but it can be done by taking ~~big~~ batches at a time.
Mini batches

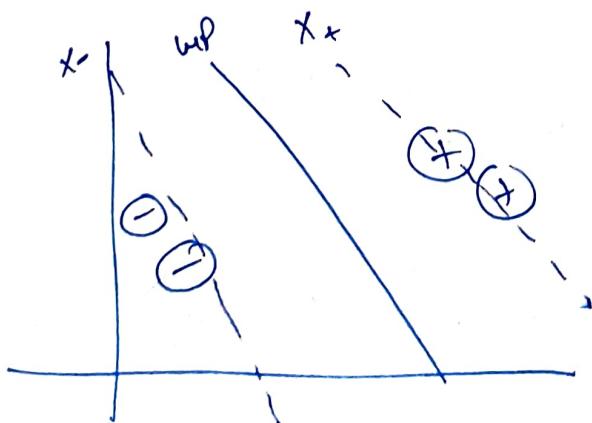
SMO \Rightarrow Sequential Minimal optimization.
some usefuls for making it fast.

up side of SUM:

Once we have trained tree data set
we don't need old feature & the classifier
for the new point is just ~~one~~ $\underbrace{(x \cdot w + b)}_{\text{Sign of}}$

L29:

$$\text{Any hyperplane} \Rightarrow (x \cdot w + b) = 1 \quad (+ \text{ class}) \text{ sv} \\ = -1 \quad (- \text{ class}) \text{ sv}$$



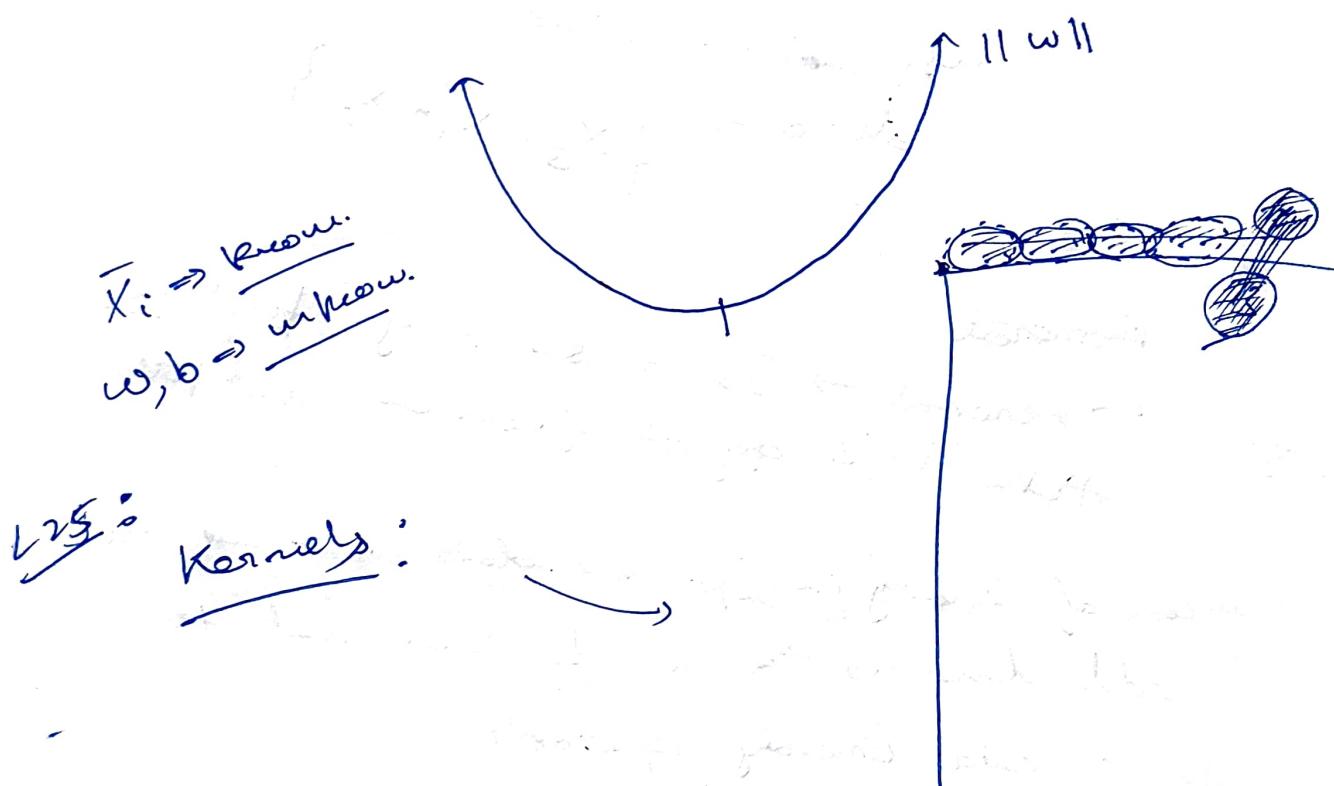
Find \bar{w} & b ?

optimization objective:

To min $\|\bar{w}\|$ & max' b'

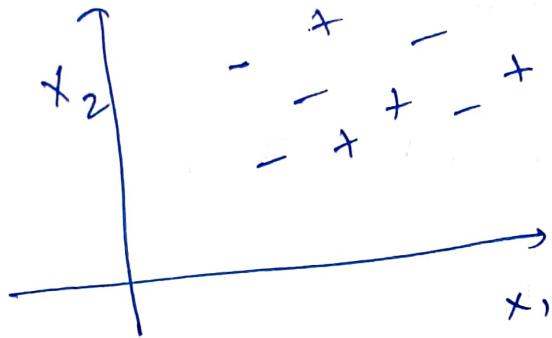
subject to $y_i(\bar{K}_i \bar{w} + b) \geq 1$

∴ a convex problem



Kernel

Non linear data



Let's add in a 3rd dimension: $\{x_3 = x_1 * x_2\}$

(use to convert) Kernel \hookrightarrow Kernel \rightarrow Is a similarity function
take input x_1, x_2 and given an output x_3 .

\Rightarrow In case of linearly inseparable data what we do is add dimensions to our feature set to make our data linearly separable.

(dot = inner) product

X space $\longrightarrow Z$ space

Unknown feature set \bar{X} ,

classification of $X \Rightarrow y = \text{sign}(\bar{w} \cdot \bar{x} + b)$

scalar

∴ we can replace \bar{X} to \bar{z} in any dimension

constraint

$$\textcircled{1} \quad \left\{ y_i (\bar{x}_i \cdot \bar{w} + b) - 1 \geq 0 \right\}$$

\bar{z}_i

\textcircled{2}

$$w = \sum \alpha_i y_i \bar{x}_i$$

$$L = \sum \alpha_i - \frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j (\bar{x}_i \cdot \bar{x}_j)$$

kernel(ϕ)

$$k(x, x') = \bar{z} \cdot \bar{z}'$$

given $\bar{z} = f_{\text{vec}}(x)$ \rightarrow same function

$$\bar{z}' = f_{\text{vec}}(x')$$

$$\{y = w \cdot \bar{x} + b\} \Rightarrow \text{using Kernel.}$$

feature set $X = [x_1, x_2]$ Correct \Rightarrow 2nd order poly.

$$Z = \underbrace{\{1, x_1, x_2, x_1^2, x_2^2, x_1 x_2\}}_{Z \text{-space}}$$

\nwarrow

2nd order

$$K(x, x') = Z \cdot Z'$$

~~(Don't do)~~

$$Z' = \underbrace{\{1, x_1, x_2, x_1^2, x_2^2, x_1 x_2\}}_{Z \text{-space}}$$

$$K(x, x') = Z \cdot Z'$$

$$\hookrightarrow 1 + x_1 x_1 + x_2 x_2 + x_1 x_1^2 + x_2 x_2^2 + x_1 x_2 x_1 x_2$$

we can actually use kernel without using the Z space
above:

$$K(x, x') = (1 + x \cdot x')^p \stackrel{\text{degree of poly}}{=} \text{polynomial kernel}$$

$x, x' \in \text{has dimension } m$

$$K(x, x') = (1 + x_1 x_1 + x_2 x_2 + \dots + x_m x_m)^p$$

\Rightarrow if we choose $m = 15$ & $p = 100$
how can we do it?

RBF \Rightarrow Radial basis Kernel.

Default
kernel

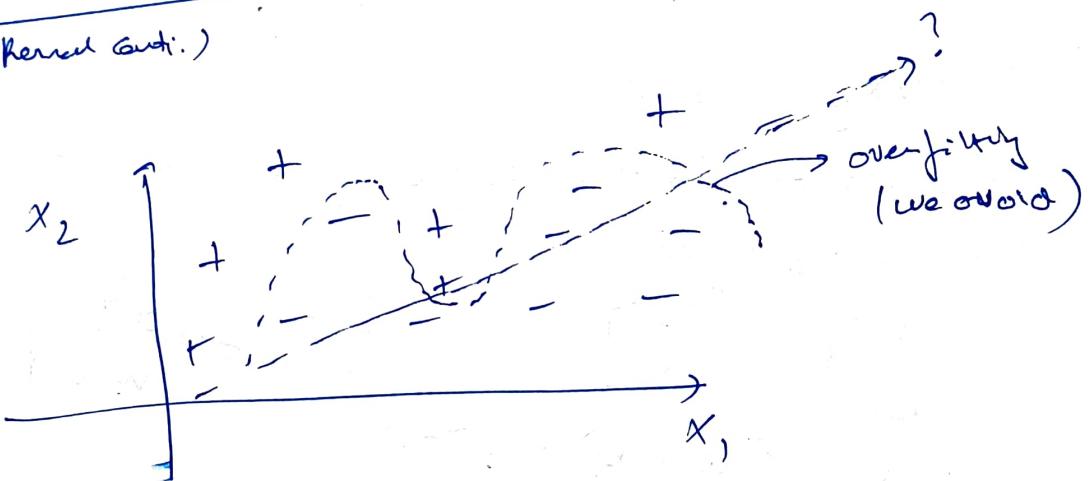
$$k(x, x') = \exp(-\gamma)$$

$$k(x, x') = \exp(-\gamma \|x - x'\|^2)$$

Can take up to ∞ dimension

L31 Soft Margins in SVM

(Kernel cont.)



$$\frac{\# \text{SV}}{\# \text{Samples}} > 10\% \quad \text{overfitting}$$

if we has

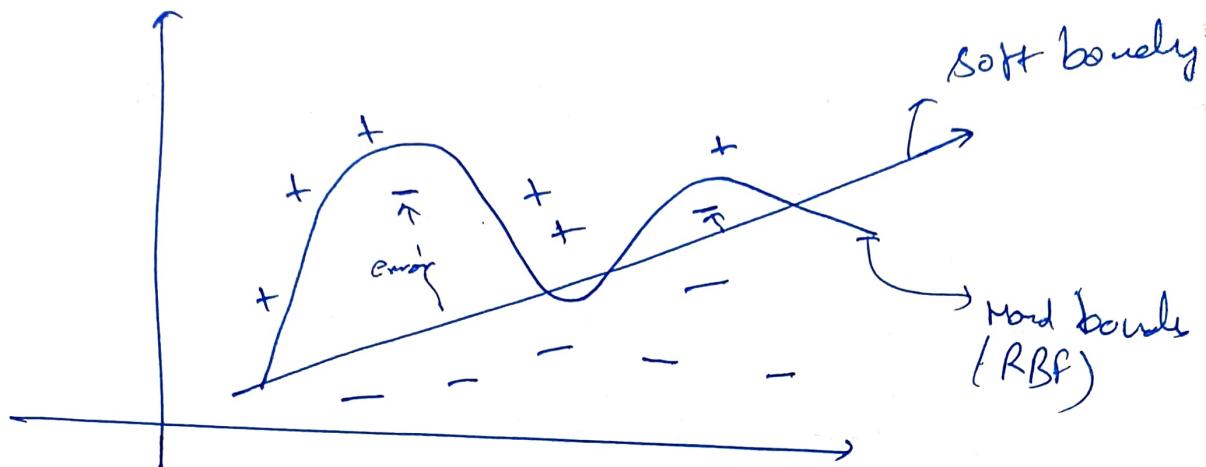
52% accuracy & 80% SV's

overfitting.

If we have [52% accuracy & 8% SV's]

then the data is probably not good. work out & \therefore try a different kernel.

Soft Margin SVM



Soft Margin Classification → allow some degree of error in the hyperplane. we allow do this using slack.

$$\boxed{\text{Slack} = \epsilon_i}$$

$$\left\{ \begin{array}{l} \text{for } \epsilon_i > 0 \\ \text{hard margin} \end{array} \right.$$

$$y_i(\bar{x}_i \bar{w} + b) - 1 \geq 0$$

$$y_i(\bar{x}_i \bar{w} + b) \geq 1 - \epsilon_i$$

$$\text{Total Slack} = \sum_{i=1}^n \epsilon_i$$

Now we want to minimize the slack,

$$\therefore \text{Now we will minimize} = \min \left(\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \epsilon_i \right)$$

Controls the Violation Margin.
i.e if $C \uparrow \Rightarrow$ we allow less violations.
if $C \downarrow \Rightarrow$ we allow Violations.

L32

"Cvx opt": Pattern Recognition.

$$, \text{mp. dist}(x, x_2)$$

→ linear Kernel ν_2

→ polynomial Kernel ν_3
($\nu_3 = \text{mp. dist}(x, y)$)

→ Gaussian Kernel

→ Radial Basis Function

$$\text{mp. exp}\left(-\frac{\|x - y\|^2}{2 * (\sigma^2)}\right)$$

~~sk~~

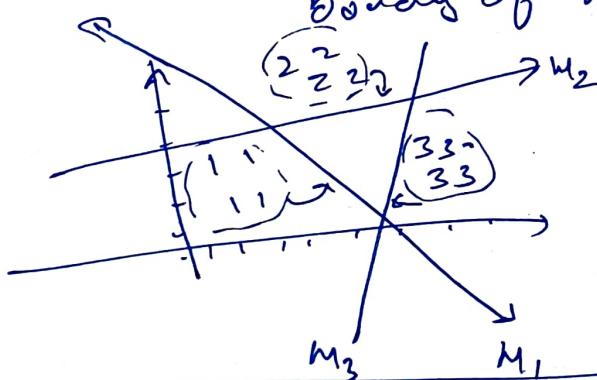
L33: Classifying more than 2 groups using SVM:

OVR → one vs rest

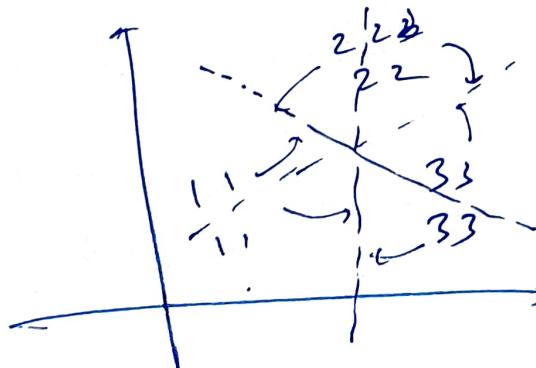
OVS → one vs one

OVF

problem: → have sort of weighting issue i.e.
(Default) ~~imbalance b/w the two~~
of data points on two sides of the
boundary of my hyperplane.



OVO



SLearn

SVM, SVC



kernel
 → 'Sigmoid'
 → 'precomputed'

"Probability" → just like "Confidence" in the K-NN

"Shrinking" & uses SMO if True, to ignore the non important features while fitting the data.

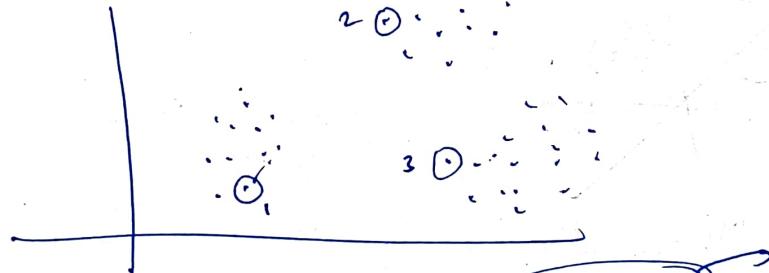
"Verbose"?

L34 Unsupervised Learning

Clustering → Flat → we give no. of clusters

→ Hierarchical → Machine figures out (meanshift) the number of clusters.

K-Means:



Find dist of all points from ①, ②, ③

random
Chosen
Centroids.

& assign to it the cluster closest to it.

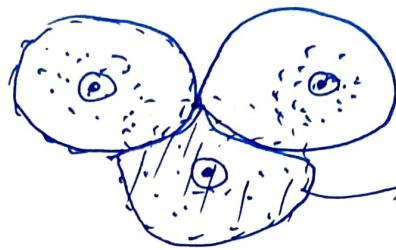
Now find ^{new} centroid of the each cluster and repeat the process until the new centroid becomes equal to the old one.

We will have confidence

We will have \rightarrow "tolerance" & "re-iter"
parameter in this.

Downside of K-mean

In like,
house
data



This must not
be like this.
This problem can
be resolved by
use of kernels.
due to Euclidean
distn}

Up Side

\rightarrow we don't have to train again & again
just like SVM.

{ Convert non-Numerical data column \rightarrow Numerical data column }

Titanic dataset : for $k=2$ (Survived / non-Survived)

Can also use SVM,

* # Scaling the features make a great changes in accuracy. (like in titanic data).

L-37

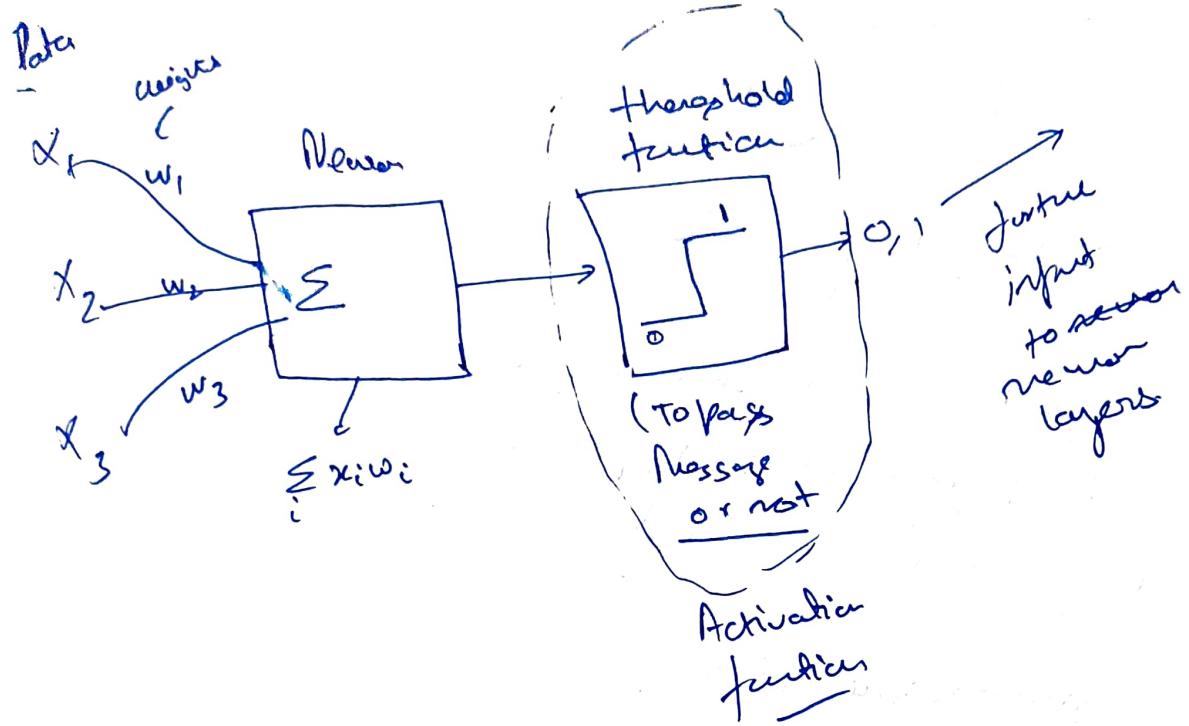
Predictive: in K-Mean we predict. here new data by seeing the distⁿ from the ~~set~~ previous centroids.

L-38 : Hierarchical Clustering:

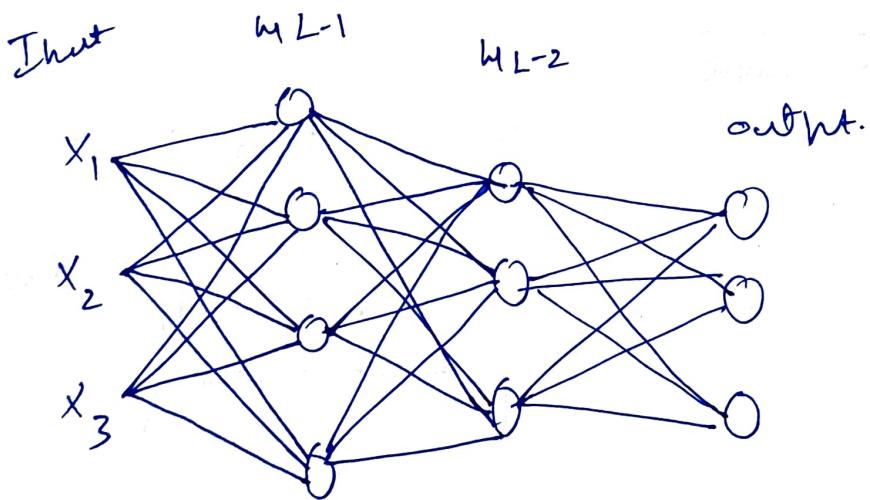
- Choose every data point one by one
- ~~Do~~ find the centroid of all the points lying in the Radius/bandwidth of 'r' with the chosen point as center.
(takes default value if not passed.)
- after finding the new centroid, do the same until the centroid converges.
- Do this for all the data points ~~BAH~~,
Bam!! we did hierarchical clustering.

Automating Radius/bandwidth values:

L 43: Neural Networks



$$y = f(\bar{x} \bar{w})$$



1 $ML \rightarrow$ Regular N.N

2 $> 1 ML \rightarrow$ Deep N.N.

Neural
Neural Networks works for doing really
good Modelling analysis.
(Modelling logics).

N.N → very impressive.

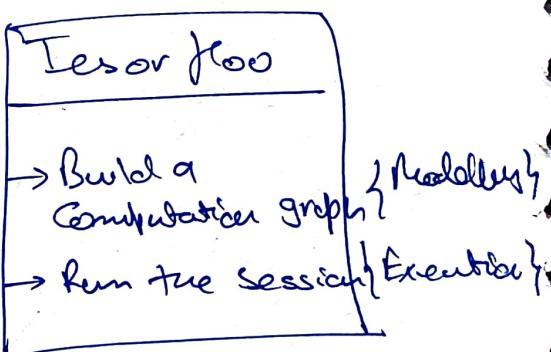
Image Data → Image net
Text Data → Wikipedia,
Speech Data → Tatoba.

Common crawl data

Luy:

Tensorflow (Framework) → Interactive
works in chunks.
(Deep learning lib), like numpy.

tf.Session()



Ans:

MNIST: 60,000 ~~training~~ training ~~samples~~ of handwritten digits. (28×28) pixels & 10,000 testing samples.

Feed forward Neural Network \rightarrow Pass the data right ~~left~~ through the last hidden layer to the output layer

\rightarrow Data in the output layer is compared to the intended output then we create a ~~get~~ loss function. (Cross Entropy).

\rightarrow Then we use optimizer to minimize the loss. Adam Optimizer, Stochastic Gradient Descent (SGD) AdaGrad etc.

\rightarrow what the optimizer does that it goes backwards ~~backwards~~ and manipulates weights (Backpropagation)

\rightarrow feedforward + backprop = epoch {epoch}

$$\text{Input data} \xrightarrow{\text{w}} \bar{x}\bar{w} \xrightarrow{\text{bias}} (\bar{x}\bar{w} + b)$$

0
0
0

Why bias?

~~If~~ $(\bar{x}\bar{w}) \rightarrow$ if all of the input data is zero
 i.e. $\bar{x} = \bar{0}$ then $(\bar{x}\bar{w} = 0)$ \therefore Input will
 be zero \therefore we add a bias 'b'. No neuron
 will ever fire/activate. \therefore we add bias 'b' to this
 to making neuron active.
 Some.

for every neuron
 Node.

of nodes in the output layer = No. of classes.

relu \Rightarrow Rectified Linear Activation Function

fit: def fit NN(x):

prediction = n_mean_model(x)

cost = tf.reduce_mean(tf.nn.softmax_

cross_entropy - with (logit prediction, y));
 optimizer

Generally,

1 hidden layer \rightarrow linear data

2 hidden layers \rightarrow Nonlinear data

{ or
!

(Buffers \rightarrow for large data)

LST: fNN (feed forward Neural Network)

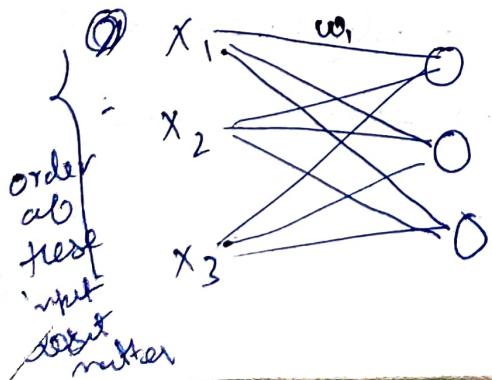
Traditionally f.f.NN (feed forward neural network)
do not take account in the ~~order~~ ^{percentage of} time or the order in which things are happening.
Ex: (A ball moving away/towards)

RNN: if ~~order~~ ^{is} used for language data.

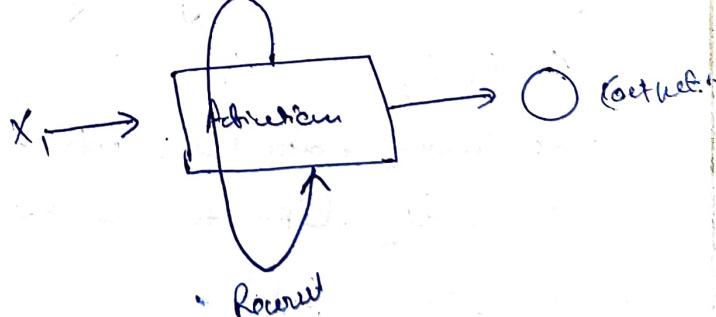
LSTM cell

Gated Recurrent Cell

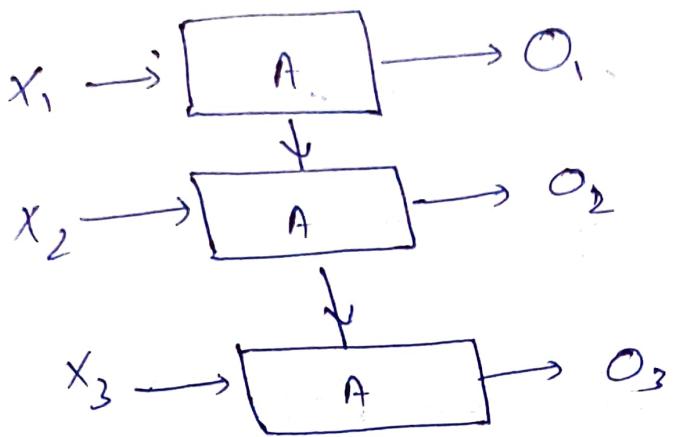
Deep NN



Recurrent NN



RNN can also be visualised like this -



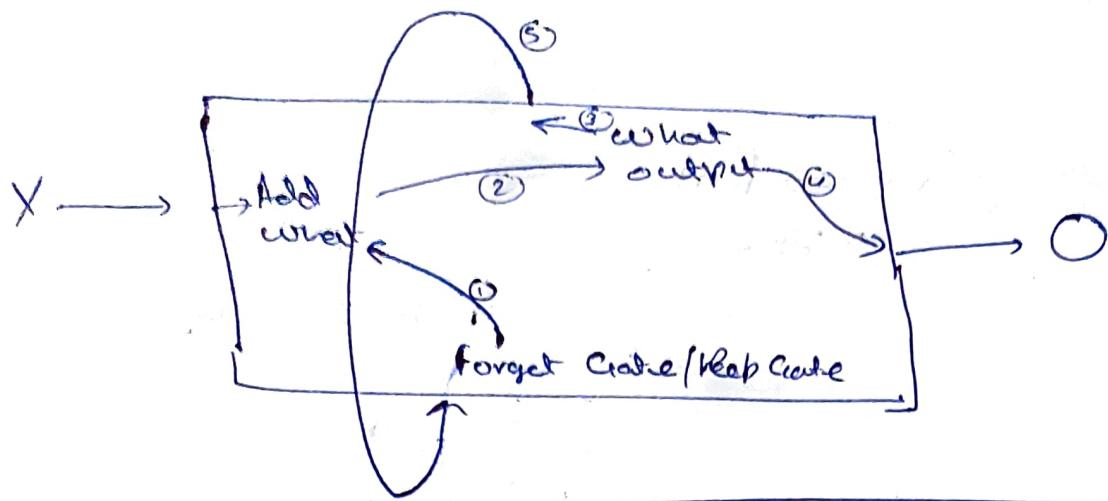
So how it is different from Deep NN,
so, if we have a sentence:

"Harrison drove the car"
for Deep Neural Network the order will not matter
& it will be same as -
"The car drove Harrison".

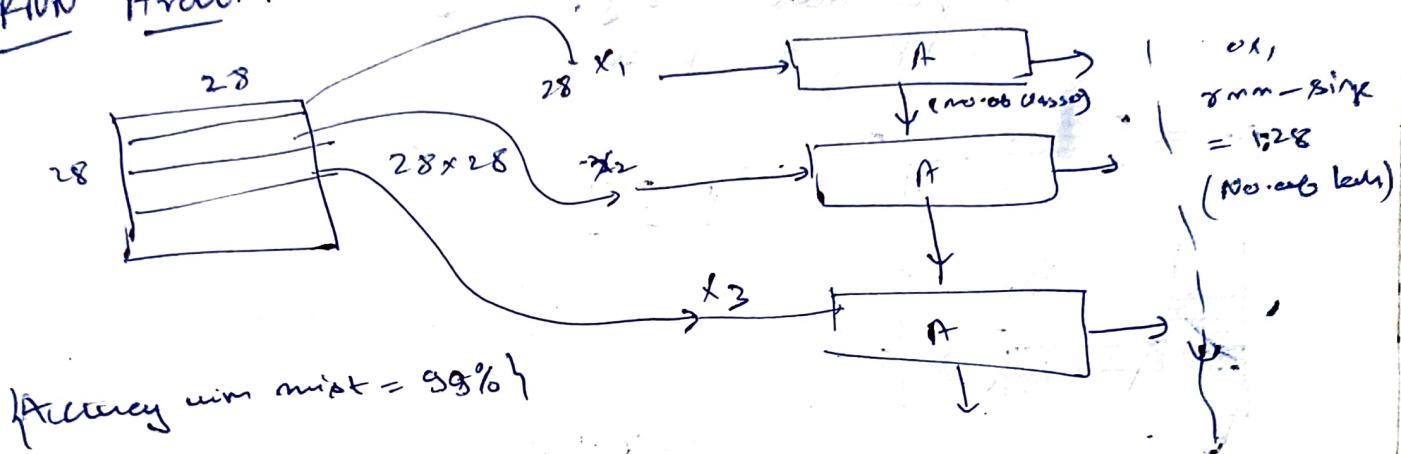
But in RNN the order/time of occurrence really matters.

LSTM cells: It is an RNN cell, ~~which do not~~
If we have a very long sequence of input it like
 $x_1, x_2, x_3, \dots, x_m$, then this series we might
have many data terms which closest really does not matter
& we only like to memorise few terms :- we
have LSTM Cell

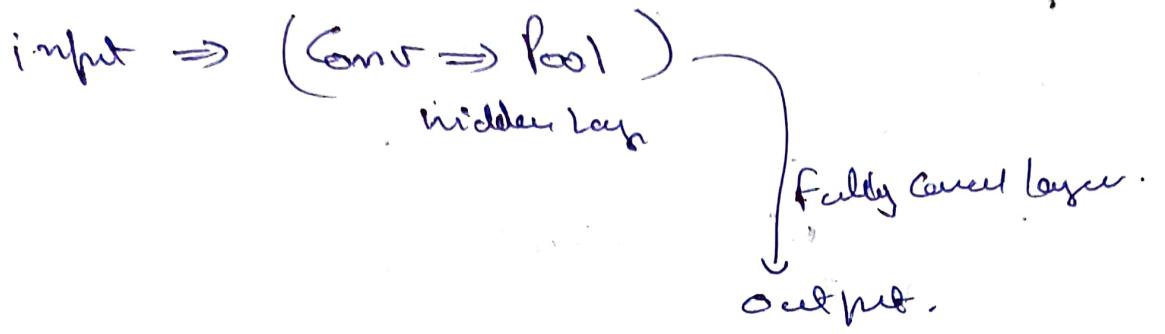
LSTM Cell



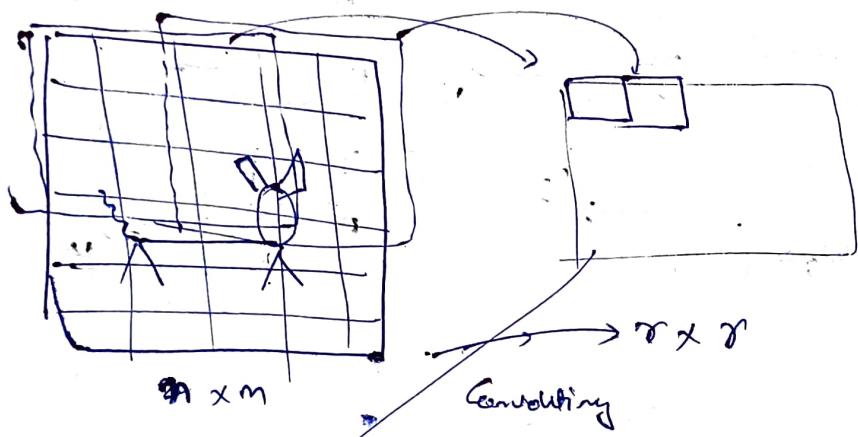
RNN Arch.:



L 56 CNN : ConvNet



Convolution \rightarrow Generating feature map from the input data. using sliding window, without skipping any pixel.



Pooling :

(No overlap)

Suppose the new feature map is

1	2	2	1	1
1	0	3	1	
1	4	1	2	2
3	2	2	0	2

Pooling

1	3
1	3

Pooling is Simple just

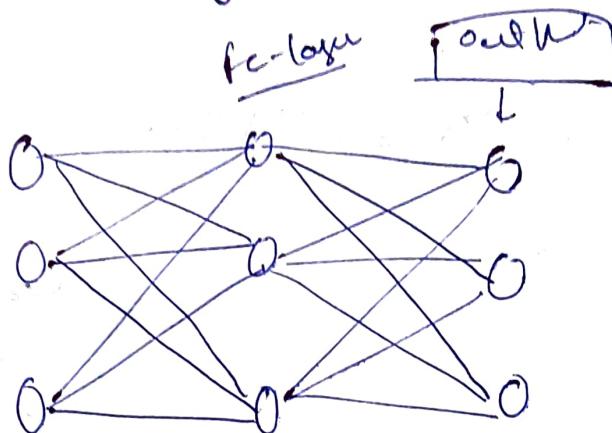
\rightarrow let's take a 3×3 pool

Max: Max pool \Rightarrow take max value in the window

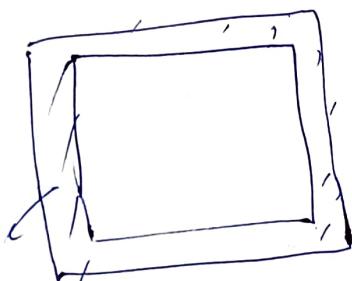
(Conv + Pool) = hidden layer

Fully Connected layer:

- ⇒ Every neuron in the previous layer is connected to every neuron in the next layer.



{
 {
 S, S } (1), (32)
 S + S
 convolution
 take 1 input
 → produce 32 outputs.



For Small Dataset RNN works
 better than CNN.
 CNN is good for large datasets.

Dropout in CNN: de-activates some unnecessary

unnecessary neurons working in the
 fc. layer.

L58° TfLearn - \Rightarrow

high level abstraction layer library on top of ~~Tensorflow~~

Tensorflow :)

" Tflearn, Keras, SKLearn, TfSlim "

↳ less error prone.

Read:
Generative Adversarial
Neural Networks

AlexNet \Rightarrow for image data.

Overlapping pooling

Was used in AlexNet
as a reduction

0.5% error was obtained

Batch Normalization: Technique used before giving input to a convolution layer to mitigate the effect of unstable gradients within deep neural network.

$$\text{input } X = \{x_1, \dots, x_m\}$$

parameters: w, b {weights & biases}

$$y_i = \text{BN}_{w,b}(x_i)$$

Normalize

$$\therefore \mu = \frac{1}{m} \sum_{i=1}^m x_i \quad \left\{ m = \text{minibatch size} \right\}$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu)^2 \quad \left\{ \begin{array}{l} \text{minibatch} \\ \text{variance} \end{array} \right\}$$

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad // \text{Normalize.}$$

$$\left\{ y_i = w \hat{x}_i + b \right\} \Rightarrow \text{Scaling \& Shifting.}$$

Optimizer: used to change attributes of NN such as weights & learning rate in order to reduce the loss

Gradient Descent: 1st order optimization.

- Play & trap local minima.
- Fast computation
- High memory

$$\theta = \theta - \alpha \cdot \nabla J(\theta)$$

Weights are changed after calculating gradient over the whole data set. So takes very long time.

Stochastic Gradient Descent:

Requirement: Randomized training dataset.

→ as the weights are updated after each training set, the error can be very much ~~as~~ noise & it's future can't help justify.

→ Mini-Batch Gradient Descent: ~~optimal~~ after every batch.

→ Momentum:
 Momentum: $\gamma \approx 0.9$ → hyperparameter
 → Convergence Optimizer

→ Adam: β_1, β_2 → hyperparameter
 (see)

→ Adagrad: η_t for each iteration and at every time step t
 → second order optimization

Model variants $\rightarrow (w, b) \rightarrow$ in ANN

(Support Vector) \rightarrow SVM

(Coefficient) \rightarrow Linear Regress & Logistic Reg

Model hyperparameters: Can be tuned for a predictive modelling problem.

→ Inputs to the function which we can tune.

Ex → learning rate in neural net.

→ C & γ in SVM.

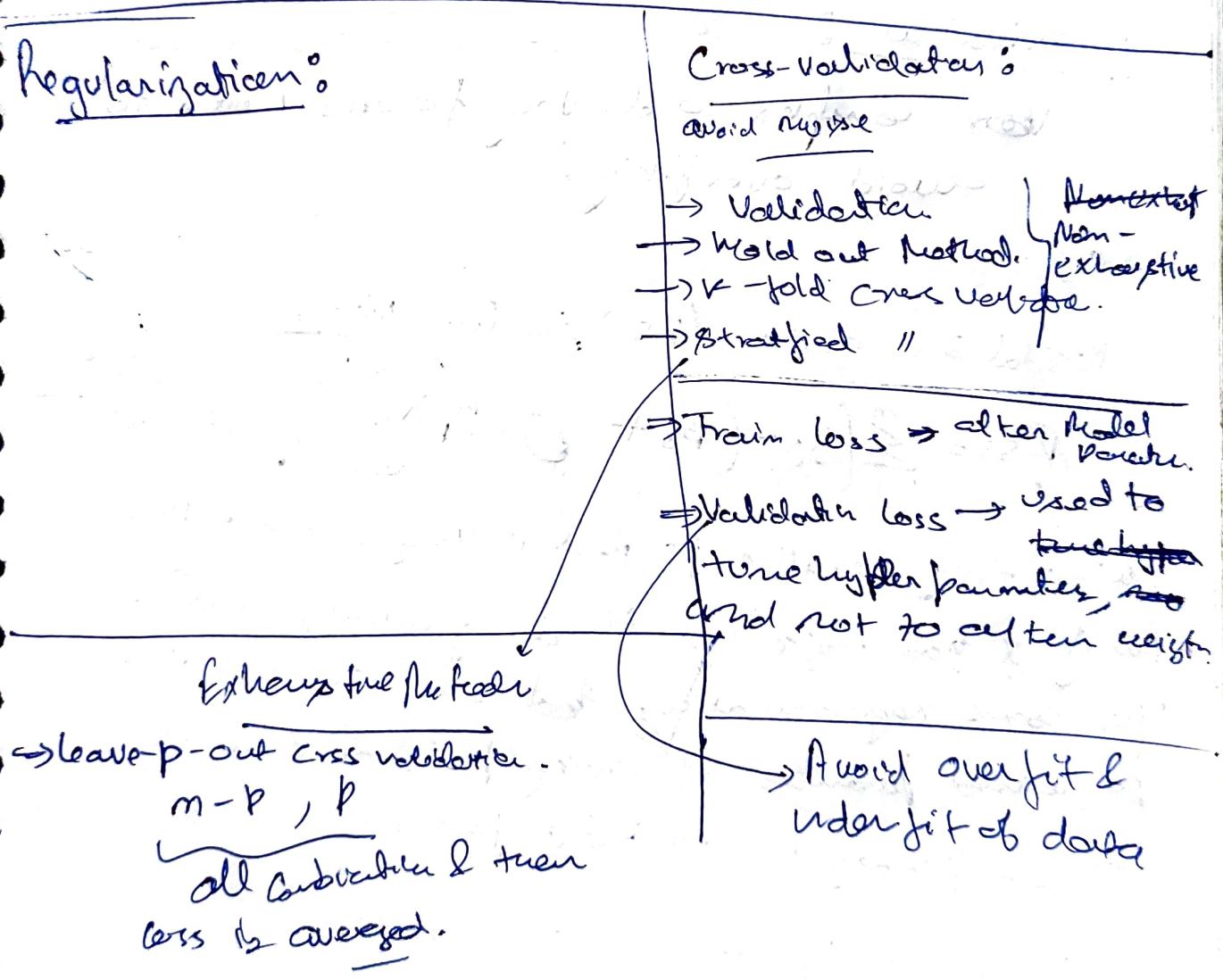
→ R in kNN

Random Search vs grid Search for hyperparameter tuning.

Grid Search:

- Pros
- Cons

Regularization:



Hyperparameter Tuning:

① Randomised Grid Search Cross-Validation

- Create a dict of ~~hyper~~ all the hyper-parameters & to all their possible/likely values as the keys.
- Our Model choose any one combination of the above hyperparameter grid and evaluate the Model.

Regularization: Discourages ~~heavy~~ from a very ~~very~~ complex and very flexible model to avoid overfitting.

L1:

Residual Sum of Squares (RSS) in linear regression.

$$RSS = \sum_{i=1}^m \left(y_i - \left(\beta_0 + \sum_{j=1}^n \beta_j x_{ij} \right) \right)^2$$

if x has only one feature

while training we will add to β_0 & β_j using Gradient Descent.

But if ' x ' has noise then it can't generalize well.
Hence comes regularization, ~~L1~~ and ~~L2~~. These ~~β~~ learned parameters at noisy points to zero by giving a penalty in the cost function.

Therefore we update the cost function.

① Ridge Regression (L2)

$$\text{RSS}' = \text{RSS} + \lambda \sum_{i=0}^m (\beta_i)^2$$

↑ no. of Parameters
↑ regularization term

Penalises for high value of

Coefficients \therefore reduce the Model Complexity & Multi-Collinearity.

Underfitting \rightarrow

- If we have very few features in data set & both training & test loss is high than its Underfitting case.

overfitting / over-generalization:

- Large no. of features
- Low training loss & high test loss.

② Lasso Regression (L1 Regularization)

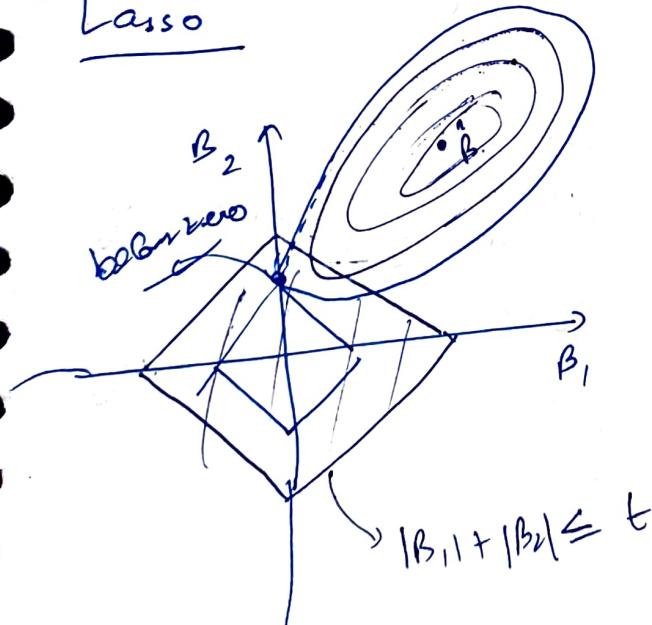
Lasso (least Abs. Shrinkage and Selection operator)

$$\text{RSS}' = \text{RSS} + \lambda \sum_{i=0}^m |\beta_i|$$

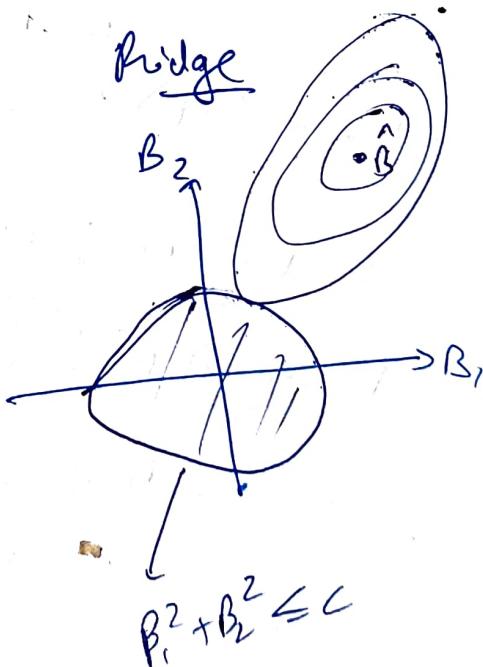
~~Ticks~~ \rightarrow Lasso can lead to zeroing some of the coefficients i.e. totally neglect them.

\therefore Lasso regression not only helps in over-fitting issue but also feature selection

Lasso



Ridge



*Some
Jobs of Model*

Performance Metrics:

Accuracy → can be decisive when dealing with unbalanced data.

Confusion Matrix, Precision, recall, F1 score. gives better intuition of prediction as compared to accuracy.

Confusion Matrix:

Matrix of size 2×2 for binary classification
with actual data on one axis & predicted on other.

In a binary classification if classes are 0, 1

		Actual	
		0	1
Predicted	0	True Neg	False Pos
	1	False Neg	True Neg
		False +ve	True +ve

True Positive (TP): Model correctly predict +ve class

True Negative (TN): Model correctly predict -ve class

False Positive (FP): Model gives wrong prediction of a -ve class (Type I Error)

False Negative (FN): Model gives wrong prediction of a +ve class.

Type II error

With the help of above terms we can calculate True Positive rates (TPR), False Negative Rates (FNR), True Negative rate (TNR) and False Positive rates (FPR).

$$TPR = \frac{TP}{Actual\ Positive} = \frac{TP}{TP + FN}$$

$$FNR = \frac{FN}{Actual\ Positive} = \frac{FN}{TP + FN}$$

$$TNR = \frac{TN}{Actual\ Neg.} = \frac{TN}{TN + FP}$$

$$FPR = \frac{FP}{Actual\ Neg.} = \frac{FP}{TN + FP}$$

→ Even if our data is imbalanced we can figure out that our Model is working ~~as~~ well or not. for that, value of TPR, TNR should be high & FNR, FPR should be as low as possible.

a) Precision & Recall:

Both precision & recall are important for information retrieval, where positive class matter the most as compared to negative class.

→ while searching on web the model does not care about something irrelevant and not retrieved (this is the TN case). Therefore only TP, FP, FN are used in precision & recall.

Precision: Out of all the "predicted", what % are truly +ve

$$\therefore \left\{ \text{Precision} = \frac{TP}{TP + FP} \right\}$$

Recall: Out of the total Positive, what % are predicted Positive.

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Same
as TPR

Example: Credit Fraud detection
(when recall is important)

		Actual	
		Fair	Fraud
Predicted	Fair	TN	FN
	Fraud	FP	TP
	Total		

In case of transaction,
→ we do not want to miss
any fraud transaction;

∴ we want FN to be as low as possible.
∴ recall should be high, but we can't
compromise with low precision.

Example: Spam detection
(when precision is important)

		Not spam	Spam
!	spam	TN	FN
	not	FP	TP

In Spam if we get any spam as non-spam it's OK, but if we miss any important. ~~very~~ due to spam it's not at all good.
∴ precision should be as low as possible.
→ In such case, FP should be as low as possible.
∴ precision is more vital than recall.

⇒ F1-Score: when comparing two models it is
used to compare which is more precise and
(low recall or vice versa). ~~Forgetting~~

It is the harmonic mean of precision & recall.
thereby taking both FP, FN into account.
∴ performs well on imbalanced data.

$$F_1 \text{ Score} = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} = \frac{2 \times (\text{Precision} \times \text{Recall})}{\text{Precision} + \text{Recall}}$$

∴ gives same weightage to precision & recall.

Weighted F1-Score: Here we can give different
weights to precision & recall. As discussed

$$\left\{ F_{\beta} = (1+\beta^2) \times \frac{\text{Precision} \times \text{Recall}}{(\beta^2 \times \text{Precision}) + \text{Recall}} \right.$$

→ weight can be given to precision & recall or
per the demand of problem.

β = represent how many times recall is more important
than precision.

Support %: No. of samples of True Positives that
lie in that class. The no. of actual
occurrences of the class in the specified dataset.

$\beta \rightarrow 1$, Precision & Recall have equal weight

$\beta = 0.5 \rightarrow \text{Precision} \downarrow$

$$\beta \in (0 - 10)$$

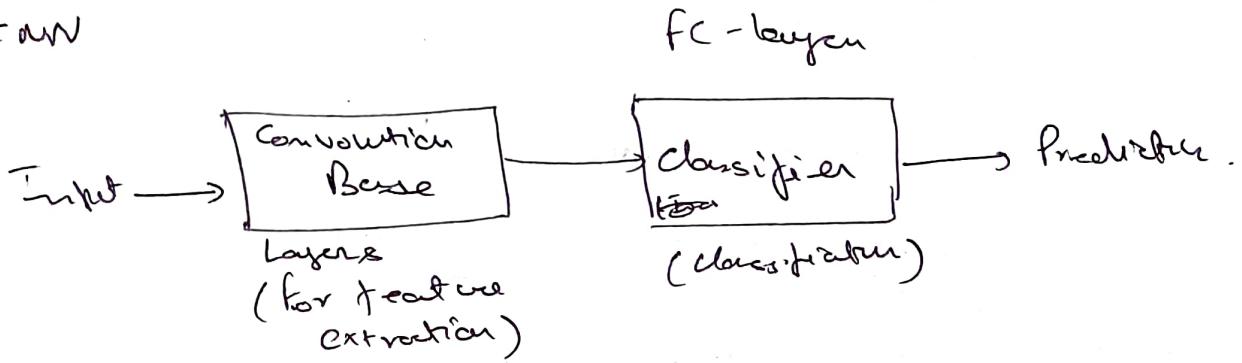
$\beta > 1 \rightarrow \text{Recall} \uparrow$

Transfer learning from pre-trained Models

① Transfer learning → To build accurate Models in a time-saving way.

pre-trained Models: VGG, Inception, Mobile Net.

② CNN



③ Repurposing a pre-trained Model:

When we are repurposing a pre-trained model for our need, we start by removing "the original classifier". Then we add a new classifier that fits our purpose, and finally we fine tune our model according to one of the 3 strategies.

④ Train the entire Model: Train from scratch.

⑤ Train some layers and leave others frozen:

As we know lower layer → Problem independent ∴ they are left frozen, while the higher layers are pre-trained also in this case too ∵ we keep the learning rate \downarrow low, because for high learning rates \uparrow may lead to risk of losing previous knowledge.

→ Now if we have small dataset $\frac{& \text{large no. of parameters}}{\text{freeze more}}$, we freeze more layers to avoid overfitting.

→ If we have large dataset, and small no. of model params we can train more layers as overfitting is not an issue

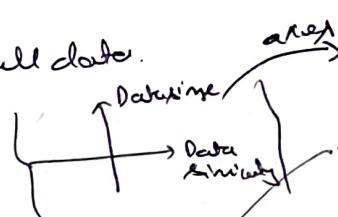
•) Freeze the Convolution base: Train/freeze trade off.
Trained Model to be used as a fixed feature extraction module.

④ Transfer learning process:

(a) Select a pre-trained Model \Rightarrow ~~be very have~~ \rightarrow VGG, InceptionV3, ResNet5.

(b) Classify your problem according to the size-similarity Matrix.
Rough outline:

if $n < 1000 \rightarrow$ small data per class.



Quadrant 1	Quadrant 2
Large data, but diff from Pre-trained data	Large data, similar to Pre-trained
Small data, different from Pre-trained	Small data, similar to Pre-trained.
Quadrant 3	Quadrant 4

⑤ Fine tune the Model:

Here we can use the size similarity Matrix to guide our choice and then refer to one of the 3 options we saw in preparing a pre-trained Model.

\rightarrow Quadrant 1:

Grayscale

Quadrant

Q1: Train whole, but we can initialize the parameters by using a pre-trained Model very its architecture & weight.

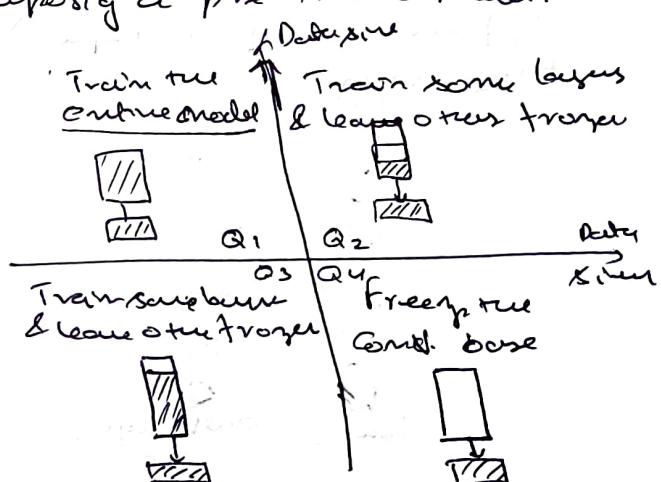
Q2: here Bingo, we can start computational (training) by freezing more no. of layers

Q3: "2-7 off & it's hard of Computer with" problem. In this case we can't do anything as if we "train" more layers then there will be overfitting & if we "train" less layers then our model won't learn anything useful.

(we can use data augmentation techniques in this case)

After doing all this we must go deeper than Quadrant 2, i.e. train more than in Quadrant 2.

Q4: use strategy 3: freeze the Convolution base.



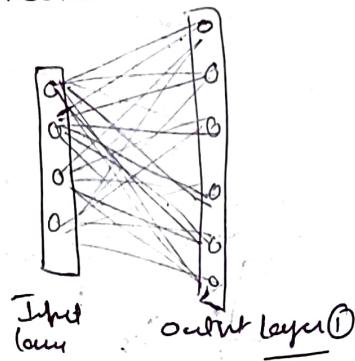
Data Augmentation :

Rocket Engine → Deep learning Models
 Fuel → Amount of Data.

Ex: of Recent Models :

	VGG Net	DeepVideo	CNNMT
Used for	Image classification	Video classification	Translation
Input	Image	Video	English text
Output	1000 classes	47 classes	French text
Params	140M	100M	380M
Datapoint	1.8M image with class	1.1M audio & class	6M sentences Sentence Pair 340M words
Dataset	ILSVRC-2012	Sports-1M	WMT'14

$$\left\{ \begin{array}{l} \text{No. of parameters} = [\text{single input} + 1] * [\text{single output} + 1] \\ \text{to train} \end{array} \right.$$



A CNN, that can robustly classify objects even if they are placed in different orientation is said to have the property called Invariance. More specifically a CNN can be invariant to translation, viewpoint, size or illumination.

When to augment in our ML Pipeline?

- ① Before feeding data into model. {offline augmentation}
or
② Augment the mini-batches, just before feeding it to the machine learning Model. {online augmentation / Augment on fly}

Offline Augmentation

→ for small datasets

Online Augmentation

→ large datasets.

Type

⇒ flip (Horizontally, Vertically)

⇒ Rotation

⇒ Scale

⇒ Crop

⇒ Translation. {Moving image in x, y direction}

⇒ Gaussian Noise {over-fitting usually happens when Net tries to learn high frequency feature (HD)}
· {Salt Noise, Pepper Noise} ∵ Gaussian Noise reduces the quality by adding noise.

Advanced Augmentation : {frozen lake, glacier etc}

Adds effect of different season onto pic, day/night etc.

⇒ Conditional GANs {Converts image from one domain to other} ex: summer → winter, night → day
↳ high computation

⇒ Neural Style transfer : It grabs Style (Texture, Ambience, appearance) of one image & mix it with Content of other. {This provides effect similar to our Conditional GAN}

Downside ⇒ output look more artistic than realistic

↳ can be overcome by using Deep photo style transfer.

Transfer learning

⑤ Classifiers on top of Deep NN

- (a) Fully connected layer: Mainly uses Softmax activation function used for multi-class regression classifier
- (b) Global average pooling: Instead of adding fully connected layers on top of the conv. base, we add a global avg. pooling layer & feed its output directly into softmax
- (c) Linear SVMs: we can use sum to ~~the~~ classifiers the features output by the conv. base, it can increase the accuracy.

Activation function :

(or Transfer function)

→ used to determine the output of neuron like Yes/No.

Type:

① Linear Activation fn

→ output not bounded b/w range $(-\infty, +\infty)$

② non-linear Activation

(a) Sigmoid / Logistic Act. fn

$$\phi = \frac{1}{1+e^{-x}}$$

$\text{logit}(0, 1)$

∴ gives probability output.

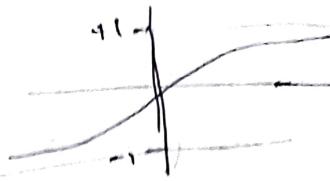
Soft Max \Rightarrow Multi-class Sigmoid Activation.

(b) Tanh / Hyperbolic tangent Activation fn



⇒ Tanh:

- Same like sigmoid but a bit better.
- Range $(-1, +1)$
- Negativ input \rightarrow slightly $-ve$
mean zero \rightarrow stays centered
- $+ve$ input \rightarrow $++ve$.

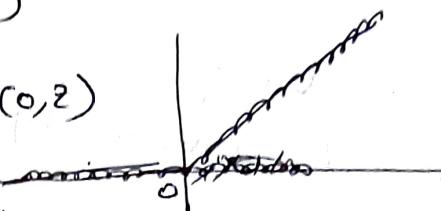


⇒ ReLU (Rectified Linear Unit)

$$R(z) = \max(0, z)$$

$$\text{Range} = (0, \infty)$$

- Not able to map $-ve$ values properly.

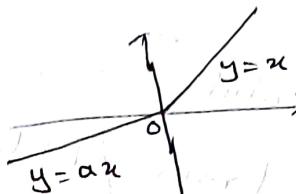


⇒ Leaky ReLU:

$$\alpha \approx 0.01$$

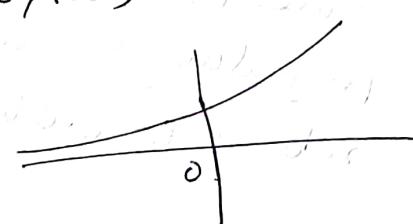
when $\alpha \approx 0.01 \Rightarrow$ Randomized ReLU.

$$\text{Range} = (-\infty, +\infty)$$



$$f = \ln(1 + e^x)$$

⇒ Softplus



Sigmoid vs Softplus

→ output probabilities are independent

$$\sum p_i = 1 \}$$

inc. likelihood of one class
dec. that of the other.

→ Trained with log loss or cross entropy regime.

→ Multiclass classification

→ Can be used for

binary classification as you can control the probability

Losses:

→ Loss functions

① Regression loss

Types ② Classification losses.

① Regression loss:

(a) MSE / Quadratic loss / L₂ loss:

$$MSE = \frac{\sum (y_i - \hat{y}_i)^2}{n}$$

→ does not consider error direction

→ Squares penalizes the outliers.

(b) Mean Absolute Error / L₁ loss:

$$MAE = \frac{\sum |y_i - \hat{y}_i|}{n}$$

→ does not consider error direction

→ use linear programming to calculate gradients.

→ ~~more~~ less robust to outliers.

③ Mean Bias Error:

$$MBE = \frac{\sum (y_i - \hat{y}_i)}{n}$$

→ Can be used to determine if model α has an overall +ve or -ve bias.

② Classification loss:

(a) hinge loss / multi-class SVM loss

- The score of correct category must be greater than the sum of all the scores in incorrect categories by some safety margin (usually one)
- Hinge loss is used for Maximum-Margin classification, i.e., SVMs.
- Convex function \therefore can be optimized easily.

$$\text{SVM loss} = \sum_{i \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

(b) Cross Entropy loss / Negative log likelihood:

- inc. as predicted probability diverges from actual label.

$$\text{Cross Entropy Loss} = -(y_i \log(\hat{y}_i)) + ((1-y_i) \log(1-\hat{y}_i))$$

↳ even if we have 2 classes (0, 1)

\therefore when actual $y_i=0 \rightarrow$ ① half discr.

& when $y_i=1 \rightarrow$ ② half discr.

- Penalizes heavily for predictions that are confident but wrong. or if it's ~~both~~ $y_i=0$ & $\hat{y}_i=1$

$$\text{if } L = 1 \times \log 0 = \infty \text{ or}$$

$$\text{if } L = 1 \times \log 1 = 0 \text{ if } y_i=1, \hat{y}_i=0$$

Dropout:

Keras : Sequential API
→ Feed forward

Functional API
→ Recurrent.

Handling Multi-collinearity in data:

Model's interpretability is affected by several problems such as :

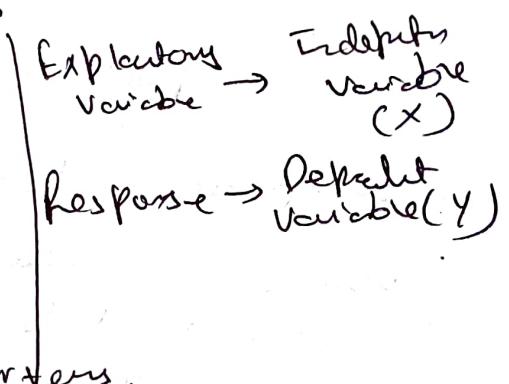
- ① Multi-collinearity
- ② Correlation of noise of error terms.
- ③ Non-linearity impact the Model's interpretability.

① Multi-collinearity:

→ When two or more explanatory variables are found to be closely co-related to each other, this correlation is referred to as Multi-collinearity.

→ It affects the interpretability of linear regression models

→ Due to this it is difficult to isolate the individual effect of explanatory variable on the response variable.



It results in following:

- ① uncertainty of coefficient (weight) estimates or
mistake variables: small changes in data can change
the coefficient.
- ② Increased Standard Error; the decrease accuracy of the estimate
& inc. the chances of detection.
- ③ Dec. Statistical Significance; Due to increased
standard error, t-stat. decline which negatively
impact the capability of detecting statistical
significance in coefficient leading to type II error.
- ④ Predicting coefficient & p-value: The influence of
the correlated explanatory variable to model
due to collinearity
- ⑤ overfitting: Leads to over fitting or is indicated by
high variance.

Multi-collinearity detection:

- Correlation matrix,
- heatmaps are definitely more intuitive & visual. However,
it ~~not~~ helps to identify correlation b/w 2 variables
strictly and fails to identify collinearity which exist
b/w 3 or more variables, for which "Variance Inflation
Factor" can be used.

Variance Inflation Factor (VIF): VIF is the ratio of
variance of coefficient estimate when fitting the full model
divided by the variance of coefficient estimate if fit on
its own. Impossible value is '1' which mean no-
collinearity, if $VIF > 5$ then collinearity must be addressed.

Dealing with Multi-collinearity:

- ① Introduce Penalization or remove highly correlated Variables: by Lasso / Ridge (L_1 / L_2 Regularization) and/or dimensionality reduction by observing VIF values.
- ② Combine highly correlated variables: using methods like PCA

~~CKD Trees~~:

Efficient alternative of KNN?

f f. K-D-Trees:

→ k-dimensional trees,

Another Alternative of KNN
→ ball tree

Decision Trees:

- A decision tree can be used to visually and explicitly represent decisions and decisions naturally.
- A decision tree is drawn upside down. Internal node contains the conditions according to which the node splits into edges.
- Non-parametric Bayesian learning used for both regression & classification.
Decision Tree $\xrightarrow[\text{algo}]{\text{to}}$ CART (Classification & Regression Trees)

Q) What's going on in the background?

- It involves deciding on which feature to choose
- What condition to use for splitting.
- Knowing where to stop.
- As a tree generally grows arbitrarily, we will need to trim it down for it to work beautifully.

Techniques used for splitting:

- ① Recursive Binary Splitting: In this process:~~all~~
 - All the features are considered
 - Different split points are tried and tested using a cost function.
 - The split with the best (lowest) Gst is selected.
 - Greedy Algorithm, as we have an excessive desire of lowering the Gst.
 - This makes the root node as best predictor / classifier

② Cost of a split:

In both classification & regression, we try to find the cost function try to find most "homogenous branches", or branches having groups with "similar responses".
 Regression: $\text{Sum} (y - \text{predict}_k(y))^2$ (interpretive)
 Cost function

Decision Trees: for nonlinear decision making with simple linear decision surface.

Classification: Entropy (E) = $-\sum_{i=class} p_i \log(p_i)$ → probab. of i class in tree & split.
 (Gini Index)

$$\text{Information Gain} = E(\text{Parent}) - \sum w_i E(\text{child}_i)$$

Model compares every possible split. & ~~Trees~~ choose those that maximize the IG (Info. Gain)

Regression: To check impurity use "Variance Reduction" technique.

$$\text{"Normalized Mean } Y \text{"} \left. \right\} \quad \left\{ \text{Var} = \frac{1}{n} \sum (y_i - \bar{y})^2 \right\}$$

$$\left. \left\{ \text{Variance Reduction} = \text{Var}(\text{Parent}) - \sum w_i \text{Var}(\text{child}) \right\} \right. \rightarrow \text{impurity more.}$$

Random forest classifier:

(Binary classifier)

⇒ Why we use Random forest when we have D-trees?

- If there is a slight change in data then the whole decision tree needs to be rebuilt.
- Highly sensitive to training data, which shows high variance & overfitting. And we need to generate our model.

⇒ Random forest ⇒ Is a collection of multiple Random Decision trees & therefore is much less sensitive to the training data.

Step 0 ^(Bootstrapping)
Build new data set from our original data.
randomly select rows ~~some~~ with replacement,
to make a new dataset.

Ex: if our original data has M rows. Then
we will create random data with M rows
& select rows with replacement.

$$\therefore \text{No. of such random trees} = M \times M \times M \times \dots \times M$$

$$\begin{array}{c}
 \text{Data} \\
 \begin{array}{c}
 \text{O} \\
 \text{O} \\
 \vdots \\
 M \\
 \text{O} \\
 \text{O} \\
 \vdots \\
 M \\
 \text{O} \\
 \text{O} \\
 \vdots \\
 G
 \end{array}
 \end{array}
 \quad
 \begin{array}{c}
 \square^M \\
 \square^M \\
 \vdots \\
 \square^M \\
 \square^M \\
 \vdots \\
 \square^M \\
 \square^M \\
 \vdots \\
 \square^M
 \end{array}
 \quad
 = \frac{M^M}{-}$$

"Random Feature Selection"

for building tree Class we do not choose all the Class we take different set of Class at a time

e.g.: original data

	x_1	x_2	y
0			
1			
2			

New data:

0
0
1

or

1
0
2

x_1, x_2

No. of rows No. of Col

$$\therefore \text{Total possible such data} = \underline{(m \times 2)}$$

Step 2°: Make tree for all such data.

Step 3°: we pass the data to be predicted through each tree one by one and make tree prediction.

Step 4°: we combine all the prediction by a Majority Voting.

Aggregation Step.

Result: $(\text{Bootstrapping} + \text{Aggregation}) = \text{Bagging}$

Random feature selection step: Helps to dec. tree & effect of Correlation b/w the features.

Initial size to select the no. of features (cols) to take into account

$$O_{out} = \log(\text{Total features}) \text{ OR } \sqrt{\text{Total features}}$$

Now Random Forest for regression problem.

In Step 4: while combining the prediction just take average.

Dealing with Imbalanced data:

- ① Change performance metric: Accuracy can be very misleading or in such cases dummy classifier can also give a high accuracy.
- ② Change the algo: Use Decision tree
- ③ Resampling Techniques:
 - (a) oversample Minority class
 - (b) undersample Majority class
- ④ Generating Synthetic Samples:
imblearn's SMOTE (Synthetic Minority OverSampling Technique)
Use Nearest neighbor to generate new and synthetic data we can use for training our model.

Categorical Encoding: changing ~~the~~ non-numerical data to numeric

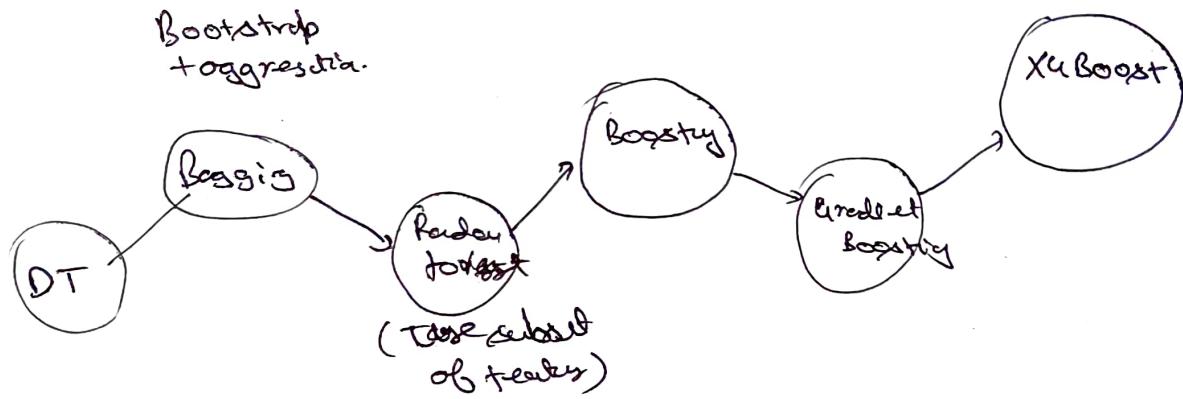
- ① Label Encoding
- ② One Hot Encoding.

XGBoost Algorithm:

(Extreme Gradient Boosting)

→ Decision tree-based ML algo. that uses gradient boosting framework.

Evolution of Decision Trees:



Boosting: Models are built sequentially by minimizing tree error from previous model cellwise increasing (or boosting) influence of high performing Models.

Gradient Boosting: Gradient Boosting employs Gradient descent algo to minimize errors in Sequential models.

XGBoost: Optimized Gradient Boosting algo through parallel processing, tree pruning, handling missing values and regularization to avoid overfitting/bias.



post office of the USSR and the Central
Committee of the Communist Party of the
Soviet Union, Moscow, Soviet Union

Post office of the USSR and the Central
Committee of the Communist Party of the
Soviet Union, Moscow, Soviet Union

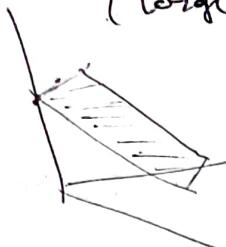
Post office of the USSR and the Central
Committee of the Communist Party of the
Soviet Union, Moscow, Soviet Union

Post office of the USSR and the Central
Committee of the Communist Party of the
Soviet Union, Moscow, Soviet Union

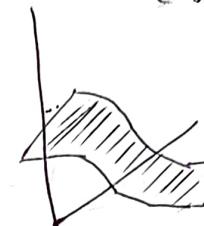
Post office of the USSR and the Central
Committee of the Communist Party of the
Soviet Union, Moscow, Soviet Union

Bias - Variance Trade Off:

- Model interpretability vs Model flexibility
 (large assumptions) (few assumptions)



ex: linear



ex: polynomial

Bias: Error introduced due to assumptions made to model a ~~too~~ complex real-life problem by an approximately simple model.

→ More flexible models take less assumptions about the data result in less bias, however they might lead to underfitting or overfitting?

Variance: Amount by which the estimates or results of the ML model change with a different training dataset.
 More flexible model has more variance. → leads to overfitting.

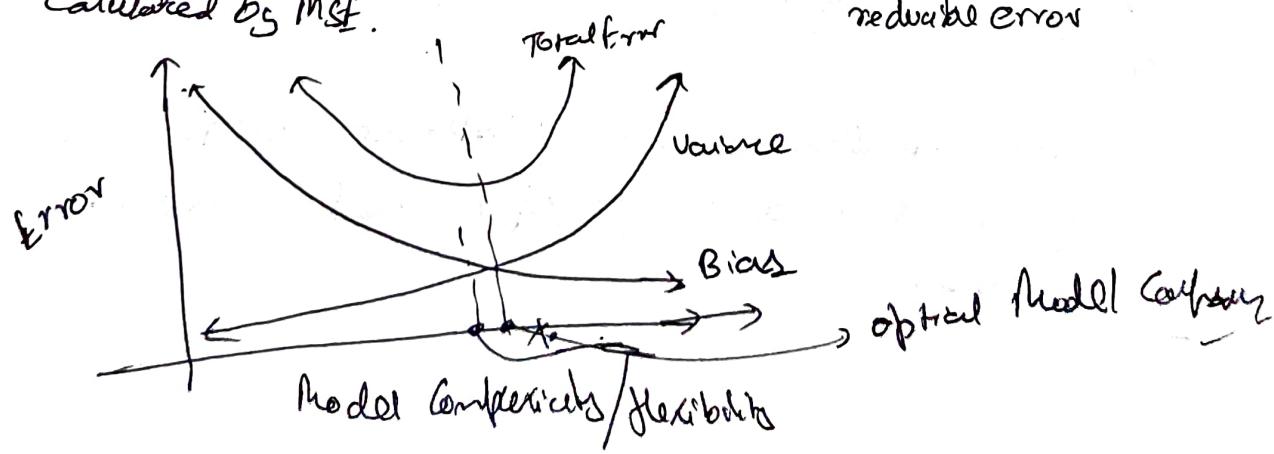
Trade off

$$\text{Expected Test Error} = \text{Irreducible Error} + \text{Bias}^2 + \text{Variance}$$

due to noise

Irreducible error
reducible error

Calculated by MSE.



Ways to reduce Bias & Variance :

- ① Bias can be reduced by inc. Model Complexity by
 - Boosting
 - Add more features or introduce polynomial feature.
- ② Boosting :
 - unlike many ML models which focus on high quality prediction done by a single model, boosting algo seeks to improve the prediction power by training a sequence of weak models, each compensating the weakness of its predecessors.
 - It is generic

Types : Based on different definition of weakness and their corresponding algorithm. Some are →

- ① Adaptive Boosting (AdaBoost) ② Gradient Boosting.
- ① Adaptive Boosting (AdaBoost) :

 - for classification Problems & regression.
 - weakness is identified by the weak estimator's error for a weak classifier C :
$$X: m \times d ; Y: m \times k, \text{ Sample weights } W: m \times 1$$

$$\text{Error} = \sum_{j=1}^m w_j I(C(x_j) \neq y_j) ; I(x) = \begin{cases} 1, & \text{if } x = \text{True} \\ 0, & \text{if } x = \text{False} \end{cases}$$

no. of classes

no. of feature

→ The weak learners in AdaBoost are Decision Trees with a single Split, called decision Stumps.

⇒

→ AdaBoost works by putting more weights on difficult to classify instances (~~features~~) & less on those features already handled well.

* Ensemble : ML techniques that just combine several base models in order to produce one optimal predictive model

→ Does same splits like tree in Decision Tree.

Each Stump choose a feature X_i & a threshold T_{ij}
& classifier checks for the condition on X_i with T_{ij} Node

→ This happen for can be tried out for all the features
& all possible Threshold look for the best accuracy or the Information Gain.

→ While it naively seems like there are ∞ no. of choices for threshold, two diff. Threshold are meaningful if it divide the data in diff. Sides of Split.

→ To try every possibility, then we can sort the examples by feature in the question & try one threshold fully b/w each adjacent pair of example.

(Although the Algo is fast if it tries all possible T's).

⇒ AdaBoost → forest of Stumps.

D-Trees	Random Forest	AdaBoost
Tries all features	Randomly selects data	Only Decision Stumps are used i.e. (one feature per tree)
To make tree best possible	Selects feature to be used to make many trees	diff. trees are given to different stumps.
+ tree	All trees are given equal weight with calculated output.	all the stump Stumps are combined with diff. weights for output
→ 1 tree is made		

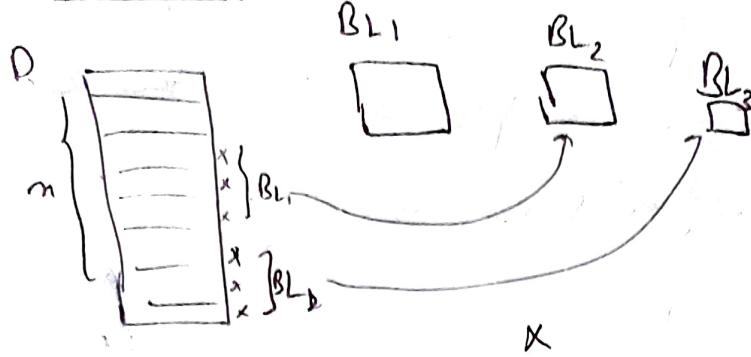
- In Adaboost order of tree + trees (stumps)
- Netto, because the Error that the prev. Stump
Make influence how the next Stump is made.
- all stumps get the same weight at the start
i.e. (y_{true}) , & change as the other stumps get
created.
 $\xrightarrow{\text{dataset}}$
- Total Error of Stump. = Sum of weights associated
(lies in $[0, 1]$) with tree in correctly classified
sample.

$$\text{Amount of Bay} = \frac{1}{2} \log \left(\frac{1 - \text{Total Error}}{\text{Total error}} \right)$$

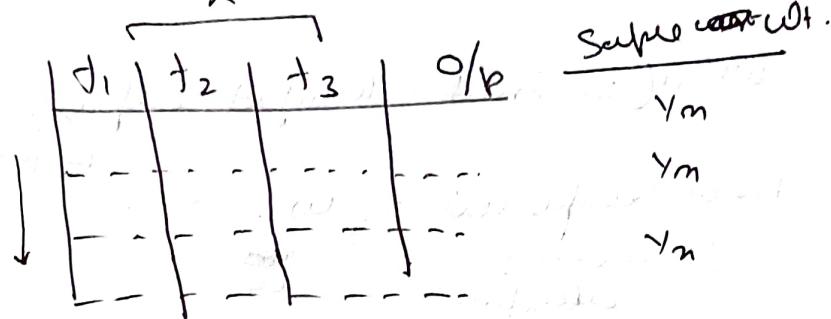
- New Sample weight (out of Bay)
 $\xrightarrow{\text{for this sample}} = \text{Sample weight} * e^{-\text{amount of Bay}}$
 $\xrightarrow{\text{theo. incorrectly classified}}$
- Do New sample weight
 $\xrightarrow{\text{Dec. the sample weight}} = \text{Sample weight} * e^{-\text{amount of Bay}}$
 $\xrightarrow{\text{to class of}}$
 $\xrightarrow{\text{Samples that correctly classified}}$
- Now Normalize the new sample weight.

see first Adaboost

Boosting



AdaBoost :

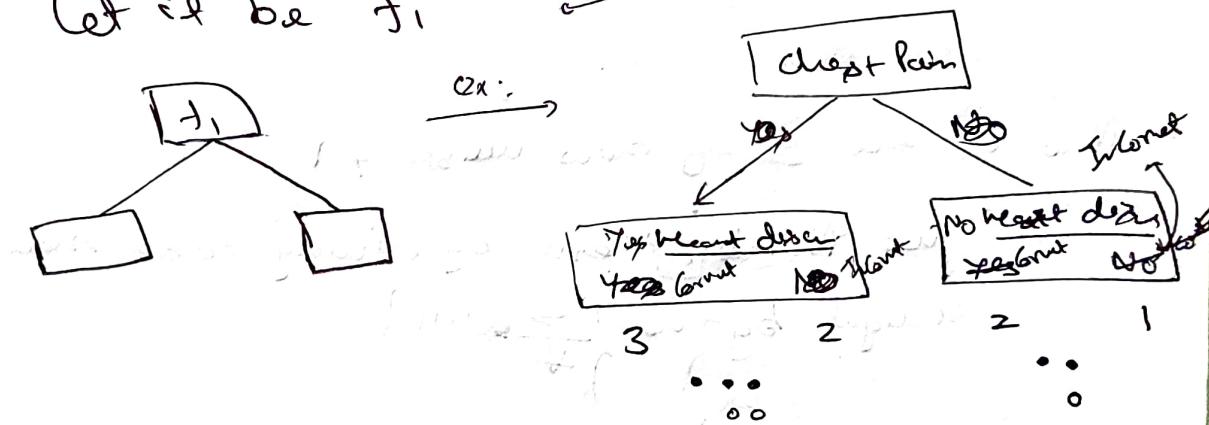


initially $\Rightarrow w = y_n$



Below Entropy Step is Selected.

Let it be 't_1'



Total error = Sum of weights associated
(t_h) with incorrectly classified each

$$\text{Total error} = \left\{ 2 \times \frac{1}{n} + 1 \times \frac{1}{n} \right\}$$

Amount say :

$$\text{Performance of Step} = \frac{1}{2} \log \left(\frac{1 - TE}{TE} \right)$$

Now as the wrongly classified records pass to the next ~~Passive~~ Step. i.e inc. the weights of wrongly classified sample & dec. the old wt. for correctly classified ones. Using the "Performance of Step" = P

① Inc. wt. of Incorrectly Classified samples,

$$\text{New Sample wt.} = \left\{ \begin{array}{l} \text{wt.} \times e^{\text{Performance}(P)} \\ \text{for Incorrectly} \\ \text{classified.} \end{array} \right.$$

② Dec. the wt. of ~~old~~ Correctly Classified samples.

$$\text{New Sample wt.} = \left\{ \begin{array}{l} \text{wt.} \times e^{-\text{Performance}(P)} \\ \text{of correctly classified.} \end{array} \right.$$

\Rightarrow Now as the \sum of new weights $\neq 1$

\therefore we Normalize them by dividing each new weight by the \sum of new wt.

$$\text{i.e. } \left\{ \begin{array}{l} t = \sum_{i=1}^n w_i \\ w_i' = \frac{w_i}{t} \end{array} \right.$$

\Rightarrow Now new dataset with n record (without replacement) is created for next step

randomly based on the weights of the sample (the wrong records have higher wts.)

\Rightarrow & Continue the same steps till we obtain purity or tolerance.

⇒ Aggregate the results of all the Stumps. ∴ Majority Vote.

Gradient Boosting

① For regression & classification:

② Regression:

			BaseModel		
			\hat{y}	f_1	f_2
Ex:	Exper	Deg.	Salary		
	2	BE	50	75	
	3	MS	70	75	
	5	MT	80	75	
	6	PhD	100	75	

Step 1: Base Model 1 → Avg. of two output

Step 2: we Compute Residuals Errors, Pseudo Residuals.

Calculate loss or Residuals (R_1)

Step 3: Construct DecisionTree $\{x_i, R_1\}$

feature score

The decision tree outputs the residual from the given input features.

If we apply pass to all the data sample features in the Decision tree it will give residuals R_2 which will be near equal to R_1 .

To

Step 4: Final prediction value is calculated by doing.

~~\hat{y}_1~~ ~~\hat{y}_2~~

$\{\hat{y} - R_2\}$

Base Model
value

DT-residum.

But this is an over fitting problem

which will have low bias (because of overfitting)

& high variance (because of overfitting).

→ What we do is that we add a learning rate to the tip

∴ we do: $\hat{y} + \alpha(R_2)$ learning rate

$$\left[\hat{y} + \alpha(R_2) \right] \quad \left\{ \alpha \in (0, 1) \right\}$$

∴ in our example: $\alpha = 0.1$

$$= 75 + 0.1(-23)$$

$$f_2 = \underline{\underline{73.7}}$$

⇒ Now we create a new Decision Tree based on the new R_2 value keeping features same.

∴ what we can do:

$$f(x) = h_0(x) + \underbrace{\alpha_1 h_1(x)}_{\text{base}} + \underbrace{\alpha_2 h_2(x)}_{\text{DT}_1} + \underbrace{\alpha_3 h_3(x)}_{\text{DT}_2} + \dots + \alpha_m h_m(x)$$

$$= \boxed{f(x) = h_0(x) + \sum_{i=1}^m \alpha_i h_i(x)}$$

XGBoost : Extreme Gradient Boosting

No Classifier -

Ex: Get loan or not \rightarrow

Salary	Credit	Approved (%)	\hat{y}	R_i
$\leq 50K$	B	0	0.5	-0.5
$\leq 50K$	G	1	0.5	0.5
$\geq 50K$	G	1	0.5	0.5
$\geq 50K$	B	0	0.5	-0.5
$> 50K$	N	1	0.5	0.5

Bad, Good, None.

\rightarrow used for both binary & multi-class classification.

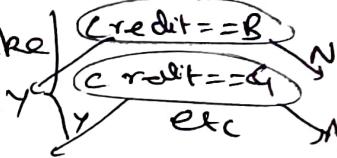
Step 1: Base Model \Rightarrow Binary classification \therefore take $\hat{y} = \frac{0+1}{2} = 0.5$

$$f_1(\text{Residual}) = \hat{y} - \hat{y}$$

$(Pr = 0.5)$
(mostly)

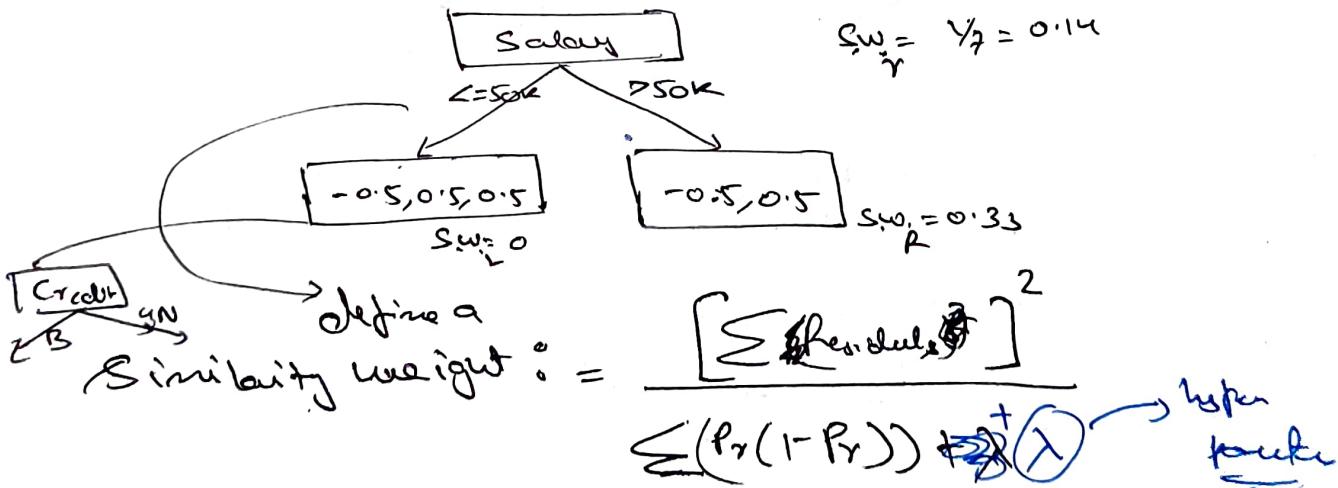
Now construct a tree with Residuals R_i .

\Rightarrow In XG boost tree node always gets divided into two classes (even if there are 3 classes) if like



Ex:

$$R \rightarrow [0.5, 0.5, 0.5, -0.5, 0.5]$$



$$\text{Gain} = S.W_{\text{left}} + S.W_{\text{right}} - S.W_{\text{root}} = 0 + 0.33 - 0.14 = 0.12$$

→ do the above recursively every time taking care of highest gain condition.

• Post Pruning:

In case of ~~nodes~~ types of values which had more than two types of values we have to decide after we should divide or not.

Calculated by easily 'Cover Value' $\Rightarrow \boxed{\Pr(1 - \Pr)}$

Now we check if of Node,

$\text{Gain} < \text{CoverValue}$, then we just cut the branch, ~~stop~~

to calculate,

Base Model output w.r.t the initial probabilities (p_0)

$$\log(\text{odds}) = \log\left(\frac{p_{i=1}}{1-p_{i=1}}\right)$$

Now New probability

The for doing predictions, we traverse the Binary tree & find the similar wt. occurred to it sw. then, the our prediction will be.

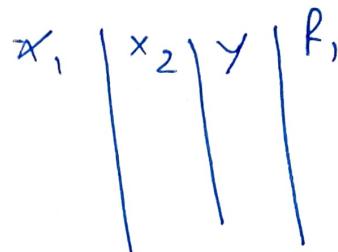
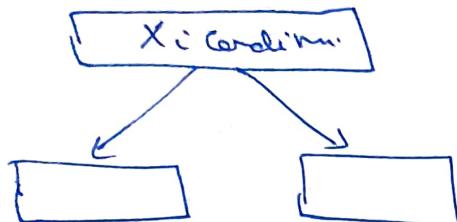
$$P_i = \frac{1}{1 + e^{-\sum_{i=1}^n \log(\text{odds}) + \sum_{i=1}^n \alpha_i(sw)}}$$

Now repeat the whole process again with new probabilities, until

XGBoost regressor

Base Model $\rightarrow f_i = \cancel{\text{Resid}}(y_i) - \text{Avg}(y_i)$

Make tree (Binary only)

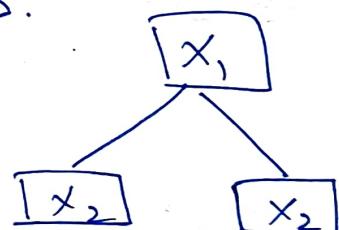


first similarity wt. = $\frac{[\sum (\text{Resid})]^2}{\text{No. of residuals} + \lambda}$ hyperparameter

$$\left\{ \text{Gain} = S.W_{\text{left}} + S.W_{\text{right}} - S.W_{\text{root}} \right\}$$

Take "Maximum Gain" Condition

\rightarrow keep doing this recursively ~~and stop~~ to fit all the max gain ~~for~~ features.



To predict the Y,

we ~~traverse~~ traverse the tree and reach the right leaf.

The ~~is~~ $T_i = \text{Avg of all the values stored in that leaf.}$

~~Output New residual~~ $\hat{y}_i = \cdot (\text{BaseModel value}) + \sum_{\text{trees}} (x_i \cdot T_i)$

\rightarrow repeat the same process until $y = \hat{y}_i$

$\Rightarrow \gamma$ (gamma parameter) } Reduces over fitting
 ↳ used for post prunig } (higher parameter)

$$\text{Cover Value} = (\cancel{\gamma - \text{train}}) C_{\text{train}} - \gamma$$

if } Cover Value < 0 } \rightarrow then we can
 post prune tree
 tree branch i.e
 cut it.

CNN Revisited :

Layer 1 : Vertical Edge detector \Rightarrow

3x3	1	0	-1
	2	0	-2
	1	0	-1

without padding : $\frac{(m-f)+1}{s} = L_1 \text{ (size)}$

Stride filter size

\downarrow next layer

with padding : $m' - f + 1 = m$

so that $m' = \underline{\underline{m + f - 1}}$
 we want \rightarrow

\rightarrow zero padding

\rightarrow second \rightarrow nearest value padding.

$P_i = \text{padding divisor}$

$$\therefore \text{next layer size} = \{m + 2p - f + 1\}$$

$(m+f-1) + 2p + 1$
 stride

Max Pooling : Based on "location invariant"

↳ uses

captures higher pixel values to refine the image.

Other pools \Rightarrow Mean Pooling, Min/Avg Pooling.

CNN \rightarrow Do dimensionality reduction by keeping important features.

F_B Metrics :

$$\text{if } F_B = \frac{(1+\beta^2) \text{ Precision} * \text{Recall}}{\beta^2 * \text{Precision} + \text{Recall}}$$

$$\left. \begin{array}{l} \text{if } \beta=1 \quad F_1 = \frac{2 \text{ Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \\ \text{Score} \end{array} \right\} \text{Harmonic Mean}$$

$$\beta = 0.5$$

$\beta = 1$ { when FP & FN both are equally important}

if $\text{FP} \gg \text{FN} \xrightarrow{\beta \in [0, 1]} \text{reduce } \beta \Rightarrow \beta = 0.5$

if $\text{FP} \ll \text{FN} \xrightarrow{\text{impossible}} \text{inc } \beta \Rightarrow \beta = 2$



ROC & AUC Curves

$$TPR = \frac{TP}{TP+FN}$$

→ Binary Classification.

→ Used for deciding the threshold probability

FPR

Carrying threshold
other than 0.5 as per requirement

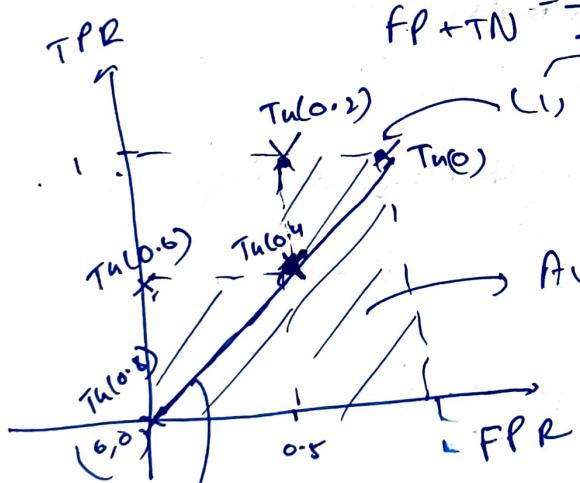
Threshold Value = [0, 0.2, 0.4, 0.6, 0.8, 1]

y	\hat{y} (probability)	$y(0)$	$y(0.2)$
1	0.8	1	1
0	0.96	1	1
1	0.4	1	1
1	0.3	1	1
0	0.2	1	0
1	0.7	1	1

$$TPR = \frac{TP}{TP+FN} = \frac{4}{4+0} = 1$$

$$FPR = \frac{FP}{FP+TN} = \frac{2}{2+0} = 1$$

$$\begin{cases} TPR = 1 \\ FPR = 0.5 \end{cases}$$



AUC = Area under the curve

More the area better our model.

→ good models area must be atleast area under the line because \rightarrow coin tossing would be better than this if it ~~has~~ prob < 0.5

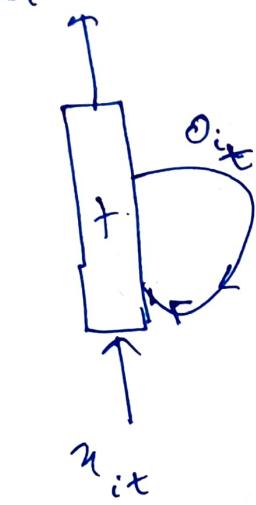
$\text{Th}(0.6) \rightarrow 0 \text{ FPR}$,

Now if we want Low FPR & high TPR we can select $\text{Th}(0.6)$ in graph.

If we don't care about FPR & just want high TPR we can select $\text{Th}(0.2)$.

If we want to focus on TPR & FPR both then we can select $\text{Th}(0.4)$.

RNN:

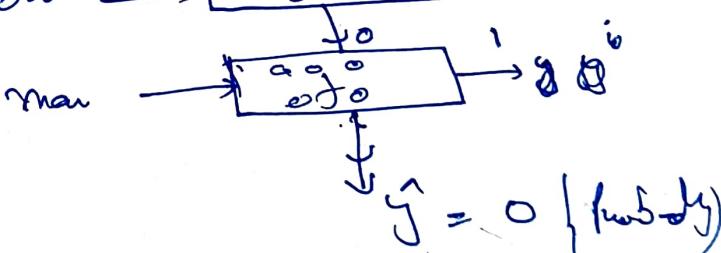
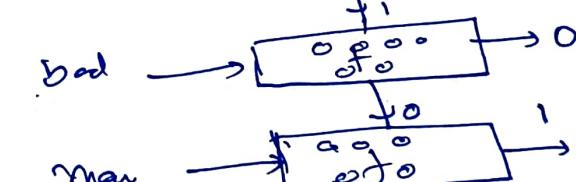
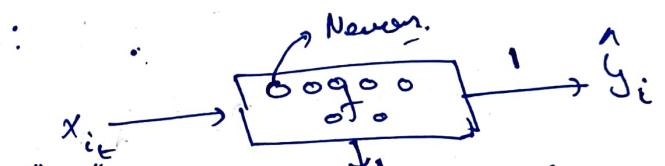


"you are bad man"

Sentiment Analysis

0 → Neg
1 → +ve

ex:

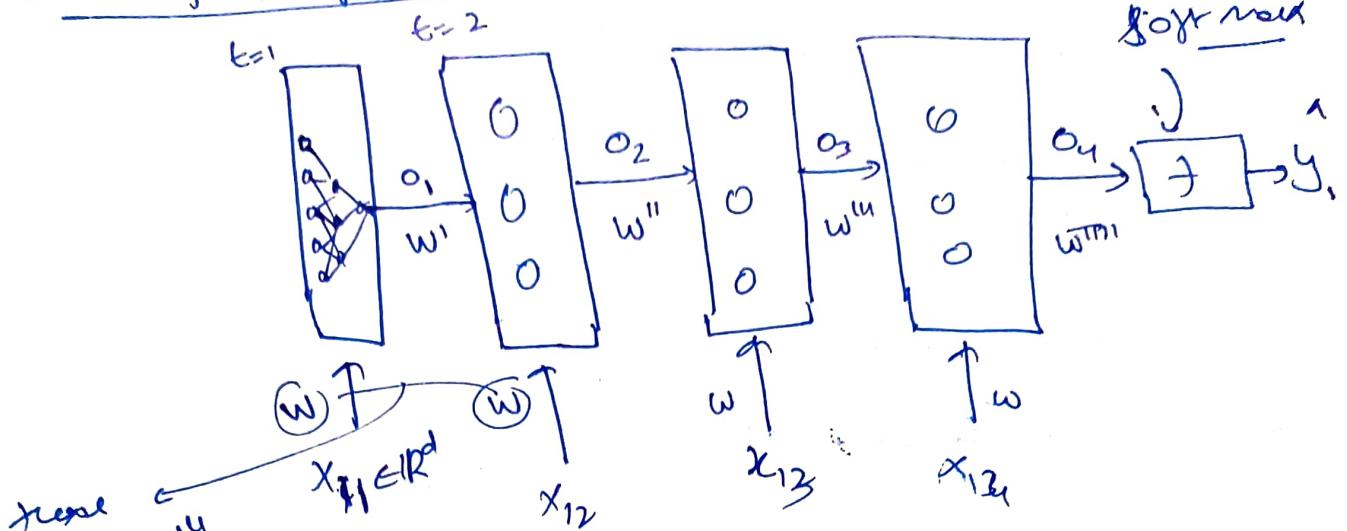


$$\hat{y} = 0 \{ \text{fwd}\}$$

$$\therefore o_i = f(x_i, w)$$

Let say we start with weights "W"

forward prop. of RNN:



These
weights will
be same
and only
change during
"Back propagation".

at time step t

$$o_1 = f(x_{11} * w)$$

$$o_2 = f(x_{12} * (w + w'))$$

$$o_3 = f(x_{13} * w + o_1 * w')$$

$$o_4 = f(x_{14} * w + o_2 * w'')$$

$$\hat{y}_1 = f(o_4 * w''')$$

$$\hat{y}_1 = f(o_4 * w''')$$

$$y_{loss} = (\hat{y}_1 - y)$$

Next reduce the loss in back prop

Backprop of RNN:

$$L_{\text{loss}} = (\hat{y} - y)$$

$$\frac{\partial L}{\partial \hat{y}_i}, \frac{\partial L}{\partial w^{'''}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w^{'''}}$$

$$\left\{ w''' = w''' - \frac{\partial L}{\partial w'''} \right\} \xrightarrow{\text{update w'''}}$$

$$\left\{ \frac{\partial L}{\partial w} = \frac{\partial L}{\partial w'''} \cdot \frac{\partial w'''}{\partial w_i} \right\}$$

!

Problem in simple RNN:

* # Vanishing & Exploding Gradient Problem -

Vanishing Gradient: When during back propagation we find the derivative of functions like "Sigmoid" the derivative lies b/w [0, 1] but as we go from last layer to the first layer of back propagating the value of the derivative being calculated again and again becomes very small: Ex: $L_1 \xrightarrow[w_1]{w_2} L_2 \xrightarrow[w_2]{w_3} L_3 \rightarrow \text{loss} = (\hat{y}_3 - y)$

$$\begin{aligned} \frac{\partial L_2}{\partial w_2} &= \frac{\partial L_3}{\partial w_3} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w_3} \\ &\quad \curvearrowleft \curvearrowleft \curvearrowleft \cdot \frac{\partial L_3}{\partial w_3} = \frac{\partial L_3}{\partial w_3} \cdot \frac{\partial w_3}{\partial w_2} \end{aligned}$$

\therefore due to this our step size will continually dec. \therefore we will not be able to reach global minimum. This is called vanishing gradient problem.

Exploding Gradient Problem:

if we use activation function like 'ReLU'
 $\text{sum derivative(ReLU)} > 1$

- The inc. the δ gradient at every step back which creates an exploding gradient problem



To avoid this we use LSTM (Long Short-term Memory)

LSTM :

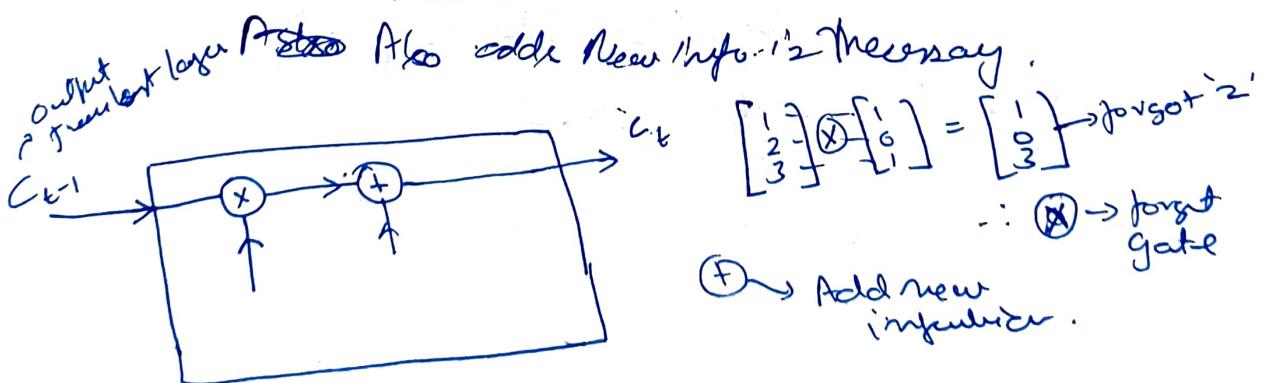
(Unated RNN)

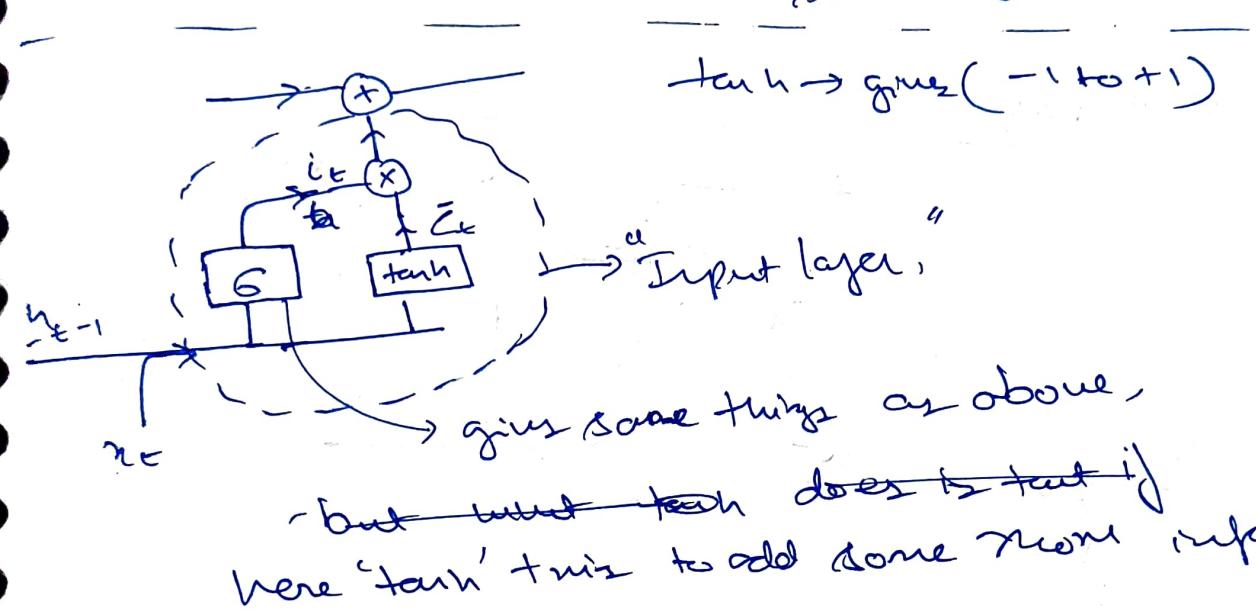
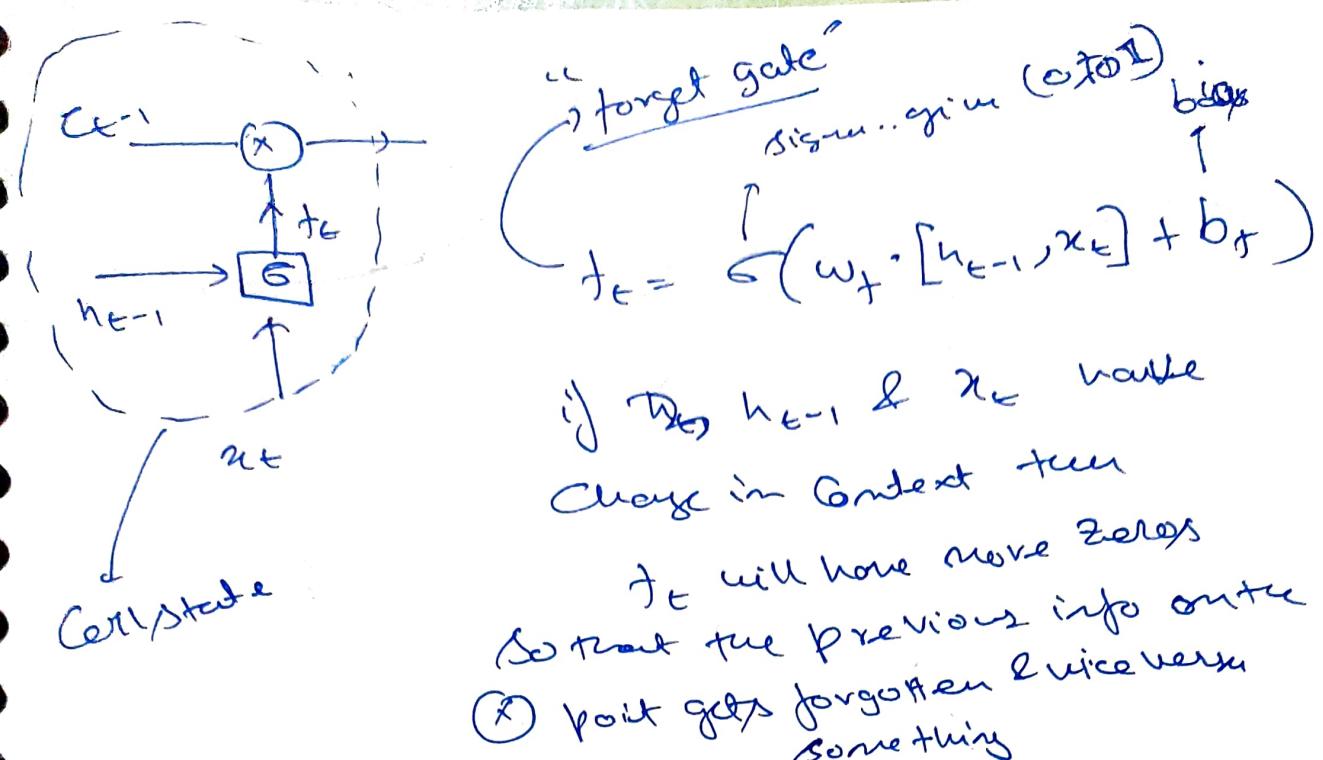
- ① Memory cell
- ② forget gate
- ③ I/P Gate
- ④ O/P Gate

 Neural Network layer
(MLP)

- Point wise operator
 - Vector traxfer
 - Concatenate
 - Copy

① Memory cell : Used to remember & forget things based on context of Input.

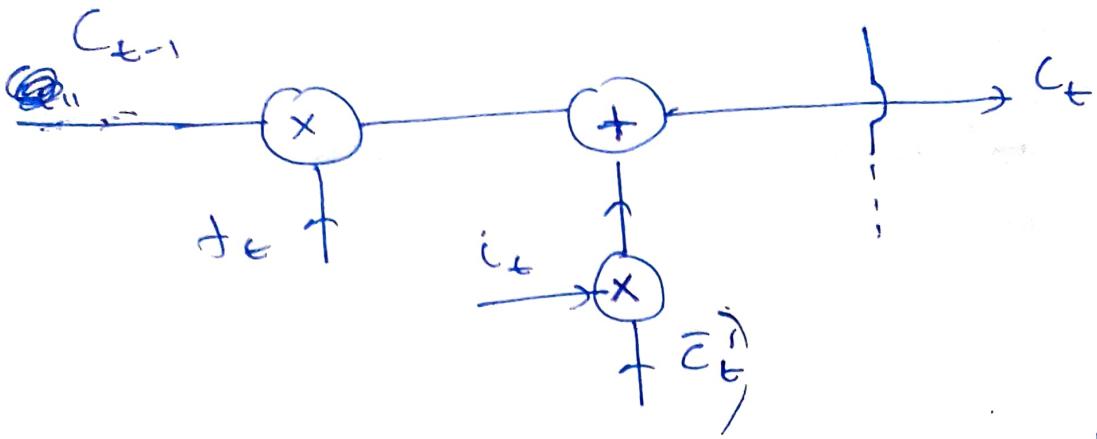




$$f_t = \sigma(w_f \cdot [h_{t-1}, x_t] + b_f)$$

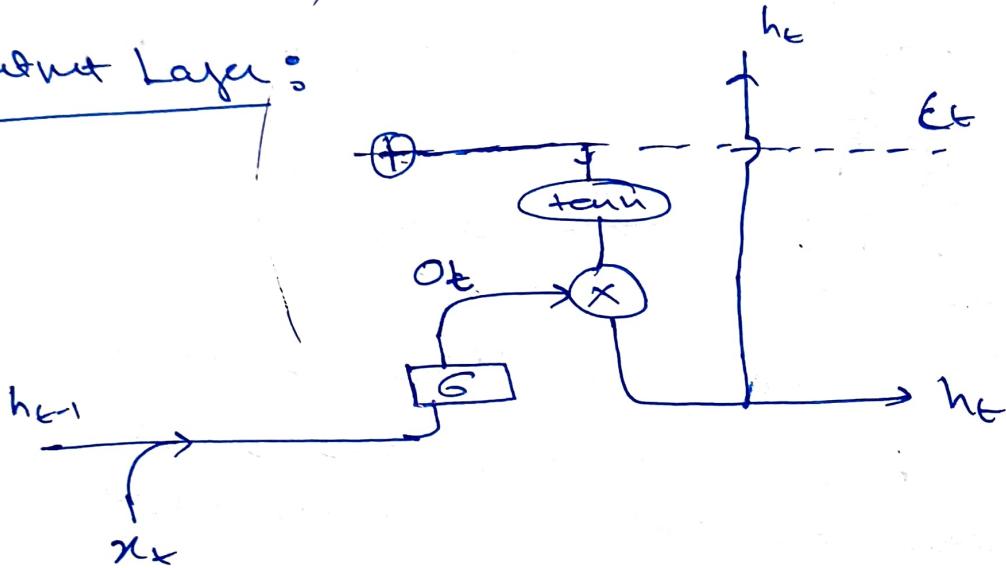
$$\bar{c}_t = \tanh(w_c \cdot [h_{t-1}, x_t] + b_c)$$

"tanh" here acts like the MaxPooling layer in CNN.



$$\left\{ c_t = f_t * c_{t-1} + i_t * \bar{c}_t \right\}$$

Output Layer:



$$o_t = \sigma [W_o [h_{t-1}, x_t] + b_o]$$

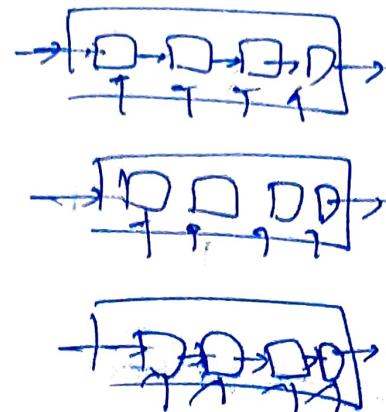
$$h_t = o_t * \tanh(c_t)$$

$\therefore h_t \rightarrow$ will have More Many Features
than that of c_t but to one
More Point wise operate.

LSTM Channels in NN

#Types of RNN

- ① One to one
- ② one to many
- ③ Many to one
- ④ Many to many.



Algorithmic Trading

M① NN-based regressor model.

M② RNN Model regression. (LSTM)

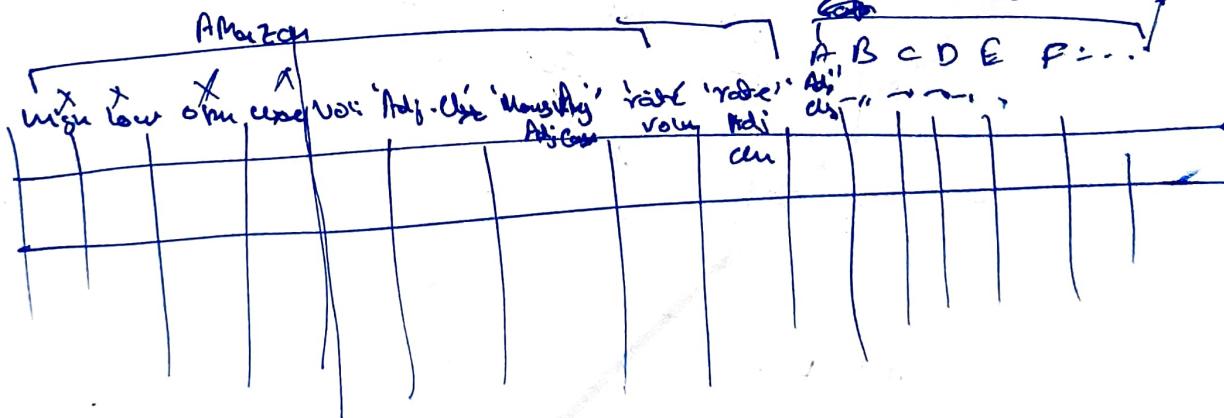
M③ Takes both the above predictor & gives the output.

S&P 500 Companies: Data to predict the stocks at Amazon & Microsoft affect the other companies affect it.

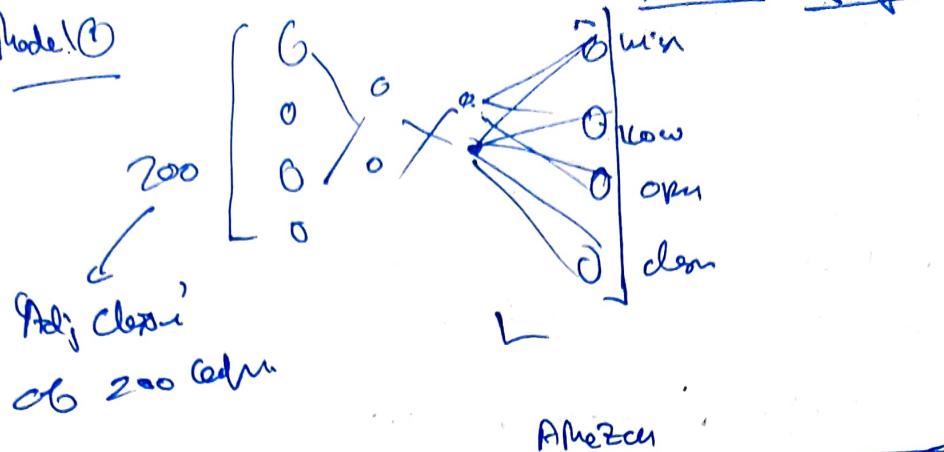
→ BeautifulSoup (Data Scrape/Get)

Using 199 company stock data to predict the 200th stock

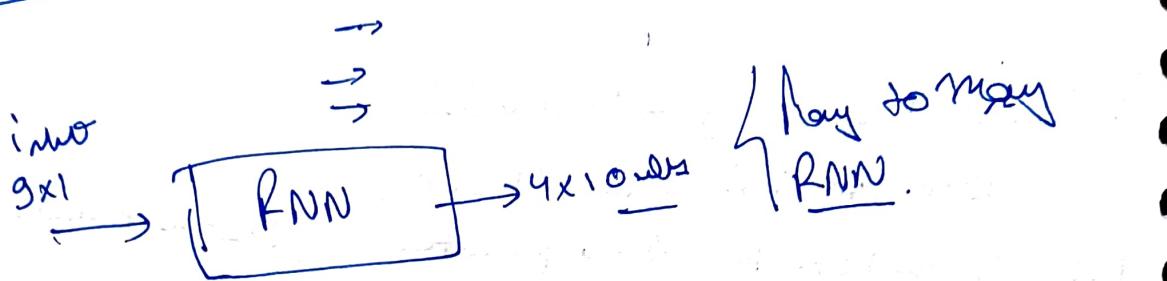
→ 'Adj Close' was used for 199 companies



Model①



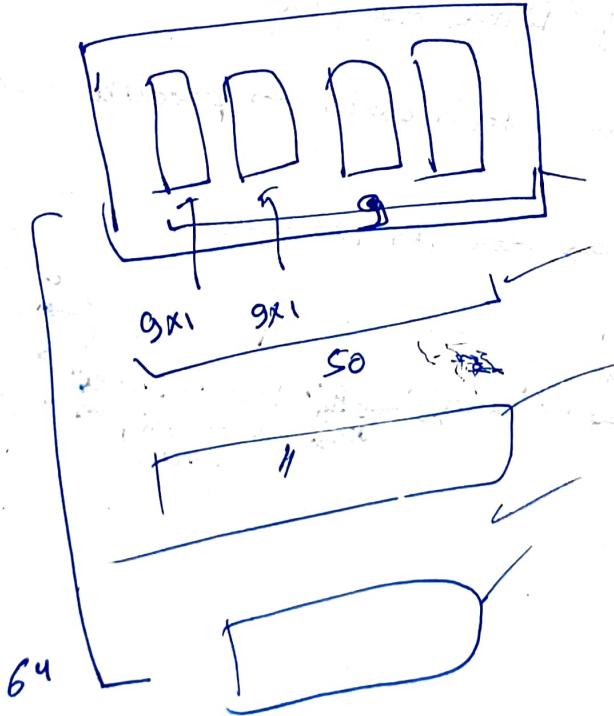
Model②: RNN:



Next set of 50 data entries.

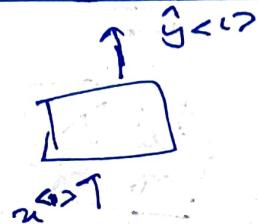
Denoising

RNN



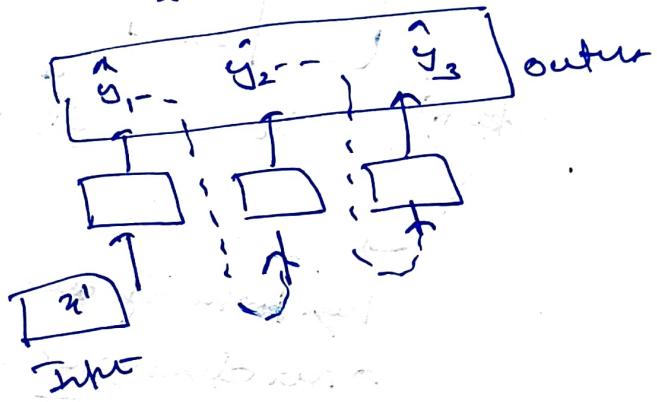
RNN's

① one to one:

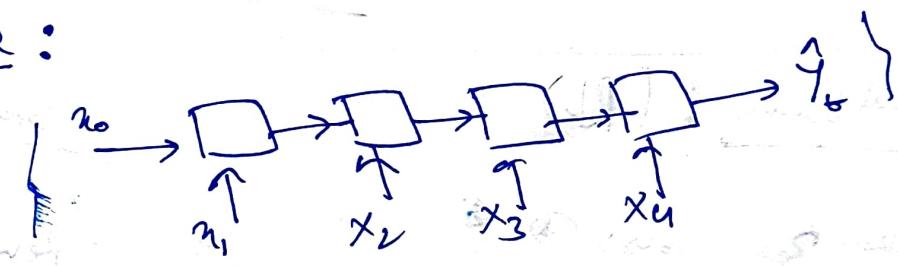


② one to many:

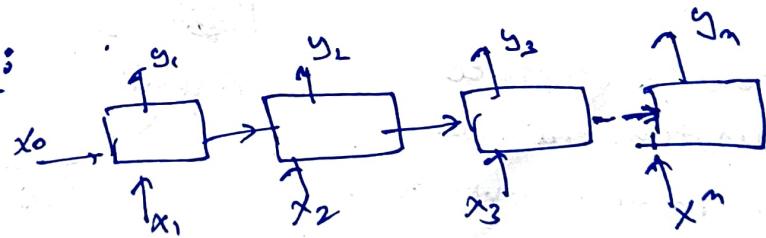
Ex.: Music generation
from single note



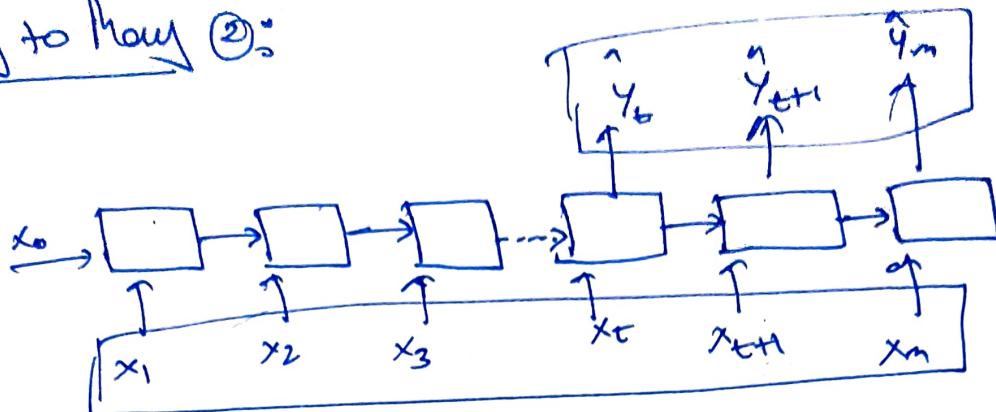
③ Many to one:



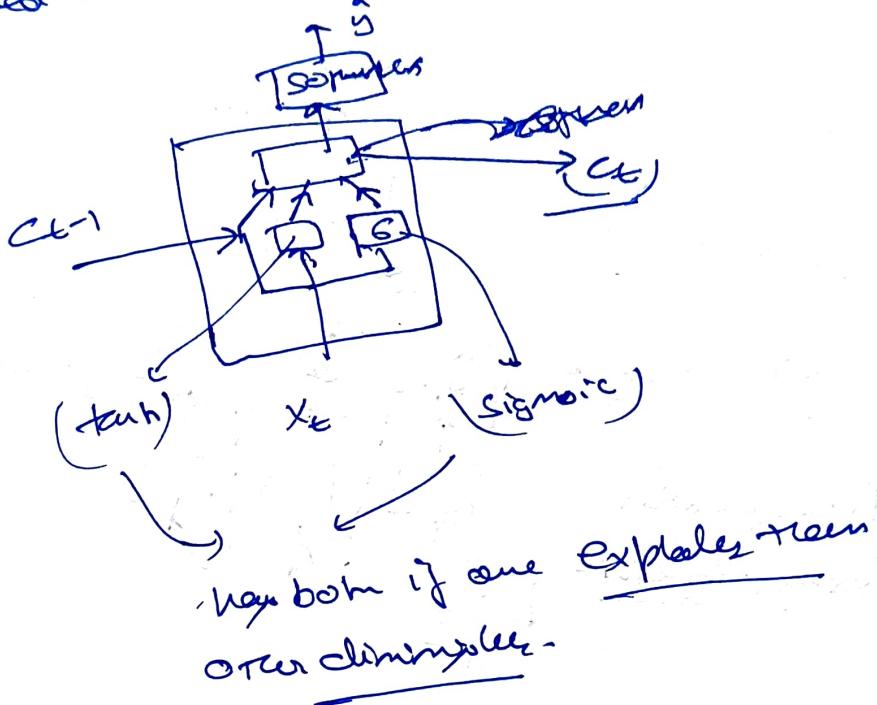
④ Many to Many ①:



Many to Many ②:

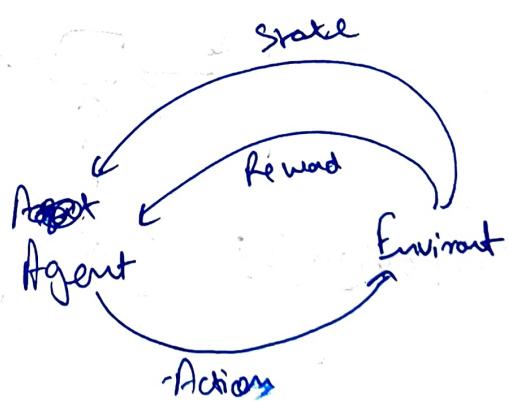


GRU (Gated Recurrent Unit)



Reinforcement Learning : (RL)

- Markov Decision Process
- learns from its experience and over time will be able to identify correct action lead to the best reward



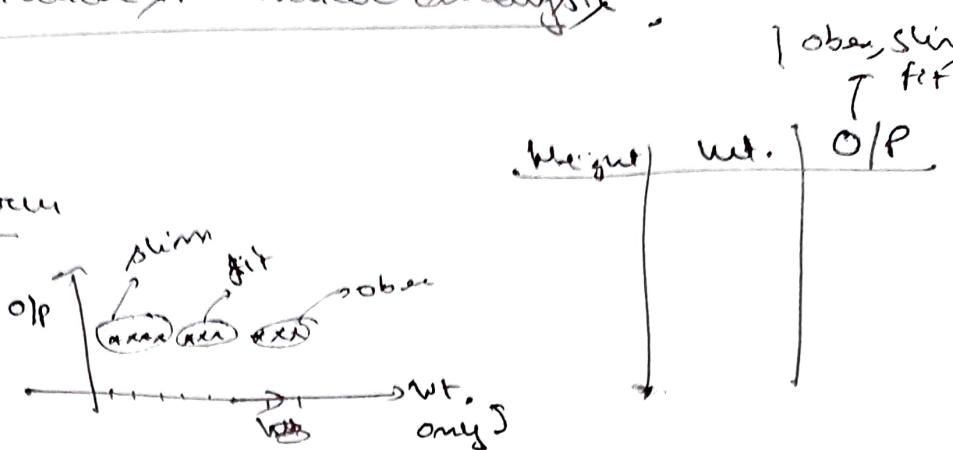
Fresh Week

Univariate, Bivariate, Multivariate analysis:

We have one feature

Univariate

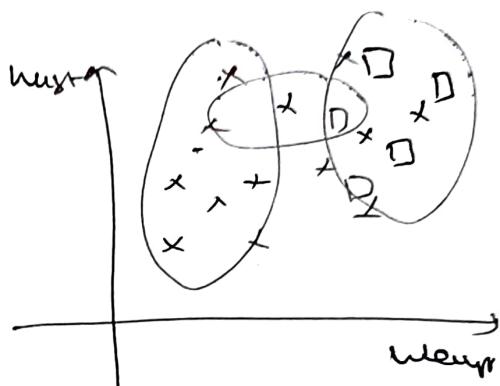
Univariate \rightarrow



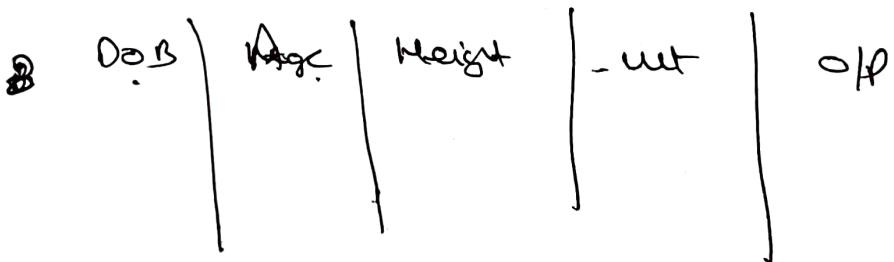
Bivariate \rightarrow Both weight, height

Non linear type

\therefore Logistic can't be used.

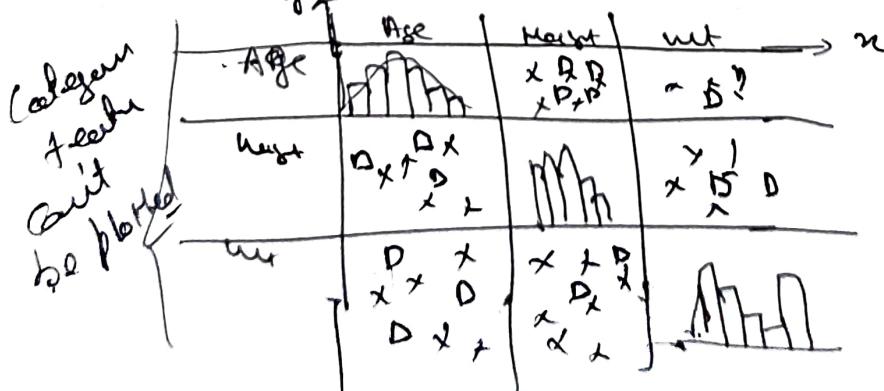


If we have 4 features: Multivariate analysis.



4 features \rightarrow DOB, Age, height, weight.

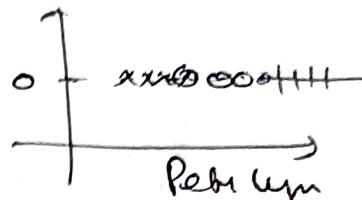
Scatter 'pair plot' \Rightarrow can be used for multivariate analysis.



P \Rightarrow Correlation

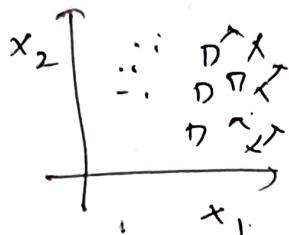
Ex: Sepal-length, Sepal-width, Petal-length, Petal-width, Species

In Classification, we use Multivariate analysis like Cen do $y=0$ for all & change color for diff. classes.



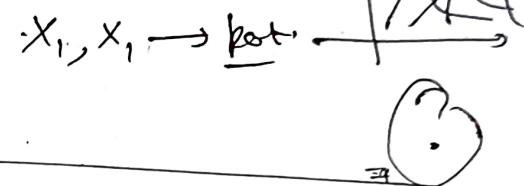
~~Bivariate Analysis~~

Sepal width \rightarrow Facet Grid ($dt, hue = "Spcl"$)

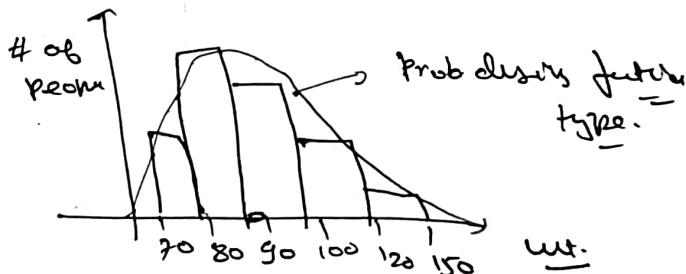
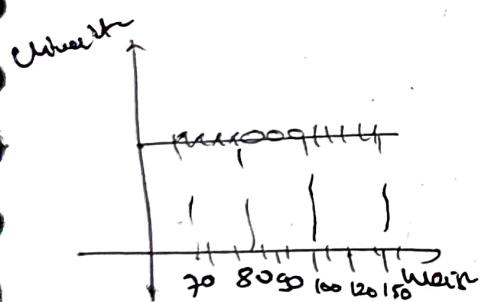


Multivariate Analysis: Pair plot

Sepal width



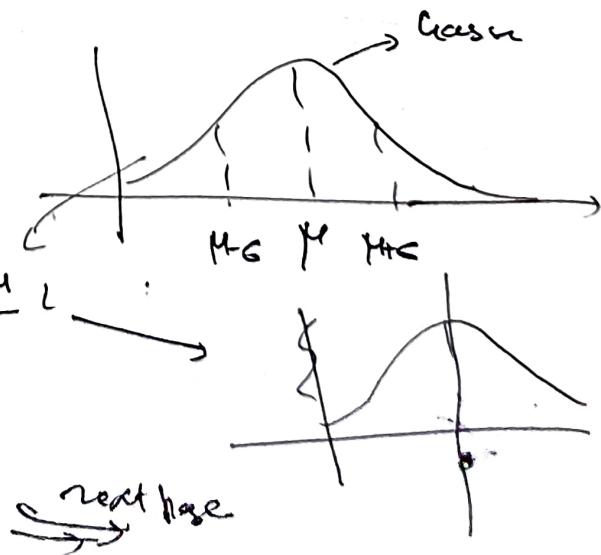
Histogram: No. of people in particular Category.



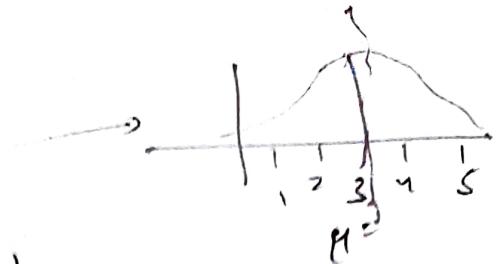
2 Score Statistics %

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

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



$$Z-\text{Score} = \frac{x_i - \mu}{\sigma}$$

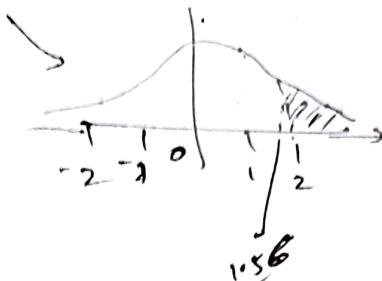


Ex. $\{1, 2, 3, 4, 5\} \rightarrow \text{mean} = \frac{1+2+3+4+5}{5} = 3$

$\sigma = \sqrt{\frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5}} = 1$ (if)

$$\text{Z-Score} = \left\{ -2, -1, 0, 1, 2 \right\}$$

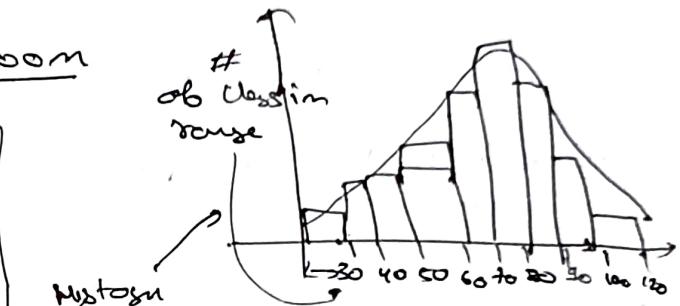
$$\frac{x_i - \mu}{\sigma}$$



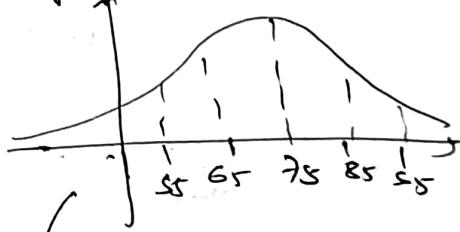
Ex: X : Population of classroom

~~Ex~~ Class No. Population (x)

1	50
2	71
3	70
4	72
5	80
6	10



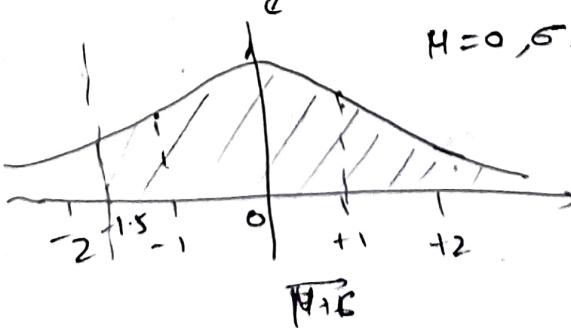
$$Pdt = \frac{\# \text{ of classroom}}{\text{total no.}}$$



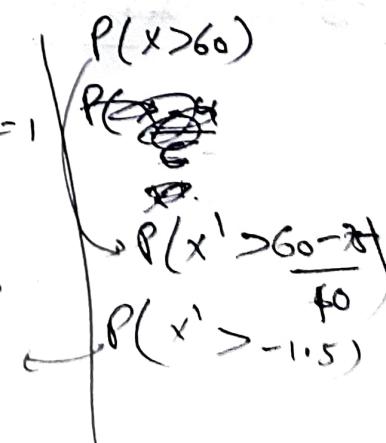
$$\mu = \frac{\sum x_i}{N}$$

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

$$\mu = 75, \sigma = 10$$



$$\mu = 0, \sigma = 1$$



$$P(X > 60)$$

$$P(X < 60)$$

$$P(X' > 60 - 1.5) = P(X' > 45)$$

$$P(X' > -1.5)$$

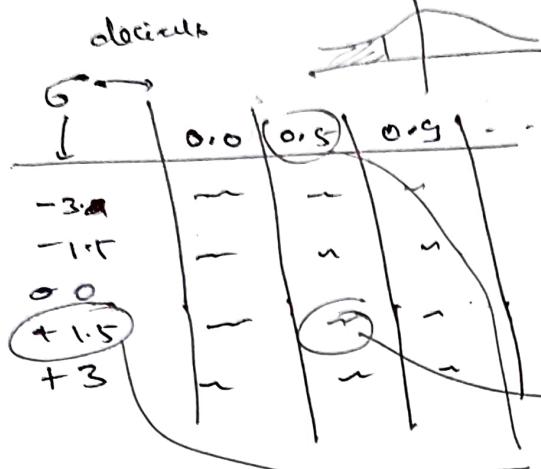
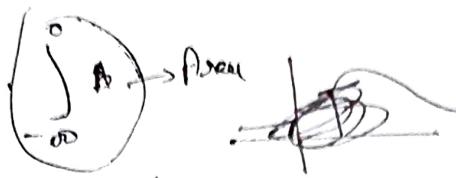
2-Sigma table $\Rightarrow P(X' \geq 1.5) = 1 - P(X' \leq -1.5)$

$$= 1 - 0.668$$

$$\approx 0.34$$

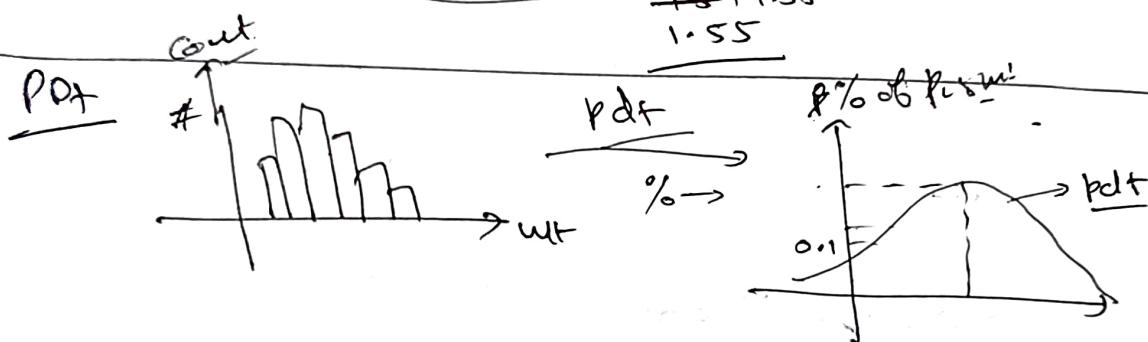
2-Sigma table

2-Sigma rule



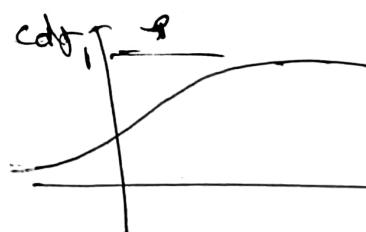
$\therefore 2.86\text{mu}$
+ 1.55 i.e.

~~+1.55~~
1.55



Regression:

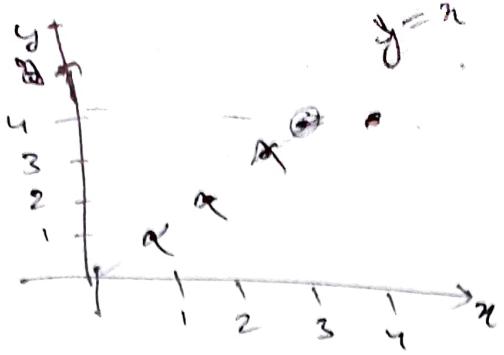
$$\hat{y} = mx + c$$



$$\text{Cost func} = \sum e_i^2 \rightarrow \text{Minimize}$$

Cost function = ~~$\frac{1}{2m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$~~

$$\frac{1}{2m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

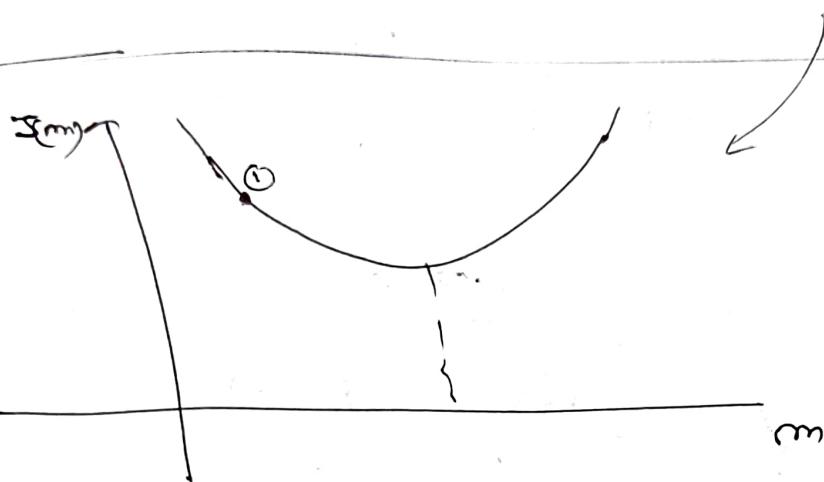
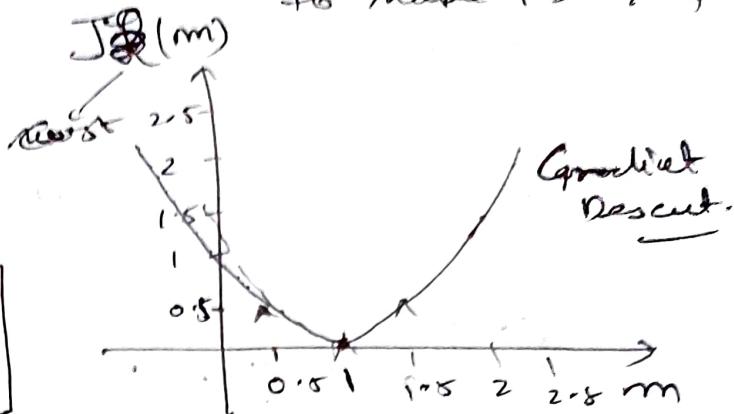


$$\hat{y} = m x + c$$

Assume $c = 0$,

to make 1D Cost fun.

$$\begin{aligned} J(m) &= \hat{y}(m) \\ J(m) &= \frac{1}{m} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \end{aligned}$$



Convex Convergence Theorem: For convex optimization.

$$m = m - \frac{\partial J}{\partial m}$$

$$m = m - \left(\frac{\partial J}{\partial m} \right) \quad \begin{array}{l} \text{Convex Rate} \\ \text{slope.} \end{array}$$

$$\frac{\partial J}{\partial m} \text{ (Convex Rate)} = \underline{\text{Small}}$$

Generalized Model \rightarrow low Bias & low Variance

Overfitting \rightarrow high variance

Underfit \rightarrow high bias

Ridge Regression

$$J = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda (\text{Slope})^2$$

High bias \rightarrow for training Data we have some amount of error
 Inc. bias but reduced variance

why?

high slope
high polarization

(large changes with x)
 \therefore over fitted.

do not go to zero

$$\text{Lasso: } J = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda |\text{Slope}|$$

\hookrightarrow helps in feature selection.

$$\text{Ex: } Y = m_1 x_1 + m_2 x_2 + m_3 x_3 + C$$

$$\geq (|m_1| + |m_2| + |m_3| + |m_C|)$$

\hookrightarrow These move toward zero

~~Multicollinearity~~

~~Pearson Correlation coefficient~~

Covariance: random variable X_1, X_2
 (Size) \quad (Price of house)

Size	Price
1200	100
1800	200
1800	300

Size \rightarrow Price

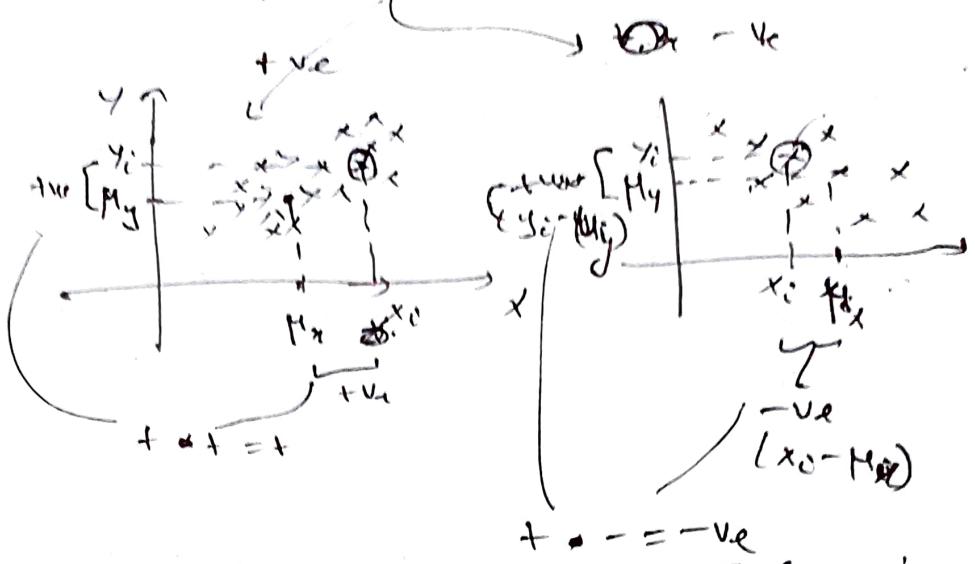
Can we quantify relation between?

$$\text{Covariance}(Size, Price) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_x) * (y_i - \bar{y}_y)$$

$$\text{Cov}(x, y) \quad \text{if } \sigma_{xx} = \frac{(x_i - \bar{x})^2}{n}$$

$$\left\{ \text{Cov}(x, x) = \text{Var}(x) \right\}$$

$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$



Pearson Correlation Coefficient:

$$PCC(P_{(x,y)}) = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y}$$

(-1 to +1)

s.d. s.d.

~~Strength~~ → Strength of Relationship
Direction

X	Y
Height	Weight

$$\text{Cov}(x,y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

$$\text{Cov}(x,y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

We get to know the.

Strength of cov(x,y)

Because here we divide

by $\sigma_x \cdot \sigma_y$.

$$\sigma_x \cdot \sigma_y = \text{S.d.}$$

$\sigma_x^2, \sigma_y^2 = \text{Variance}$

$$\sigma_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)^2$$

$$\text{Cov}(x,y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

Spearman's Rank Correlation Coefficient

(Better than r_{xy} , as Pearson checks for linear relationship but this checks for non-linear)

$$r_s(x,y) \in [-1, +1] \quad \left. \begin{array}{l} \text{Pearson of } \text{Rank}(x) \\ \text{& Rank}(y) \end{array} \right\}$$

$$\left\{ r_s = \rho_{R(x), R(y)} = \frac{\text{Cov}(R(x), R(y))}{\sigma_{R(x)} \cdot \sigma_{R(y)}} \right\}$$

$$r_s = 1 - \frac{6 \sum d_i^2}{m(m^2-1)}$$

Accordly 'Y'					
X	Y	x_{rank}	$R(x)$	$R(y)$	$d_i(d_i)$
10	50	5	1	2	-1
5	100	10	2	1	1
80	400	20	3	3	0
20	300	50	4	-1	0
			5	4	

$$r_s = 1 - \frac{6 \times 2}{4 \times 155} \quad \rightarrow \sum d_i^2 = 2$$

$$r_s \Rightarrow \{4, 5\}$$

Variance & S.D.

Bias & Variance:

overfitted \rightarrow understand like with degree of polynomial

example

$$P_m=1 \quad P_m=2 \quad P_m=n \}$$

\rightarrow overfitted.

Regression:

Underfitted

\rightarrow high Bias

\rightarrow high variance

error of training data

error of test data

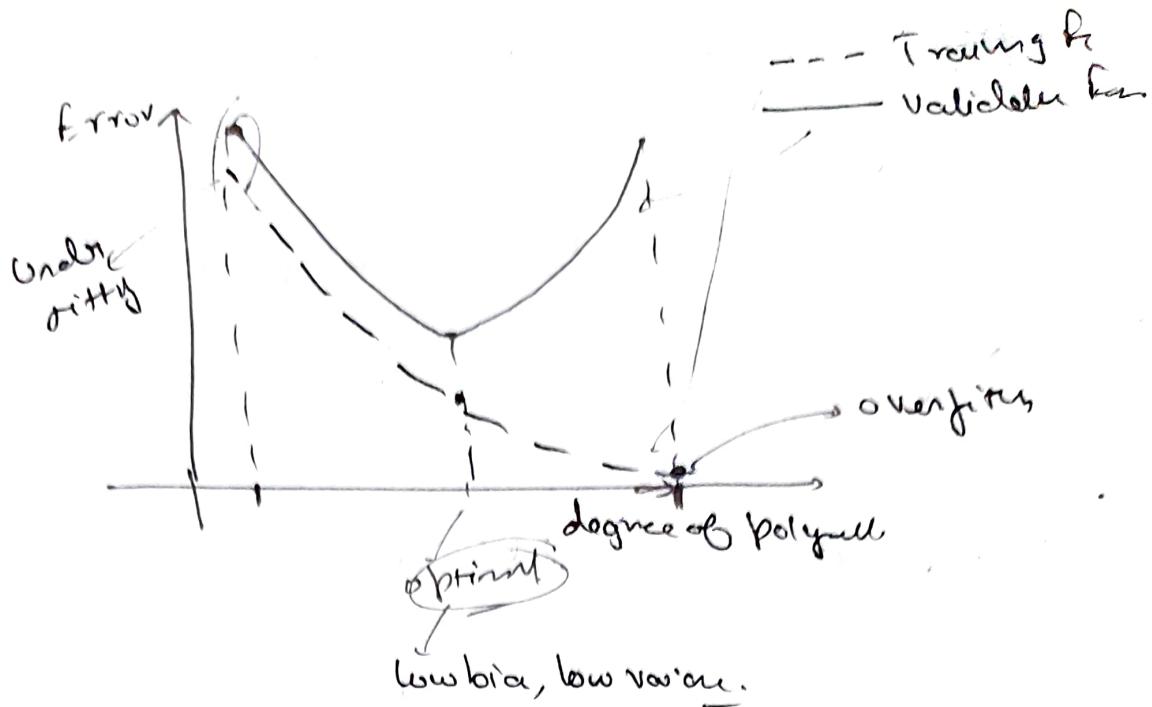
overfitted

\rightarrow low Bias

\rightarrow high variance

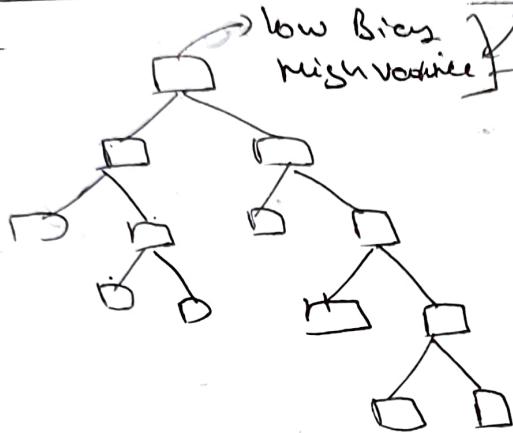
Classification

Similar



Overfit & underfit in Decision Trees & Random Forest

① DT



low Bias
high variance

→ To understand
the Tree is made
in accordance with the
training data.

→ To Avoid High variance
we do (pruning)

we do not make
DT for whole depth.

② RF → Low Bias

→ high variance

Initial

but due

to
Boosting

(Bootstrap
& Aggregating)

low variance}

as the training is splitted.
Different DT's are formed.

Standardization & Normalization

Feature Scaling \rightarrow ~~frequencies~~

Normalization \rightarrow Scales the feature b/w $(0-1)$

Standardization \rightarrow Scales features based on Standard Normal distribution $(\mu=0, \sigma=1)$

Ex: Nonzeroes

Min max normalization

$$X_{\text{norm}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}}$$

Standardize \rightarrow (standard scalar) \rightarrow learn.

Scaling can be ignored in tree based ~~and~~ models.

When to use Normalization & Standardization?

Normalisation

- \rightarrow Deep, ANN
- \rightarrow we scale down b/w $(0-1)$
- \rightarrow This helps them to learn weights quickly

Standardization

- \rightarrow In Models where we use Euclidean distance
- \rightarrow when we need some type of optimizatrices

\rightarrow Mostly performs well.

Hypothesis testing :

\rightarrow either one is true

\rightarrow we Evaluate 2 mutually Exclusive statements \rightarrow On the population \rightarrow using SAMPLE DATA

Ex: D: Defendant \rightarrow Guilty
I: Innocent \rightarrow Innocent

Steps: ① Make initial Assumption (H_0) \rightarrow Null hypothesis

Ex: H_0 : Dope test is Innocent \rightarrow

Step②: Collected Data }
 { Evidence

Ex: Fingerprint → fingerprints
→ DNA test
etc

Step③: Gather evidences to either REJECT / ACCEPT

NULL hypothesis (H_0),

or Status Quo Hypothesis
(hypothesis w/o No effect)

Alternate hypothesis: H_1 (Opposite of NULL hypothesis H_0)
(Research hypothesis)

Type I Error:

If I know, by some way that the NULL hypothesis H_0 is true but I do not have enough evidence to prove that, then H_0 will get rejected.

Confusion Matrix:

		(Truth)	
		H_0 (False)	H_1 (True)
(Test)	ACCEPT	OK	Type 2 Error
	REJECT	Type 1 Error	OK

Type 2 Error: $H_0 \Rightarrow$ Market will crash

Hypothesis is always on parameters & NEVER ON STATISTICS ($\mu, \sigma^2, \text{etc}$)

$H_1 \Rightarrow$ Not Crash
Evidence \Rightarrow Not enough \rightarrow
Evidence \Rightarrow Enough

But the H_0 gets accepted in reality i.e. Market crashes \rightarrow due to lack of evidence

is Type 2 Error

T Test, Chi Squared Test, ANOVA Test

(One categorical factor)

* One Sample Proportionality Test

- choose 1 categorical factor (Gender)
- Q: Whether there is a diff. b/w proportion of M & F? in the population

$$H_0 \Rightarrow \#M = \#F$$

$$H_1 \Rightarrow \#M \neq \#F$$

Gender	Age Group	Weight (kg)	Height (cm)
M	Elderly	70	144
F	Adult	65	122
M	Adult	65	144
M	Child	20	111
F	Adult	75	133
M	Elderly	80	133

Test: Considering H_0 is true
what is the prob. that H_1 is true?

Sample

$$\text{if } P \leq 0.05 \rightarrow \text{Significance level} (\alpha)$$

→ Reject H_0 & Accept H_1 . \rightarrow Chi Sq. test

* Two Categorical factors [Gender, Age group]

- Q: Is there diff. b/w #M & #F based on Age group?

$$H_0 \Rightarrow M = F \text{ acc to Age group}$$

$$H_1 \Rightarrow M \neq F$$

Test \rightarrow Chi Sq. $\rightarrow P \leq 0.05 \rightarrow$ reject $\circlearrowleft H_0$

T-test

$H_0: \text{No relation } \rho = 0$

$H_1: \rho \neq 0$

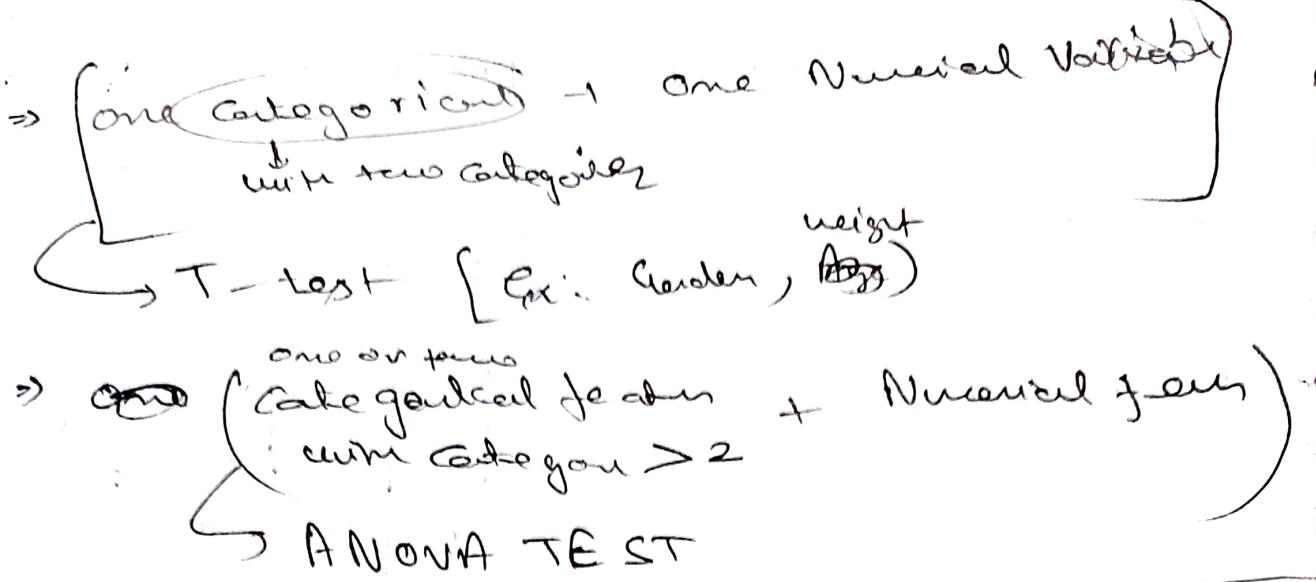
* One Continuous Variable of Weight

→ we only talk about sample statistics.

- Q: Based on ~~some~~ ^{some} ~~the~~ sample is there a diff. in mean weight of two samples

$$H_0: \mu_A = \mu_B \quad H_1: \mu_A \neq \mu_B$$

Test: T-test



#. T-test : Type of Inferential Statistics
 which is used to determine if there is a significant difference b/w the means of two groups which may be related in certain features

T-test has 2 types

- ① One-Sampled t-test
- ② Two-Sampled t-test

① One-Sampled T-test : Tells us whether means of the sample and the population are different.

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \quad \text{where} \quad s_{\bar{x}} = \frac{s}{\sqrt{n}}$$

$H_0 \Rightarrow$ No difference.

$H_1 \Rightarrow$ Difference

if $P \leq 0.05 \rightarrow$ Reject H_0
 Accept H_1

$\mu \Rightarrow$ Proposed Constant for population mean.
$\bar{x} \Rightarrow$ Sample mean
$n \Rightarrow$ Sample size (# of observations)
$s \Rightarrow$ Sample Standard deviation
$s_{\bar{x}} \Rightarrow$ Estimated S.D. of the mean

→ Two Sample T-test: The independent samples t-test or 2-sample t-test compares the means of two independent groups in order to determine whether there is statistical evidence that the associated population means are significantly different. The independent sample t-test is a parametric test. This test is also known as: ~~as~~ independent t-test.

$$\left\{ \begin{array}{l} t = \frac{x_2 - x_1}{\sqrt{s^2 \left(\frac{1}{m_1} + \frac{1}{m_2} \right)}} \\ ; \quad s^2 = \frac{\sum_{i=1}^{m_1} (x_i - \bar{x}_1)^2 + \sum_{j=1}^{m_2} (x_j - \bar{x}_2)^2}{m_1 + m_2 - 2} \end{array} \right\}$$

→ Paired T-test: To check how different samples from the same group are.

$$W_{t-1} = \{ \quad \text{↳ weight of students this year}$$

$$W_{t-2} = W_{t-1} + \{ \quad \text{↳ added distances from a} \\ \text{seed } N(1.5, \sigma) \quad \text{↳ mean} = 1.8 \text{ m}$$

Now. Are the W_{t-1} & W_{t-2} → statistically different.

$$\begin{aligned} H_0 &\Rightarrow \text{Same} \\ H_1 &\Rightarrow \text{different} \end{aligned}$$

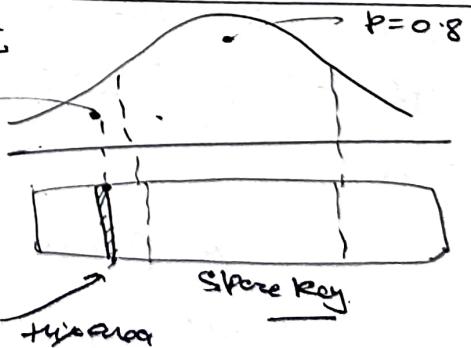
Correlation Test:

P-value:

$$\{ P = 0.01 \leftarrow$$

↳ if we repeat this experiment of pressing the space key the no of times we will be pressing is 1.

2 tailed test

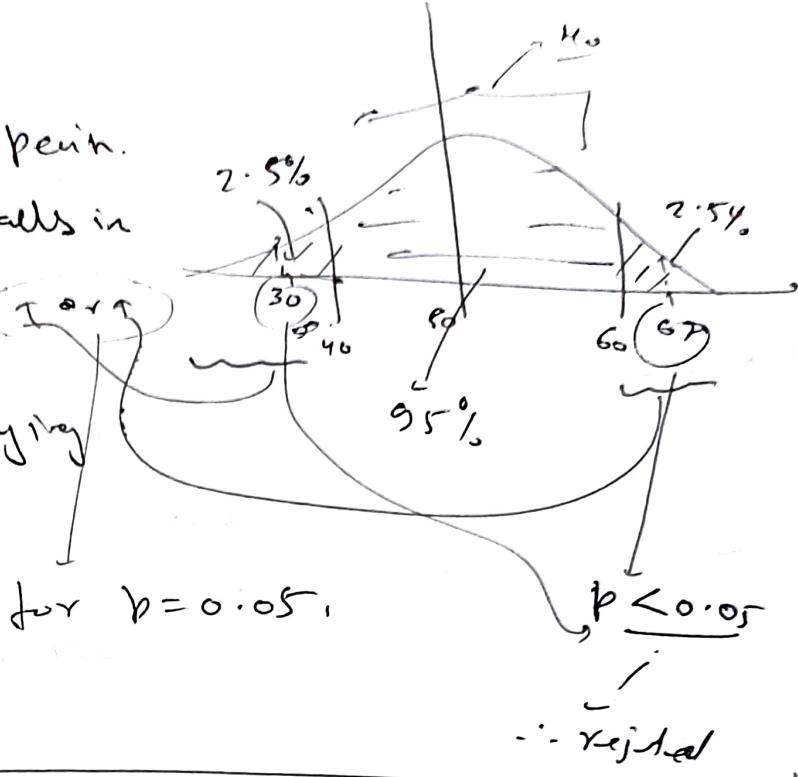


P-Value \Rightarrow It is the probability for the NULL hypothesis to be true.

H_0 (NULL hypothesis) \rightarrow Treats Everything Same or Equal.

Ex:

If I do an experiment & my outcome falls in this region (I am extremely reject my H_0 (NULL) saying that



Multicollinearity:

In linear regression

\Rightarrow If our data has multicollinearity issue then we will have a large "Standard Error" for the variables.

\rightarrow Check the p-value is > 0.05 then see co-relation plot, or drop the variable with higher P-value

OLS \rightarrow Ordinary Least Squares
 \rightarrow Total variation in $y = c + \beta_1 x_1 + \beta_2 x_2$

R-square & Adjusted R-squared

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{total}}}$$

$$R^2 \in [0-1]$$

$$R^2 \rightarrow 1 \quad \{ \text{good model} \}$$

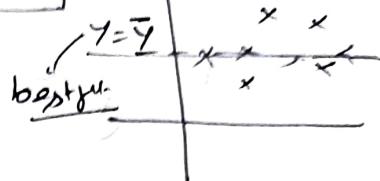
$SS_{\text{res}} = \text{sum of residuals}$

or
error.

$$SS_{\text{res}} = \sum (y_i - \hat{y}_i)^2$$

$SS_{\text{total}} = \text{sum of average total}$

$$SS_{\text{total}} = \sum (y_i - \bar{y})^2$$



Q: Can the $R^2 < 0$?

Yes

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{total}}}$$

$\therefore SS_{\text{res}} > SS_{\text{total}}$

Can happen even the line fitted is worsenthan the $y = \bar{y}$ line.

R^2 checks Goodness of best fit line

Adjusted R²

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 -$$

(Penalizes features that are uncorrelated)

as the no. of features inc. the

$$(R^2 \rightarrow \text{inc.})$$

why? Intuitive

There's a problem
that is the new feature added has
no corr with the

y (output) then
also due to inc. in
features the R^2 inc.

to avoid this we use Adjusted R^2

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

$R^2 \Rightarrow \text{sample } R^2$

$p \equiv \# \text{ of Predictors}$

$N \equiv \text{Total sample Size}$

$$\text{as } (y_i - \bar{y})$$

$$\frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

$$\therefore R^2 \uparrow$$

due to 6th feature
optimization
by L.R.

inc. shows
of R^2 (not
check if
check if
 $\beta_1, R^2_{\text{adj}} \downarrow$)
because total
when feature are not
(or related)

If feature added is correlated to
added then R^2 will be very inc.
then $(1-R^2) \rightarrow$ slow. : (R^2_{adj}) more

Dif. b/w R^2 & R^2_{adj}

→ Every time we add an independent variable to a model, the R^2 -sq. inc., even if independent variable is insignificant. It never decreases. Whereas Adjusted R^2 -sq. inc. only when independent variable significantly affects dependent variable.

$$\Rightarrow R^2_{adj} \leq R^2$$

Chi-squared test → Applied when we have two Categorical Variables from a single population.

→ It's used to determine whether there is a significant association b/w two variables.

{ Sometime like correlation of categorical variables }

Ex: $X_1 = \text{Gender } \{M, F\}$

$X_2 = \text{Smoker } \{\text{YES}, \text{NO}\}$

		= Smoker		⇒ given data (observed)
		YES	NO	
Gender	M	60	97	Expected
	F	33	54	

$$\{(Dof)_{(\text{Degree})} = (\# \text{ rows} - 1) \times (\# \text{ cols} - 1)\}$$

Not completely understood:

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

Metric: $\begin{cases} \text{Chi-Sq. statistic} \\ \text{OR P-value} \end{cases}$

Imbalanced $\Rightarrow \{70 - 30\%\}$

Confusion

		Action	
		1	0
Predicted	1	TP	FP
	0	FN	TN

Type 1 \rightarrow False Pos Rate
 $FPR = \frac{FP}{FP + TN}$

Type 2 Error
 $FNR = \frac{FN}{TP + FN}$

Total Predicted

Recall $\Rightarrow \frac{TP}{TP + FN}$
 or
 (TPR)
 or
Sensitivity

Precision $\Rightarrow \frac{TP}{TP + FP}$
 } Positive Predictive Value

} Out of the total +ve Predictions
 how many were correct

Ex: Spam 0 \rightarrow Not Spam
 1 \rightarrow Spam

Precision $\Rightarrow \frac{TP}{TP + FP}$ reduce +ve.

} Not a spam but classified as spam

Ex:

		Action	
		1	0
Predicted	1	TP	FP
	0	FN	TN

Recall

Show \rightarrow Not.
~~0~~ 1 \rightarrow 0

Precision
 $\frac{TP}{TP + FP}$ Show

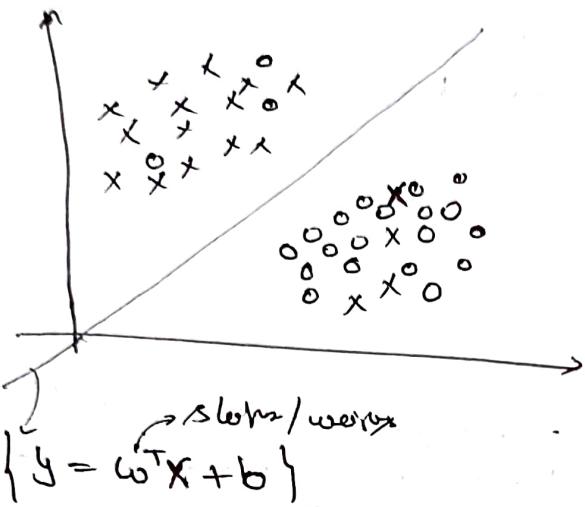
		Action	
		1	0
Predicted	1	TP	FP
	0	FN	TN

1 \rightarrow N
 0 \rightarrow O
 0 \rightarrow 1

Logistic Regression:

(fixes Binary classification)

→ for linearly separable classes.



Let's take as only one feature x_1

$$y = mx + c$$

Coefficient is not calculated here

that in Linear regression case

Assume

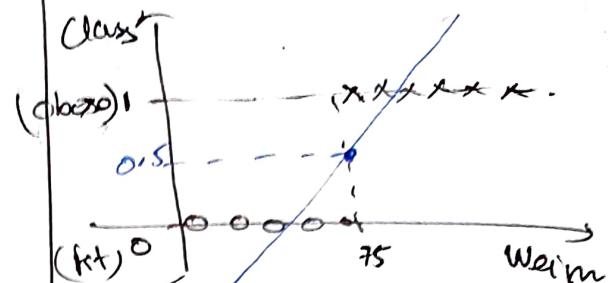
$$\begin{cases} x \rightarrow +1 \\ 0 \rightarrow -1 \end{cases} \quad \text{take } c=0$$

$$y_i = w^T x_i + b$$

$$\begin{cases} y_1 = 1 \\ y_0 = -1 \end{cases} \quad d_i = w^T x_i$$

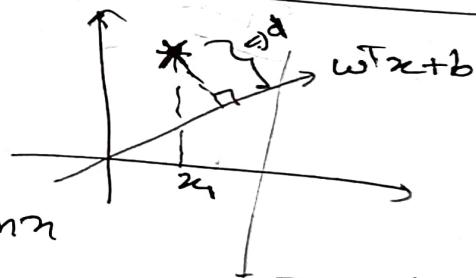
$$d = \sum_{i=1}^n d_i = \sum_{i=1}^n w^T x_i$$

Why Not do we use
Binary classifier using
Linear Regression?



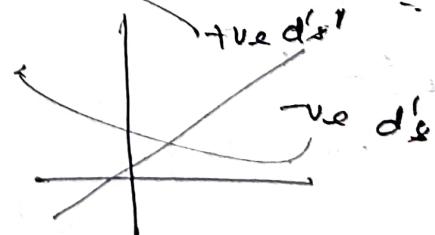
Regression
if $y > 0.5 \rightarrow \text{class 1}$
if $y < 0.5 \rightarrow \text{class 0}$

→ But due to one huge outlier like $w^T = 1000b$
the regression line may be affected badly
∴ we do not use linear in classification



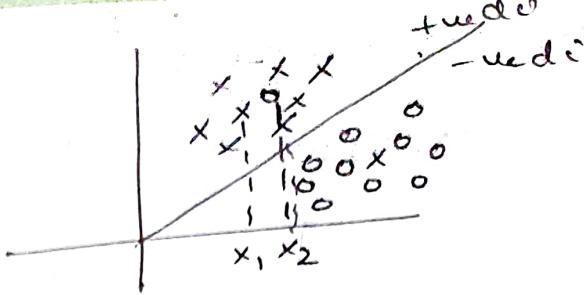
$$d = \frac{|w^T x_i + b|}{\|w\|}$$

Assume unit vector
 $\therefore \|w\| = 1$



$$y_{+} = +1$$

$$y_{-} = -1$$



d_i oben $\Rightarrow +$

d_i unten $\Rightarrow -ve$

Case ① \rightarrow Correct classification
if we have 'x' above x_1 , then $\Rightarrow y_i = +1$

$$d_i = +ve$$

Case ② \rightarrow wrong classification

If we have 'o' at x_2 ,

$$y_i = -1$$

$$d_i = +1$$

$$y_i \cdot d_i = y_i \cdot w^T x_2 < 0$$

wrongly classified

$$\begin{aligned} y_i \cdot d_i &> 0 \\ y_i \cdot w^T x_2 &> 0 \end{aligned}$$

Cost function: Max.
in logistic regression

$$\sum_{i=1}^n (y_i w^T x_i)$$

This should be Maximum

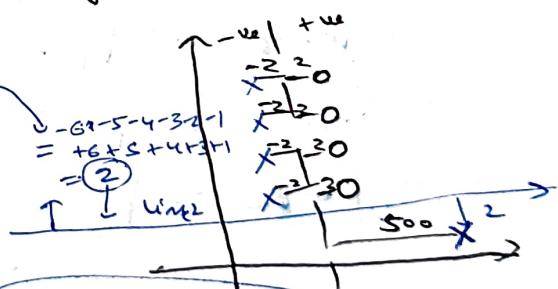
This can be optimized to get maximum Cost, for best fit.

Dealing with outliers in Logistic Regression

$$\sum_{i=1}^n y_i w^T x_i$$

$$\begin{aligned} &\Rightarrow 2+2+2+2 \\ &-(-1+2+2+2) \\ &-500 \end{aligned}$$

$$\Rightarrow 20 - 500 = -480$$



: we will skip this line much is a best fit
 $x \rightarrow +1$
 $0 \rightarrow -1$

\therefore line(2) gets selected with value (+2) rather than line1 with value (-480).

How to overcome this?

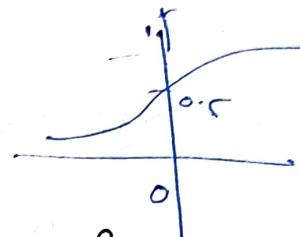
we apply sigmoid function.

Cost function =

$$\sum_{i=1}^m \text{Sig}(y_i + w^T x_i)$$

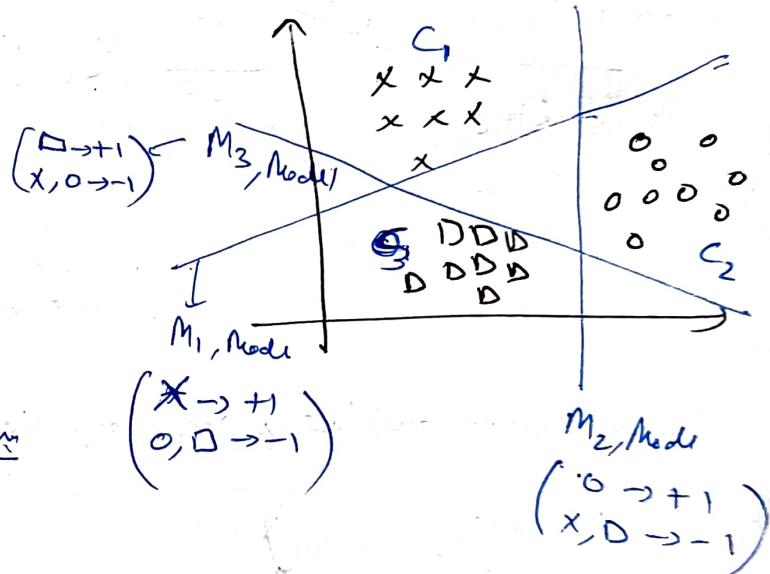
removes the effect of outliers by squashing the value. b/w 0 to 1

$$f(z) = \frac{1}{1+e^{-z}}$$



Multiclass classification using "Logistic regression"

(One Vs All)



Now if I take
a new point (x_{new})
Probability
 $M_1, \text{Model} \rightarrow p_1$
 $M_2, \text{Model} \rightarrow p_2$
 $M_3, \text{Model} \rightarrow p_3$

("rr") \rightarrow sklearn

take highest probability
i.e My point belongs to
C3 category i.e

ROC & AUC

y	y_{prob}	\hat{T}_{p_1}	T_{p_2}
1	0.8	1	1
0	0.36	1	1
1	0.4	1	1
1	0.3	1	1
0	0.2	1	0
1	0.7	1	1

→ Threshold Problem
 $\hat{T}_p = \{0, 0.3, 0.4, 0.6, 0.8, 1\}$

Actual

1	TP	FP
0	FN	TN

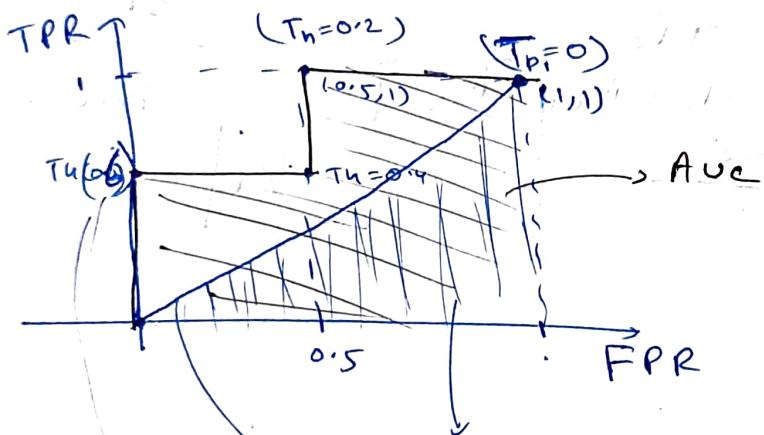
$$\begin{aligned} TPR &= \frac{TP}{TP+FN} \\ FPR &= \frac{FP}{FP+TN} \end{aligned}$$

TPR, FPR

$$TPR = \frac{TP}{TP+FN} = \frac{4}{4+0} = 1$$

$$FPR = \frac{FP}{FP+TN} = \frac{2}{2+0} = 1$$

$$\begin{aligned} TPR &= 4/4+0 = 1 \\ FPR &= 2/2+0 = 0.5 \end{aligned}$$



Model
not + model
AUC → Area
under diag line
except a dumb model
with $p=0.5$

- Now, if problem wants higher TPR (lower the FN (like in cancer case)) ; we can only focus on TPR and can take $T_h = 0.6$ as then $FPR = 0$
- But if I want high TPR & don't care about FPR, Then we can take $T_h = 0.2$ (high TPR)

Ex:

Performance Metric for multiclass classifier

$$\left\{ \text{Accuracy} \Rightarrow \frac{TP_0 + TP_1 + TP_2}{\text{Total test cases}} \right\}$$

		Action		
		0	1	2
Model	0	TP ₀		
	1		TP ₁	
	2			TP ₂

Precision & Recall

O-Cat
1 → Dog
2 → Fox

Precision:

$$P_{\text{new}} = \frac{TP + FP}{TP + FP}$$

0	FN	TN
TP	FP	

→ Precision for class '0'-Cat is
the correctly predicted Cat out
of all Predicted Cat

$$\left\{ \text{Precision(Cat)} = \frac{TP_0}{TP_0 + M_{01} + M_{02}} \right\}$$

similar for other.

similar for Recall.

Inspection

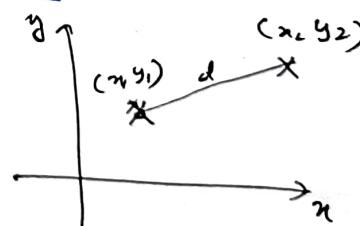
Precision.

Actu

Euclidean & Manhattan distance:

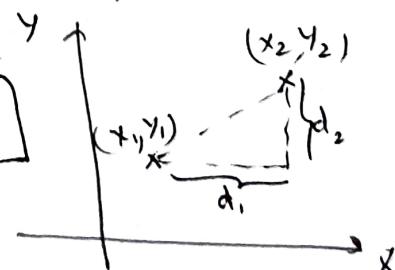
$$\text{Euclidean dist} \Rightarrow \sqrt{(y_2 - y_1)^2 + (x_2 - x_1)^2}$$

(by Pythagoras)

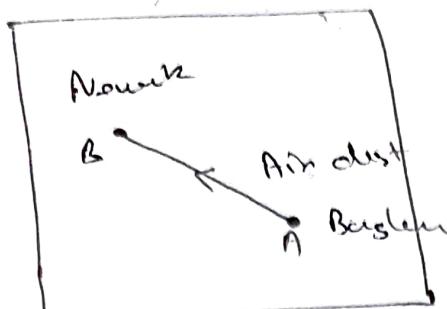


$$\text{Manhattan dist} = d_1 + d_2$$

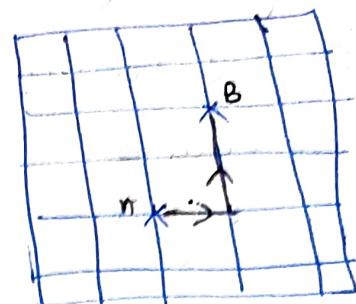
$$= [(x_2 - x_1) + (y_2 - y_1)]$$



Where to use Euclidean & Manhattan distance?



Euclidean dist case.



City S

dist b/w A & B use Manhattan dist.

(Euclidean & Manhattan)

are key parameters. While trying we can use both and see which work well.

Types of Cross Validation

① Leave One Out Cross Validation (LOOCV)

1000 Records



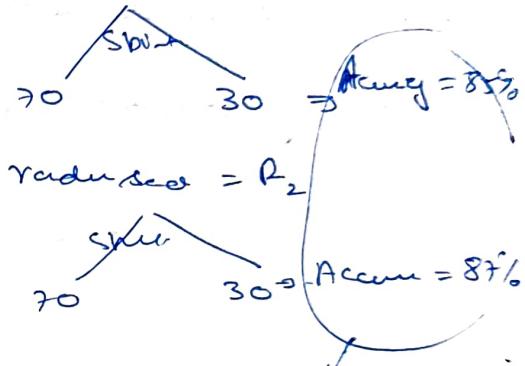
Test 1 → 1 Test, rest Train

Test 2 → 1 Test, rest Train

→ High Computation required

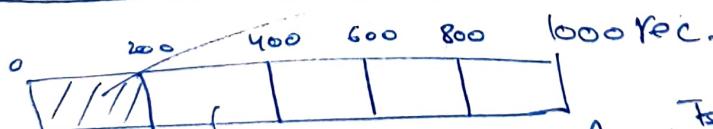
→ Lead to Low Bias

1000 Records
Use random state = R₁



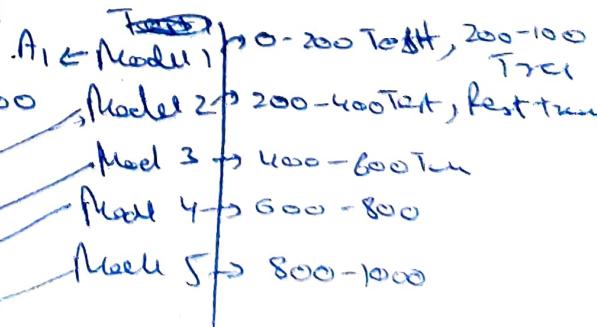
∴ use split to into Train Test + validate

② K-fold Cross Validation:



K=5 → 5 experiments →

$$\frac{1000}{5} = 200$$



$$\text{Mean Acc} = \frac{A_1 + A_2 + A_3 + A_4 + A_5}{5} \quad (\text{Acc}_1, \dots, \text{Acc}_5)$$

$$\text{Min Acc} = \min(A_1, A_2, \dots, A_5)$$

$$\text{Max acc} = \max(A_1, \dots, A_5)$$

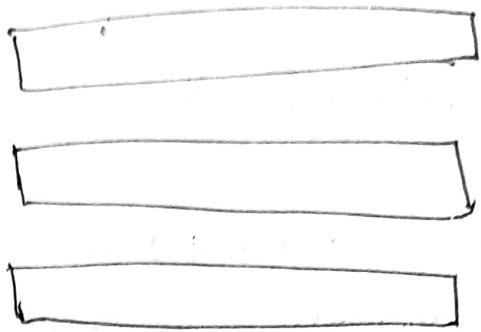
→ In R-fold Cross Validation. If the data is Imbalanced then there can be some problems of Biased learning.

∴ we have →

③ Stratified Cross Validation {For balanced data}

If $k=3$, 1200 records.

3 Models with 400 Test,



out of the 400 test, stratified CV makes sure that both Train & Test splits contain same proportion

(like if data has $\frac{60\%}{40\%}$ class imbalance,

then Test \rightarrow 60% 40% class

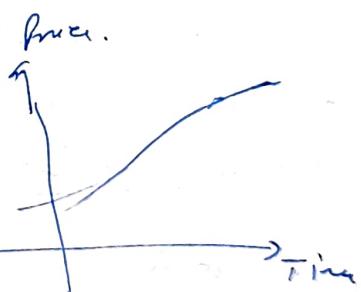
Train \Rightarrow //

④ Time Series Cross Validation's

- B Day 1
- D2
- D3
- D4
- D5

D1 D2 D3 D4 D5	D6
D2 D3 D4 D5 D6	D7

Predict [D6
D7]



BAYES' Theorem:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

↑
Posterior
Probability

↑
Prior
of
A

↑
Likelihood
Prob

⇒ Marginal
Prob

Naive Bayes Classifier:

→

$$P(Y/x_1, x_2, x_3, x_4) = \frac{P(Y) \prod_{i=1}^4 P(x_i|y)}{P(x_1) P(x_2) P(x_3) P(x_4)}$$

(0 - 1) →

$$P(Y/x_1, x_2, x_3, x_4) \propto P(Y) \prod_{i=1}^4 P(x_i|y)$$

$$P(Y/x_1, x_2, x_3, x_4) = \arg \max (P(Y) \prod_{i=1}^m P(x_i|y))$$

Ex: Naive Bayes → Yes / No ; Good/Bad Classification

→ Baseline for Text Classifier (NLP) ①

Scatterer → 4000
Bad

After preproc

we get Vector of words ⇒ much
of 1000 features

d ₁	d ₂	d ₃	d ₄	0/P	or Bad
True	Food	Delish	Bad	0	1
1	1	1	0	1	0
1	1	0	1	0	1

ex: bestest Review

$$P(Y=Good / \underbrace{\text{Scatterer}}_{[x_1, \dots, x_m]}) = P(Y=Good / (x_1, \dots, x_m))$$

$$\Rightarrow \frac{P(y) \cdot P(x_1|y) \cdot P(x_2|y) \cdot P(x_3|y) \cdots P(x_m|y)}{P(x_1) \cdot P(x_2) \cdot P(x_3) \cdots P(x_m)} \Rightarrow$$

$$= \arg \max \left(P(y) \cdot \prod_{i=1}^m P(x_i | y=y_{\text{yes}}) \right)$$

Last table $P(y_{\text{yes}}) = \frac{1}{3}$

$$P(x_1 = \text{the} | y=\text{yes}) = \frac{1}{2}$$

$$P(x_2 = \text{bad} | y=\text{yes}) = \frac{1}{3}$$

$$P(x_3 = \text{delicious} | y=\text{yes}) = \frac{1}{3}$$

$$P(x_4 = \text{bad} | y=\text{yes}) = \frac{0}{2} = 0$$

So output \rightarrow The food is delicious

$$P(y=\text{yes} | \text{sent}) = \frac{P(y=\text{yes}) \cdot P(\text{the} | y=\text{yes}) \cdot P(\text{bad} | y=\text{yes}) \cdot P(\text{delicious} | y=\text{yes})}{P(\text{the}) \cdot P(\text{bad}) \cdot P(\text{delicious})}$$

$$= \frac{\frac{1}{3} \times \frac{1}{2} \times \frac{1}{3} \times \frac{1}{3}}{1} = \frac{1}{18}$$

$$P(y=\text{No} | \text{sent}) = \frac{1}{18} = \frac{1}{18}$$

calculator similarly

$$\therefore P(1 - P(y=\text{yes} | \text{sent}))$$

\Rightarrow Naive Bayes fails in case when the sentence ~~contains~~

\Leftrightarrow Vector has words that are not in feature

\therefore Many $P(\text{word} | y) = 0$

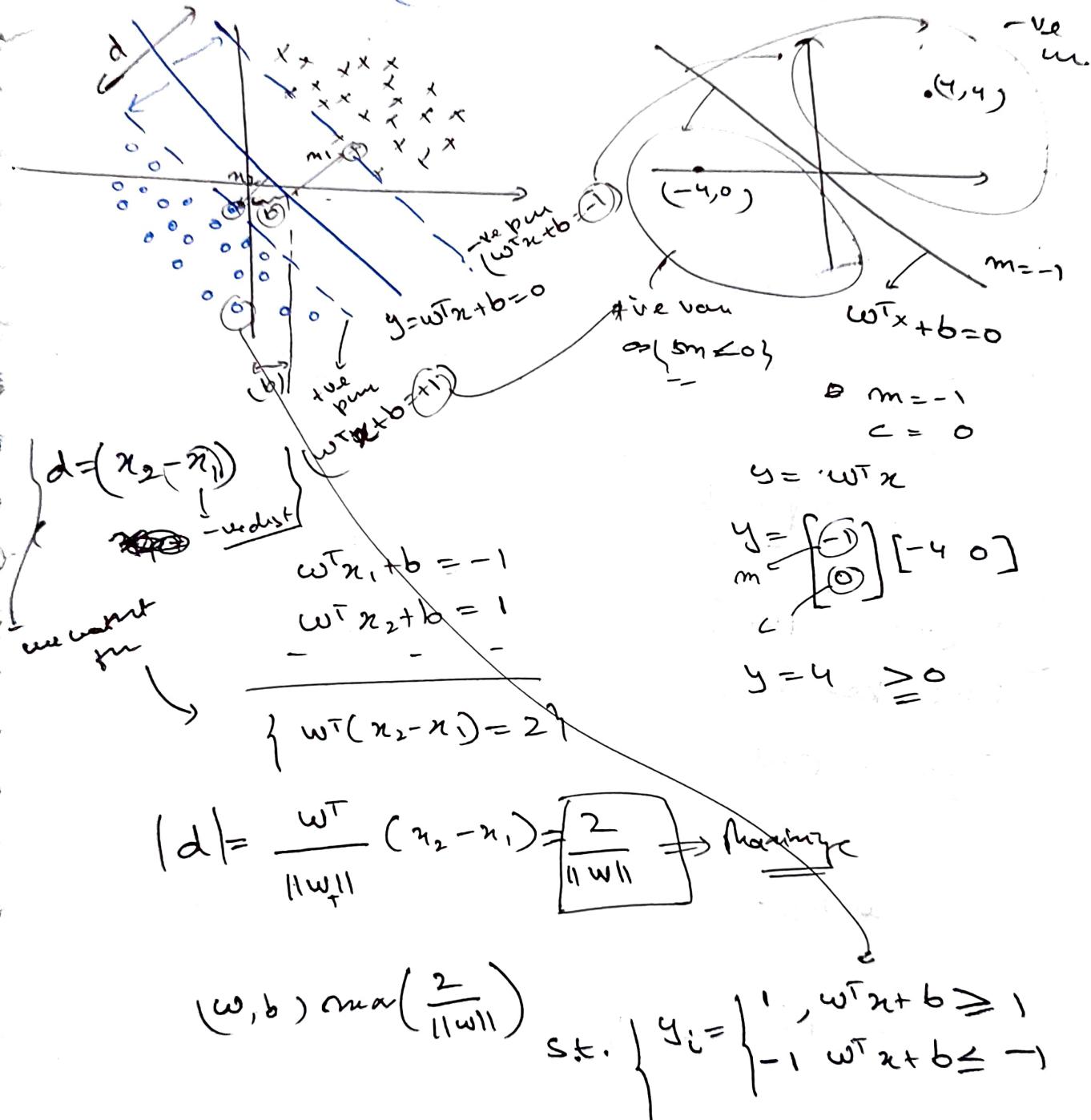
\therefore giving tree $(-\text{ve})$ output always

Support Vector Machines (SVM)

→ Both for Regression & Classification

Maximize \rightarrow Marginal distance.

SVM Kernel \rightarrow (Low dimension to high dimension)



$\xrightarrow{\text{Carried}} \boxed{y_i \cdot (\mathbf{w}^\top \mathbf{x}_i + b_i) \geq 1}$

Correctly classified.

Finally

$$(\mathbf{w}^*, b^*) = \min_{\mathbf{w}, b} \frac{\|\mathbf{w}\|}{2} + C \sum_{i=1}^m \max_{\epsilon_i} \frac{2}{\|\mathbf{w}\|} \epsilon_i$$

How many
error
my model can
take

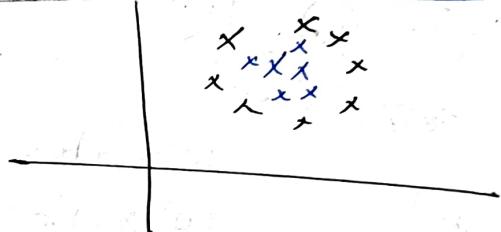
error
sum
(dist-front line)

Regularizations

- ① Soft margin
 - ② Hard margin
- frequency

Non linear separable

- ① Polynomial Kernel
- ② RBF Kernel
- ③ Sigmoid Kernel.



ROC curves:

Box 6x knist Naiz: