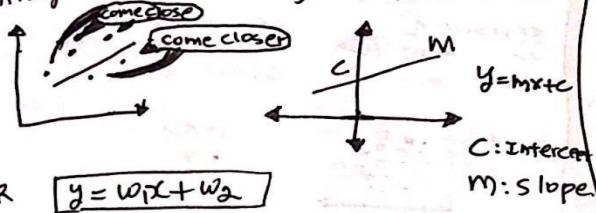


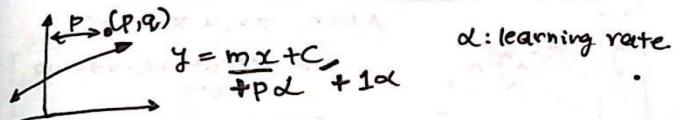
Supervised Learning:

TERM-1

Fitting a line through Data:



Methods: ① Absolute Trick:

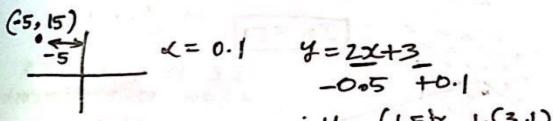


∴ New equation: $y = (m + pd)x + (c + \alpha)$

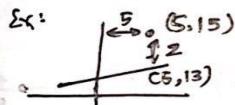
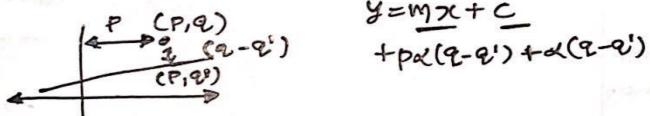
But what if point is Below the line?

→ Instead of Adding we Subtract \ominus

$$\therefore y = (m - pd)x + (c - \alpha)$$

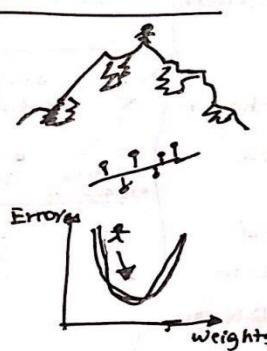


② Square Trick:



Gradient Descent:

$$w_i \rightarrow w_i - \alpha \frac{\partial}{\partial w_i} \text{Error}$$



Error Functions: (For linear Regression)

① Mean Absolute Error

$$\text{Error} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

(x_i, y_i)

(\hat{x}_i, \hat{y}_i)

$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$ NO. of points

② Mean Squared Error

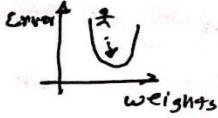
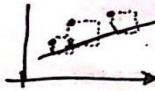
$$\text{Error} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

(x_i, y_i)

(\hat{x}_i, \hat{y}_i)

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

Ex:



Batch & Mini-Batch & Stochastic:

• Come close MSE by All points \odot



① Batch

Come closer!!

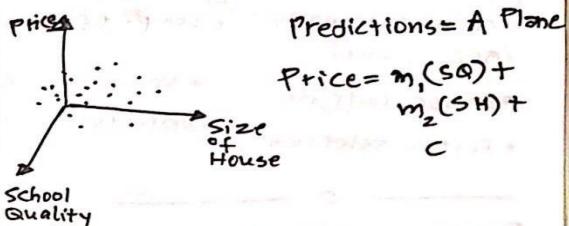
② Stochastic.

Come closer!

③ batch.

Higher Dimensions:

Price



$$y = m_1x_1 + m_2x_2 + m_3x_3 + \dots + m_nx_n + C$$

Linear Regression Works best if the data is Linear

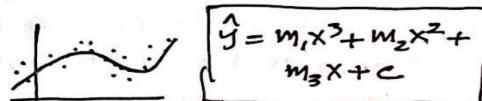


• If outliers are present; \rightarrow

Dramatic changes in Model's behavior

∴ Polynomial comes In!!!

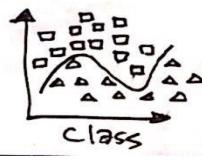
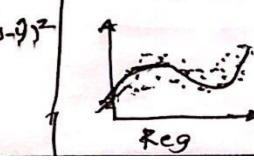
Polynomial Regression:



Polynomial can be used to both

1. Classification

2. Regression



Simple or Complex Model?

Depends on Application:



Requires low error;
OK with its complexity



Requires simplicity
and fast.

λ -Parameter:

Miss classified X Added Error \Rightarrow
Complex Model Wins

L1 & L2 Regularization:

L1 \rightarrow Absolute value

L2 \rightarrow Squared value.

L1

- computationally inefficient.

- Sparse outputs

- Feature selection

L2

- comp-efficient

- Non-Sparse O/Ps.

- No feature selection

Feature Scaling:

FS is a way of classifying data into a common range of values. 2 common scaling:

1. Standardizing

2. Normalizing

1. Standardizing: is taking each value of the column and divide by Std of the column

2. Normalization: Data is scaled between 0 and 1.

When to use FS?

In Most ML algo the result will change depending on the units of the data. This happens in 2 cases:

1. When algo uses distance based metrics to predict ex: Euc_2 , SUM etc.

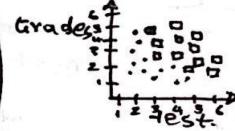
2. When you incorporate regularization ex:

Classification:

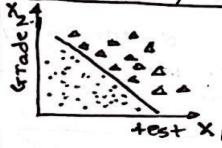
Example: Student Acceptance to college:

Student 1
test: 9/10 \bigcirc
Grade: 8/10

Student 2
test: 3/10 \times
Grade: 4/10



A Boundary line:



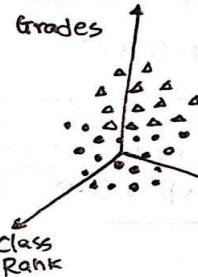
$$\text{Line: } 2x_1 + x_2 - 18 = 0 \\ \text{i.e. } 2(\text{test}) + \text{grade} - 18 = \text{score}$$

Prediction:

$$m_1 x_1 + m_2 x_2 = 0$$

Score > 0 : Accept
Score < 0 : Reject

What if there are 3 features?



IN 3D

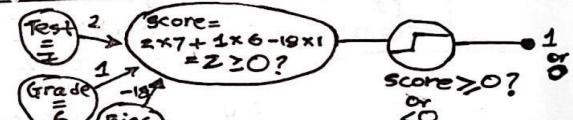
Now Boundary is
Not line But
a Plane

If n features
then $n-1$ dimension
hyperplane

$$\begin{aligned} \text{Line: } & m_1 x_1 + m_2 x_2 + m_3 x_3 + b = 0 \\ & + m_n x_n + b = 0 \end{aligned}$$

$$\text{Predictions: } y = \begin{cases} 1 & ; Wx+b \geq 0 \\ 0 & ; Wx+b < 0 \end{cases}$$

Perceptron:



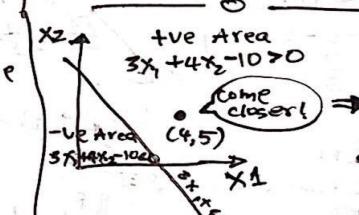
$$\text{Score} = 2 * \text{test} + 1 * \text{grades} - 18$$

Now take small steps \rightarrow use learning rate: 0.1

$$\begin{aligned} & \begin{matrix} 3 & 4 & -10 \\ -4 & 5 & 1 \end{matrix} \\ & \therefore \begin{matrix} 3 & 4 & -10 \\ -0.4 & -0.5 & -0.1 \end{matrix} \end{aligned}$$

A line $\text{eq} \approx$ which will be closer to point

If miss classified point is in Area 0
Add \oplus
If in 1 Area then \ominus



Decision Trees :

less Entropy + High Knowledge Gain

$\log \rightarrow$ will convert Product to Sum

$$\therefore \log(ab) = \log(a) + \log(b)$$



$$\text{Entropy} = -\frac{5}{8} \log_2\left(\frac{5}{8}\right) - \frac{3}{8} \log_2\left(\frac{3}{8}\right) = 0.9544$$

or

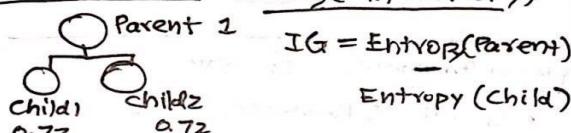
$$\text{Entropy} = -\frac{m}{m+n} \log_2\left(\frac{m}{m+n}\right) - \frac{n}{m+n} \log_2\left(\frac{n}{m+n}\right)$$

$$\text{if } p_1 = \frac{m}{m+n}, p_2 = \frac{n}{m+n}$$

$$\text{Entropy} = -\sum_{i=1}^n p_i \log_2(p_i)$$

Information Gain: (Change in Entropy)

i.e.

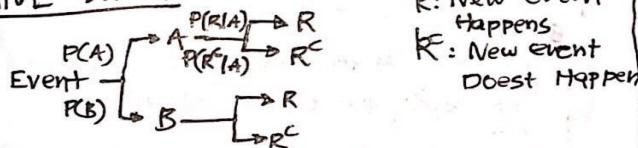


$$IG = 1 - 0.72 = 0.28$$

Hyper parameters:

- ① **Max_depth:** Max No. of levels in the tree.
- ② **Min_samples_leaf:** Min. No. of samples Allowed in leaf
- ③ **Min_samples_split:** Min. No. of samples Allowed in split
- ④ **Max_features:** The Number of features to consider when looking for best split.

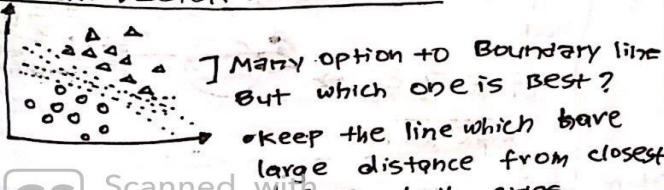
NAIVE BAYES:



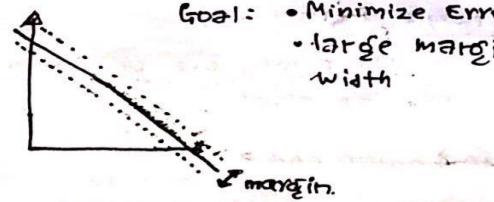
$$\text{W.K } P(R|A) = P(R^C|A) = P(A) \cdot P(R|A)$$

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

SUPPORT VECTOR MACHINES (SVM)

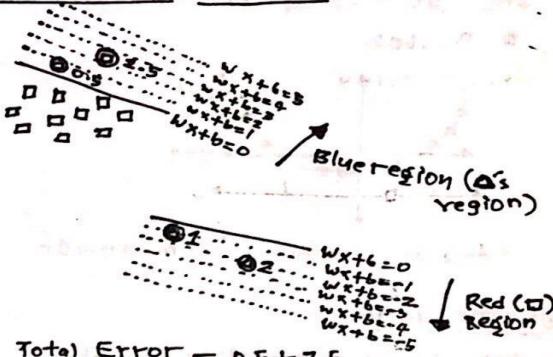


- Goal:
- Minimize Error
 - Large Margin width



ERROR = Classification Error + Margin Error

Perceptron Algorithm:



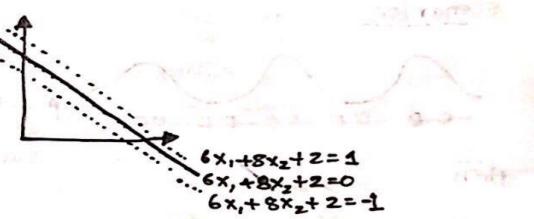
Idea is to use Gradient Descent to give the best possible cut.
→ Perceptron Algorithm.

Margin Error:

- Small margin → Large Error
- Large margin → Small Error

Ex:

$$w_1 x_1 + w_2 x_2 + b = 0 \quad \|w\|^2 = 3^2 + 4^2 = 25 \\ 3x_1 + 4x_2 + 1 = 0 \quad \# \text{Error} = 25$$



This gives same error as L1 + L2
Regularisation.

C-Parameter:

It's a constant that attaches to Error

$$\text{Error} = C * \text{Classification Error} + \text{Margin Error}$$

1. If C is very large → we are focusing on Classification Error
 2. If C is small then we are focusing on Margin Error
- Small C → Large Margin → May make classification Error
 - Large C → May have small margin → classifies points well.

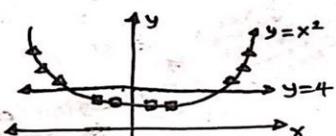


① Polynomial Kernel :-1.

- What if the points are on a line; not on 2D plane?



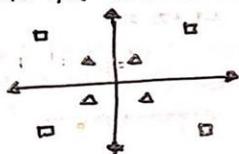
To separate this with a boundary line we need to use a model → "Polynomial Kernel"; it converts the 1D → 2D ; puts all datapoints on a Parabola.



- And then convert from 2D → 1D.

② Kernel trick-2:

How to fit a boundary to this.

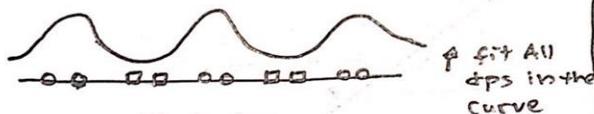


1. Fit All (Δ) in a circle.
2. Send (Δ)s to Another Dimension (lets say Z-Axis).

RBF KERNEL :-

- Radial Basis Function (RBF)

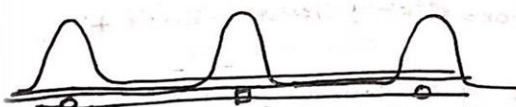
Scenario:



then



- But the question is how to do this?
- Build mountains on top of every DP.



this technique is called "RBF"

- Now combine all the Mountains

In Higher Dimension:

- Draw a Gaussian Paraboloid

the separate with a plane ; this plane will cut this Gaussian Paraboloid at some height & And this is a circle

THIS CIRCLE WILL GIVE THE BOUNDARY EQUATION

Boundary Plane...

IN 3D / N-D

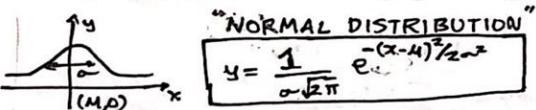
But How to decide the width of the Mountain?



∴ We use "r Parameter"

• Large r → Narrow curve ②

• Large r → wide curve ①



$$r = \frac{1}{2\sigma^2}$$

Hyper Parameter: In SVC

• C : The C Parameter

• Kernel : The Kernel : The most common ones are 'linear', 'poly' and 'rbf'

• degree: If kernel is Polynomial ; This is the Maximum degree of the monomial in the Kernel

• gamma : If kernel is RBF ; This is 'r'

ENSEMBLE :

1. Bagging (Bootstrap aggregating)

2. Boosting

Bagging

BAGGING:

Use Many Model for a Job ; combine / Avg / Add (etc.) the output of all

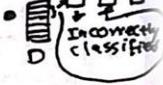
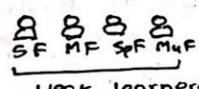


• Science friend
Math friend
Sports friend
Musician friend

• which Answer have Highest vote
that will be chosen.

BOOSTING:

• combine the Models → to make a strong model



BIAS:

• If a Model have high bias → it won't fit well with the Data

VARIANCE:

• If a Model have High Variance → means it changes drastically



High Bias

High variance.

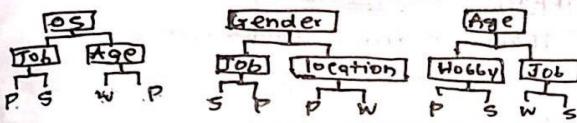
2 Methods to solve:

- ① Bootstrap the Data : Sampling the Data
- ② subset the features:

RANDOM FOREST: (BAGGING)

gender	Age	location	os	Job	Hobby	APP
--------	-----	----------	----	-----	-------	-----

are features → then Randomly PICK Any 3 features for Many times
→ the most Appeared Answer will be used as output.



Bagging: Row subsets given to Models and voting of the Models is considered. (K-Fold).

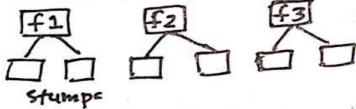
ADABOOST: (BOOSTING TECHNIQUE)

- ① Decision trees are created ; with only one depth

$$f_1 \ f_2 \ f_3 \ o/p$$

② sample weight $w = \frac{1}{n}$
 $\frac{1}{n}$
 $\frac{1}{n}$
 $\frac{1}{n}$
 $n: \text{no. of records}$
 rows

- ③ create Stumps / 1 level trees for 3 features



Compare the Entropy of all Stumps
→ Choose less Entropy Stump

- ④ calculate total Error of Stump:

$$\rightarrow \text{Total Error(TE)} = \sum \text{weight of Incorrectly classified records.}$$

- ⑤ calculate the Performance of Stump:

$$\rightarrow \text{Performance} = \frac{1}{2} \log_e \left[\frac{1-TE}{TE} \right]$$

- ⑥ Increase the weights of Incorrectly classified records

- Decrease the weights of correctly classified.

$$\rightarrow \text{New Sample weight} = \text{old weight} \times e^{-\text{Performance}}$$

$$\rightarrow \text{New sample weight} = \text{old weight} \times e^{-\text{Performance}}$$

$$f_1 \ f_2 \ f_3 \ o/p \quad \text{sample Updated weight} \quad \text{Normalized weight}$$

$$\text{Normalized weight for row}(i) = \frac{\text{Updated weight of row}(i)}{\sum \text{updated weights}}$$

NOTE: $\sum \text{Normalized weights} = 1$

- ⑦ consider Normalized weight further

f1	f2	f3	o/p	Normalized weight
----	----	----	-----	-------------------

• create buckets

Normalized weights	Buckets
0.07	0 → 0.07
0.51	0.07 → 0.58 (0.51)
0.07	0.58 → 0.6
0.07	0.65 → 0.72

- ⑧ New Dataset is Created :-

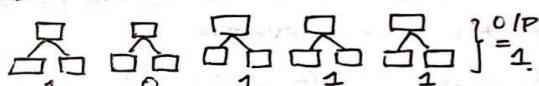
• This Dataset will consider some random Number (say 0.43)

• It takes row belonging to the bucket value range 0.43

• Repeats this with other random values (0.31...) till 'n' number of Iterations

- ⑨ Repeat from step 1 to 6 for all the Learners

- ⑩ Majority of votes from all the stumps of all learners will be considered as output.



hyperparameters: In sklearn

① base_estimator: The model utilized for weak learners

② n_classifier: The Maximum Number of weak learners used.

MODEL EVALUATION METRICS:

• when an Algorithm is working well with the Data And.

• How to tweak them?

• Training set and Testing set → to test

Golden Rule: Thou Shalt Never use your Testing Data for training

CONFUSION MATRIX:-

-Medical Model

		Diagnosed as SICK	Diagnosed as healthy
		TP	FN
SICK	TP		
	FP		TN
Healthy	FN		
	FP		TN

Q. Accuracy:

$$Acc = \frac{\text{Correctly classified}}{\text{Total records}}$$

③ Precision and Recall:

In Medical Model
FP is OK
FN is NOT OK

In Spam Doctor Model
FP is NOT OK
FN is OK

Precision:

		Diagnosed as Sick	Diagnosed as Healthy	} Recall
Sick	1000	200		
Healthy	800	9000		

$$\text{Precision} = \frac{1000}{1000+800} = \frac{TP}{TP+FP}$$

Recall:

$$\text{Recall} = \frac{1000}{1000+200} = 83.3\%$$

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

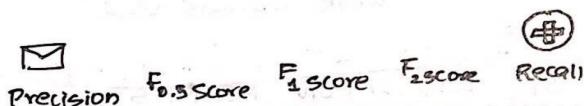
④ F1 Score:

- Its used to combine Precision & Recall by taking its Average
- Although its Not Pure Avg.

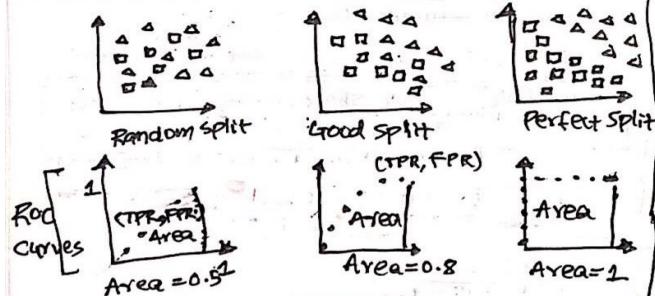
$$\therefore F1 \text{ Score} = 2 \cdot \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

⑤ F_B Score:

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



⑥ ROC CURVE:



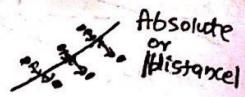
Calculate True Positive Rate and False Positive Rate

$$TPR = \frac{TP}{\text{All Positives}}$$

$$FPR = \frac{FP}{\text{All Negatives}}$$

Regression Metrics:

- Mean Absolute Error



Absolute or Distance

- Mean Squared Error

- R²-Score:

$$R^2 \text{ Score} = 1 - \frac{\text{Mean squared Error}}{\text{Simplest model squared Error}}$$

Simplest model:

If R² score close to 0 } BAD Model

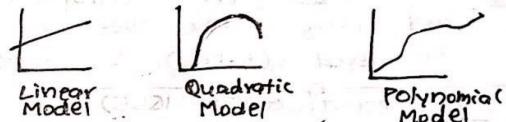
If R²-score close to 1 } GOOD Model.

TRAINING AND TUNING:

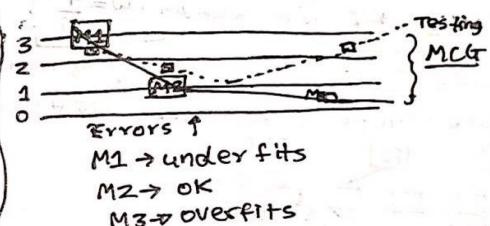
- Training error : Underfitting

- Testing Error : overfitting

Model Complexity Graph:



We calculate training and testing Error and plot that in graph → called → MCG



K-Fold Cross Validation:

Split 'D' into 'K' number of splits and then train and test with it.

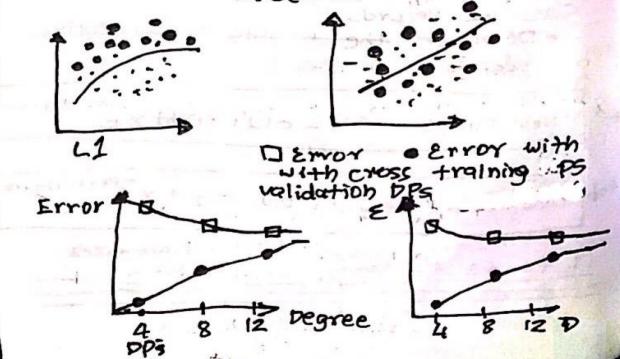


LEARNING CURVES:

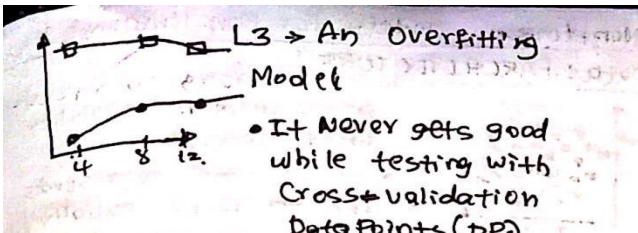
Used to check overfitting, underfitting of Model

Learner-1 : Linear Model

Learner-2 : Polynomial Model with degree=2.



Scanned with
CamScanner



GRID SEARCH CROSS VALIDATION:

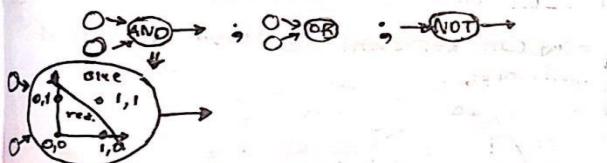
- For decision Tree try all depths (ex: depth $\geq 1, 2, 3, 4, \dots$) Pick the one with Highest F1 score,
- For polynomial models try all degree \rightarrow Pick Highest F1 scored model
- This is called "Grid search"

DEEP LEARNING



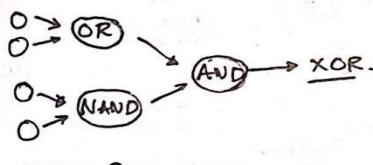
- University Student Admission classification from Perceptron Algorithm.

- Perceptron as logical operator.



XOR:

- We Need Multilayered Perceptron.



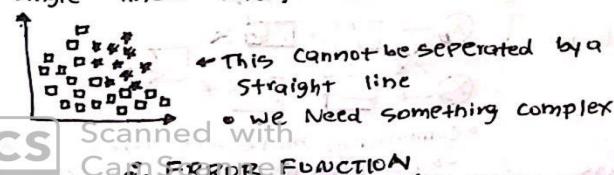
↳ DP Miss classified.
↳ DP wants the line to come closer!

Perceptron trick :- (In S. Learning section)

Perceptron Algorithm:

Non-linear regions:-

- What if the Data cannot be separated by a single linear line?



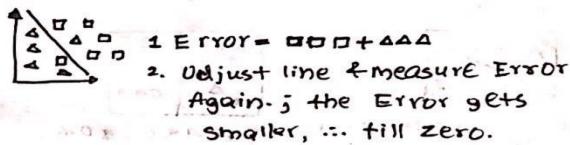
ERROR FUNCTION

- To calculate the error

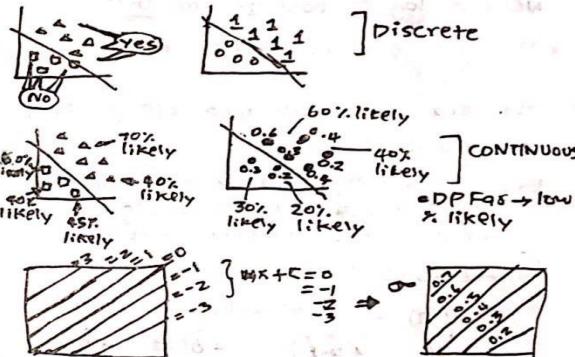
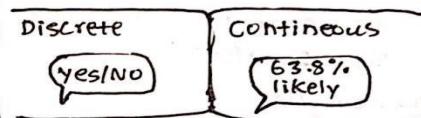


- Look around all the directions \rightarrow (se) Decrease Error
- But small steps (learning rate η) \rightarrow so we use "Gradient Descent"

Log-loss Error Function:



- But we Need continuous values in Prediction / dependant variable to do Any Gradient Descent; so If it is Discrete \rightarrow convert to \rightarrow continuous



$$\begin{aligned} &\therefore \text{Sigmoid Applied. } \\ &y = \sigma(mx+c) \\ &\text{Step}(mx+c) \end{aligned}$$

SOFTMAX:

- What if we want to classify multi-class data points like cat, dog, horse etc... ?
- Exp: This will convert -ve to +ve Number
- Bcz Probability should not have any value outside the $0 \rightarrow 1$

ONE hot ENCODING:

- Categorical \rightarrow Numerical

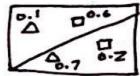
Animal	Cat?	Dog?	Horse?	...
Cat	1	0	0	
Dog	0	1	0	
Horse	0	0	1	

Maximum Likelihood:

$$\begin{array}{c} \Delta 0.1 \\ \Delta 0.2 \\ \Delta 0.7 \end{array}$$

Maximal likelihood.
 $0.6 * 0.2 * 0.1 * 0.7 = 0.0084$
 $0.7 * 0.9 * 0.8 * 0.6 = 0.3024$

- The likelihood that the Model is correct
 - A Good Model will give good Probability
- But How to Maximize it??



$$0.6 * 0.2 * 0.1 * 0.7 = 0.0084$$

0.8	0.7
0.6	0.9

$$0.7 * 0.9 * 0.8 * 0.6 = 0.3024 \leftarrow \text{Max}$$

meaning: (0.8) The model is 0.8 or 80% confident that the dp is 0.8

Red Point

CROSS ENTROPY:

But Products will take long if lot of DPs present. ∴ We convert it to sums by using log() function.

$$\log(ab) = \log(a) + \log(b)$$

We use log of base 10 or $\ln()$

The $\ln()$ -ve of Prob of All dp's is called

-ve bcz $\ln()$ will give 0/P in Perfect -ve Numbers.

$$\text{Ex: } 0.7 * 0.9 * 0.8 * 0.6 = 0.3024$$

$$\ln(0.7) + \ln(0.9) + \ln(0.8) + \ln(0.6)$$

$$-0.36 \quad -0.1 \quad -0.22 \quad -0.51$$

∴ -ve is used

$$-\ln(0.7) - \ln(0.9) - \ln(0.8) - \ln(0.6)$$

$$+0.36 \quad +0.1 \quad +0.22 \quad +0.51$$

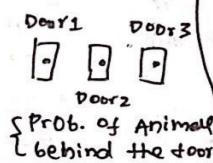
This is called "Cross entropy"

- Good Model: less entropy
- Bad Model: large entropy

$$\text{Cross Entropy} = -\sum_{i=1}^m y_i \ln(p_i) + (1-y_i) \ln(1-p_i)$$

Multi-Class - Cross - Entropy :-

Animal	DOOR1	DOOR2	DOOR3
CAT	$P_{11} 0.7$	$P_{12} 0.3$	$P_{13} 0.1$
DOG	$P_{21} 0.2$	$P_{22} 0.4$	$P_{23} 0.5$
HORSE	$P_{31} 0.1$	$P_{32} 0.3$	$P_{33} 0.4$



If the Actual thing is: DOOR1 → CAT
 DOOR2 → Horse
 DOOR3 → Horse

then the Prob table Made by Model

$$P = 0.7 * 0.3 * 0.4 = 0.084$$

$$\text{Cross Ent} = -\ln(0.7) - \ln(0.3) - \ln(0.4)$$

$$\therefore \text{CE} = -\sum_{i=1}^n \sum_{j=1}^m y_{ij} \ln(p_{ij}) \text{ for Multi class}$$

Non-linear Models for Non-linearly separable Data: [ARCHITECTURE]

we are going to separate it by using a probability based separator

This line represents the DP to be equally likely to both (red) and (a) blue

But how do we do this?!

combine 2 linear models

$$\begin{array}{c} 0.7 \\ \text{Prob} \end{array} + \begin{array}{c} 0.8 \\ \text{Prob} \end{array} = 0.7 + 0.8 = 1.5 \Rightarrow \text{Prob.} = 0.82$$

same dp in 2 different linear model have 2 different values; what do we do when we combine them?

∴ we use 'Sigmoid' (bcz $0.7 + 0.8 = 1.5 \Rightarrow 0.7 + 0.8 = 0.82 \Rightarrow$ this is the Prob. of that dp is resulting Model)

But what if we want more of first Model and little less of second Model? ∴ we add weights $0.7, 0.8$

$$\begin{array}{c} 7 \\ \text{Prob} \end{array} + \begin{array}{c} 5 \\ \text{Prob} \end{array} = 7 \times 0.7 + 5 \times 0.8 = 2.9 \rightarrow 0.95$$

we can even add the bias

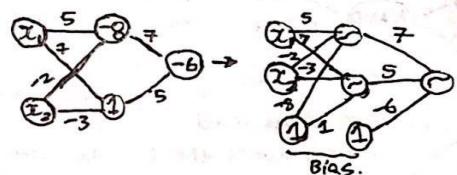
We can represent the above in Neural Networks:

$$\begin{array}{c} 5 \\ \text{Prob} \end{array} = \begin{array}{c} (x_1 5) \\ (x_2 -8) \end{array} \cdot \begin{array}{c} 0.7 \\ 0.5 \\ -6 \end{array} \rightarrow$$

$$\begin{array}{c} 7 \\ \text{Prob} \end{array} = \begin{array}{c} (x_1 7) \\ (x_2 -3) \end{array} \cdot \begin{array}{c} 1 \\ 1 \\ -3 \end{array} \rightarrow$$

To Add these both with a -6 bias

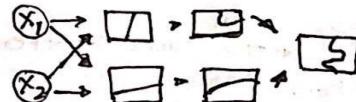
lets join this:



Input layer, Hidden layer, Output layer.

If Input layer have 'n' Nodes then the DPs will be in n Dimension space

The Node at output layer represents one class of classification [Ex cat, dog etc..]



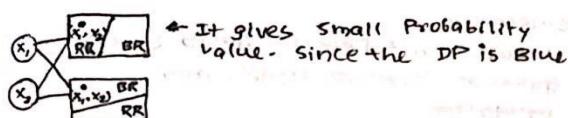
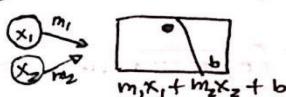
Combining more non linear model → Complex Non-linear model

Multiclass Classification:

- ① : NN to find Cat ② : NN to find Dog ③ : NN to Identify Horse

combine all this and Add softmax function to obtain the Probability of O/P.

Feed Forward :

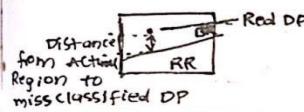


-- Blue DP
BR - Blue Region
RR - Red Region

ERROR FUNCTION:

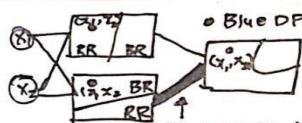
$$E(w) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1-y_i) \ln(1-\hat{y}_i)$$

$$\hat{y} = \sigma(Wx+b)$$



-- Back Propagation.
by using Gradient Descent

Backpropagation:



• Increase the weight coming from bottom Model Bcz it classified opp correctly (or have large Prob value)
• And decrease the weights of model which gave Miss classification (or less Prob./or value).

• In first Model the DP wants or Bcz of DP ; the line will come close to OP

• In second Model the DP wants the line to go away Bcz it correctly classified

• So the output model will also be affecting from this change

CHAIN RULE:

$$(1) f \rightarrow A \rightarrow B \\ A = f(x) \quad B = g(f(x))$$

$$\frac{\partial B}{\partial x} = \frac{\partial B}{\partial A} \frac{\partial A}{\partial x}$$

Scanned with
CamScanner

- Gradient Descent, local minima
- to avoid Local minima add Momentum
- Momentum is like replacing the man on top of the hill with a giant ball ; so it can avoid Local minima
- what if we don't know the weights w_1, w_2 ? :: we use Error Function $E = \sum (y_i - \hat{y}_i)^2$

$$\hat{y} = \text{New o/p}$$

$$y = \text{Actual o/p}$$

Sum of the squared Error

$$E = \frac{1}{2} \sum (y_i - \hat{y}_i)^2$$

$$w_i = w_i + \Delta w_i$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

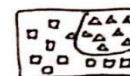
$\eta \rightarrow$ learning rate

TRAINING OPTIMIZATION:

- overfitting ; underfitting ;



underfitting (high bias)



Just Right



overfitting

STOPPING CRITERIA:

- Stop epochs before it becomes overfit

MODEL COMPLEXITY GRAPH:



underfitting



Just Right



overfitting



• TE - Testing Error
• TRE - Training Error

Regularization:

$$(1, 1) \quad (4, 9) \quad \text{which can we use?} \\ \text{or } x_1 + x_2 \\ \text{or } 10x_1 + 10x_2$$

$$\text{SOLUTION 1: } x_1 + x_2$$

$$\sigma(1+1) = \sigma(2) = 0.88(2) \\ \sigma(-1-1) = 0.12(0)$$

$$\text{SOLUTION 2: } 10x_1 + 10x_2$$

$$\sigma(10+10) = 0.99...79(0) \\ \sigma(-10-10) = 0.00...21(0)$$

- But the problem is we need to train for large amount of epochs because it have 10 ; which the gradient descent should attain from a smallest number ; with small steps ; And it overfits

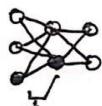
∴ To avoid overfitting we use regularization

$$\text{L1 Error Function} = -\frac{1}{m} \sum_{i=1}^m (1-y_i) \ln(1-\hat{y}_i) + y_i \ln(\hat{y}_i) + \lambda(|w_1| + \dots + |w_n|)$$

$$\text{L2 Error Function} = -\frac{1}{m} \sum_{i=1}^m (1-y_i) \ln(1-\hat{y}_i) + y_i \ln(\hat{y}_i) + \lambda(w_1^2 + \dots + w_n^2)$$

DROPOUT:

- Turning some of the nodes "OFF" randomly



- Specified in probability value

EX Dropout Prob. 0.2.

Dropped out for an epoch

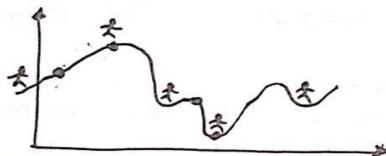
LOCAL MINIMA:



Solution 1:

RANDOM RESTART:

- Start from many random points & descend



VANISHING GRADIENT:

- With sigmoid function weight updation is

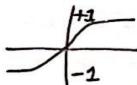
$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial h_i} \frac{\partial h_i}{\partial w_{ij}} \quad] \text{Product tiny No. will very tiny member.}$$

∴ Descending from error mountain takes a lot of time.

How to fix it? Best way to fix is changing activation function

1. tanh:

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



2. Relu: Rectified Linear Unit:

$$\text{relu}(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$



STOCHASTIC GRADIENT DESCENT:

Take some batch and feed forward + back propagate & take another batch; repeat this step.

Learning rate decay

Large $\eta \rightarrow$ may give less accuracy; take tiny η

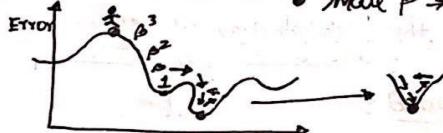
Big $\eta \rightarrow$ Big steps while descending from the mountain.

MOMENTUM:

- Take avg. of previous steps.

β : BW 0 & 1

- Big $\beta \rightarrow$ Big steps/Roll
- Small $\beta \rightarrow$ Small Rolls.



UN-SUPERVISED LEARNING:

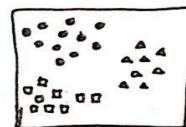
Clustering

Hierarchical & Density based clustering
Gaussian Mixture Models and clustering validation

Principal component Analysis.

Random Projection and Independent component Analysis.

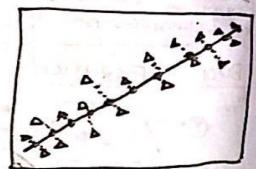
Clustering & Dimensionality reduction:



unsupervised learning task concerned with putting similar data into groups "CLUSTERING"

Dimensionality Reduction:

Unsupervised task concerned with reducing the number of variables used.

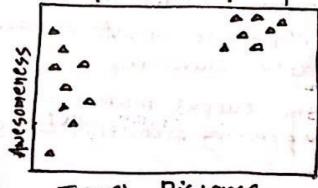


K-Means:

use cases:

- Netflix & grouping similar movies for recommendation
- Grouping books/things for recommendation
- Cluster the group of customers who are similar to each other.
- Similar movies/music

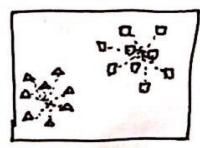
K-Means: to see family Travelling to work.



In K-Means we group the DPs to k-number of groups

- How to choose the value of 'k'?

Elbow Method:



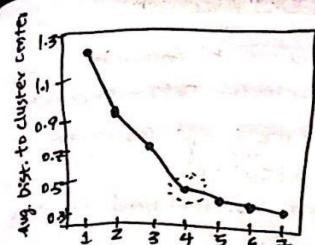
K=1, avg.dist=1.261

K=2, avg.dist=0.923

K=3, avg.dist=0.639

K=4, avg.dist=0.382

⋮



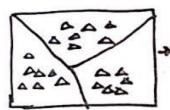
$K=7$; avg dist = 0.298

- we don't need smallest avg. dist value;
- bcz as K increases Avg. distance from center to dp increases.
- so we need a point where 'Elbow' here it is at (4)

- after which the Avg. dist significantly decreased.

How K-Means Works?

1. Randomly place K-centroids amongst your data
2. Assign each point to the closest centroid
3. Move the centroid to the center of the points assigned to it.
4. Repeat 2,3 until convergence



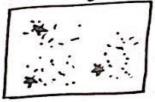
Is that the optimal solution?:

- Start K-means with different starting points

Starting set 1



Starting set 2



Starting set 3



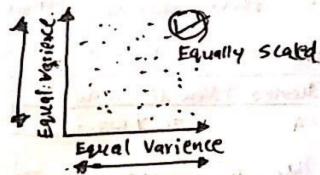
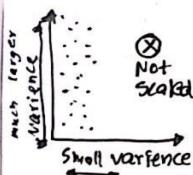
- use repeated runs to protect against Local Minima
- choose best grouping sets; (less Avg. Dist to centroids).

Feature Scaling:

- Standardizing / Z-score scaling \rightarrow
- Normalizing / Min-Max \rightarrow

- we use Normalizing \rightarrow when scaling the color of an image

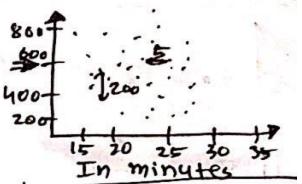
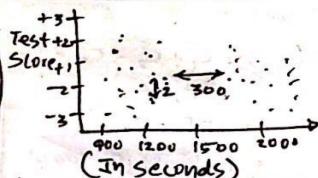
- the features with much larger variance will dominate in the clustering; we do not want this to happen; so we scale.



- Standardizing: moves variables to mean=0 & $Std=1$

- Normalizing: moves variable b/w 0 and 1

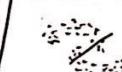
Feature Scaling Example:



| 5K completion time(s) |

| 5K completion time (m) |

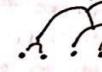
But how to cluster for some datasets that look like this



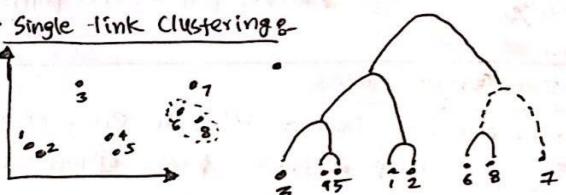
so we use

1. Hierarchical clustering:

- This shows how clusters relate to each other



→ Single-link clustering:



- every DP is a cluster at initial
- cluster 2 DPs of closest distance; continue this with all DPs.

- In (6,8) cluster; what about 7?; calculate distance b/w closest point in the cluster (Here 6) to 7; make a new cluster of (6,8) and 7 as shown above Dendrogram
- We can decide how many clusters we want
∴ if Cluster = 2 then there will be no cluster b/w (6,8) to (3,4,5,1,2)

→ Complete link clustering:

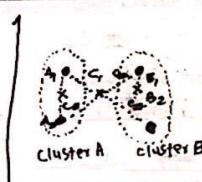
- Everything is same as single linked clustering. But in 6,8 to 7; calculate longest dist. between (6,8); make a cluster to longest dist DP to 7; here 8 and 7 (6,8,7)

→ Average link clustering:

- Everything is same as single link but in case of (6,8) and 7; calculate every possible distance from every DP to every other DPs in other cluster

[Avg. of Distance is the Distance between 2 Clusters]

1. Ward's Method:



$$\Delta(A, B) = C_1^2 + C_2^2 + C_3^2 + C_4^2 - A_1^2 - A_2^2 - B_1^2 - B_2^2$$

• whichever is least we merge those 2 clusters

• This is the Default hierarchical clustering in sklearn.

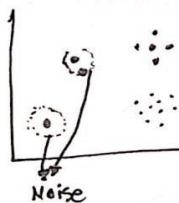
CONS of All Hierarchical Clustering:

- Computationally intensive $O(N^2)$
- Sensitive to Outliers.

2. DBSCAN:

takes i/p: ϵ : search distance Around Point

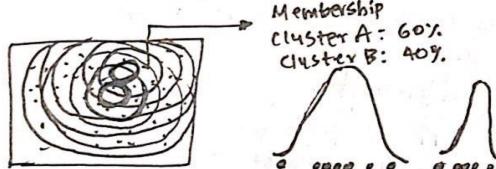
Min Pts: Minimum No. of Points required to form a density cluster



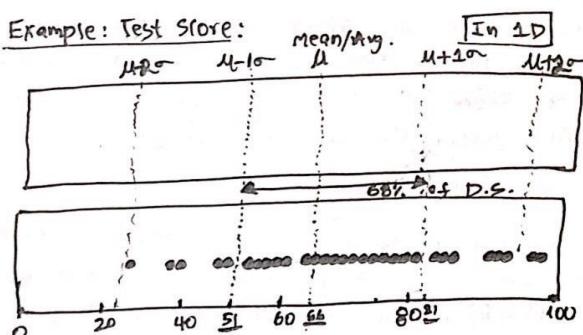
- It checks the Neighbour DP counts for all DPs and if Min No. of pts. (i.e. Here 5) present; makes it as a cluster.
- Border points; core point

Gaussian Mixture Models:

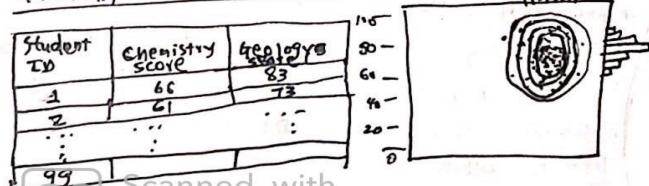
- Every points in Dataset belongs to every Cluster in DS. ; but have different levels of membership.



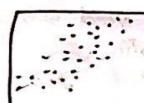
Example: Test Score:



Gaussian Distribution in 2 Dimension:



Expected Maximization (EM) Algorithms



- Step 1: Initialize 'K' Gaussian Distributions
- Step 2: Soft-cluster Data - "Expectations"
- Step 3: Re-estimate the Gaussians - "Maximization"
- Step 4: Evaluate log-likelihood to check for convergence

Repeat from Step#2 until converged

Explanation:

1. Initial epoch

Cluster	μ_1	σ^2_1
A	64.763	100
B	46.512	57

• either by μ & σ^2 of the DS or

• By using K-means (use the centroids as mean)

2. Soft Cluster DPs:

Point #	feature1	feature2	Cluster A	Cluster B
1	62	71	0.99976	0.00024
2	58	81		
3	52	74		
...		
N	52	78		

• i.e. the Membership Percentage of DP to each of Clusters. (Expectation)

= Expectation (of Membership) of point 1 belonging to cluster A' is denoted as follows:

$$E[z_{1A}] = \frac{N(X_1 | \mu_A, \sigma^2_A)}{N(X_1 | \mu_A, \sigma^2_A) + N(X_1 | \mu_B, \sigma^2_B)}$$

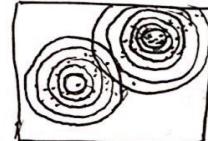
$$N(X | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

$$E[z_{1A}] = \frac{0.001288}{0.001288 + 0.0000038}$$

• x is row in the DS.

$$= 0.99976 \rightarrow \text{This means that we are 99.976% sure that the point}$$

$$E[z_{1B}] = 1 - E[z_{1A}] \rightarrow \text{belongs to cluster A.}$$



Step#3: Make New (μ) and New (σ^2) by using Above Gaussian clusters

New Gaussian parameters:

Cluster	New μ	New σ^2
A	(64...76...) 103, ...	
B	(46...51...) 67, 10...	

$$\text{New } \mu_A = \frac{\sum_{i=1}^N E[z_{ij}] X_i}{\sum_{i=1}^N E[z_{ij}]} = 0.99937 \times (62.71) +$$

$$= 0.99937 \times (62.71) + 0.9998 \times (58.51) + \dots$$

$$0.99937 + 0.9998 + 0.55818 + \dots$$

Point #	feature 1	feature 2	cluster A	cluster B
1	62	71	0.99937	0.00063
2	58	81	0.9998	0.0002
3	52	74	0.55218	0.44182
...
N	52	78	0.99133	0.00867

$$\text{New } \mu_A = (64.4872457, 76.3074590)$$

$$\text{New } \sigma_A^2 = \frac{\sum_{i=1}^N E[Z_{iA}] (X_i - \mu_A^{\text{new}})(X_i - \mu_A^{\text{new}})^T}{\sum_{i=1}^N E[Z_{iA}]}$$

$$= \underline{103.92494596}$$

Now we have new Gaussians, the cluster will move up or down; do this until it converges.

Step#4: Evaluate log-likelihood

$$\ln p(X | \mu, \sigma^2) = \sum_{i=1}^n \ln \left(\sum_{k=1}^K \pi_k N(X_i | \mu_k, \sigma_k^2) \right)$$

$\pi_k \rightarrow$ mixing coefficient.

The better gaussian \rightarrow good or Maximum log likelihood it have; we repeat entire process until we get better likelihood.

CLUSTER VALIDATION:

- External validation Indices:
- Internal validation Indices:

1. External validation Indices:

$$\text{Rand Index} = \frac{a+b}{\binom{n}{2}}$$

a: Number of pairs in the same cluster C and in the same cluster K.

Actual label (C)

clustering result (K)

b: number of pairs in a different cluster C and in a different cluster in K.

n: number of samples or points

$$\therefore \text{Rand Index} = \frac{2+4}{\binom{10}{2}} = \frac{6}{45} = 0.6$$

$$ARI = \frac{RI - \text{Expected Index}}{\max(RI) - \text{Expected Index}}$$

High ARI \rightarrow Good cluster (Near to 1)

Low ARI \rightarrow Bad cluster (Near to 0)

INTERNAL INDICES:

- Working with D.S.'s that are labelled.
- Silhouette coefficient:
- Between -1 and +1 related with

$$S_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

for i-th point.

a: avg. distance to other samples in the same cluster

b: avg. distance to samples in the closest neighboring cluster

$$S = \text{average}(s_1, s_2, \dots, s_n)$$

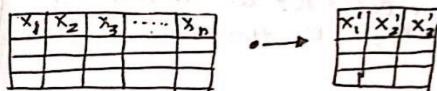
Principal Component Analysis (PCA)

- PCA is a technique aimed at dimensionality reduction.
- This technique is about taking the full dataset and reducing it to only the parts that hold the most information.

Have you had Breakfast?

→ So, I was running late and.. I stopped by the coffee shop, and I had a coffee, and then.. I grabbed a pastry.

→ "I had a coffee and a Pastry"



size of house
lot size x
number of rooms x
Floor Plan size

Neighbourhood Quality
local crime rate
no. of schools in 5 mil.
Property tax rate

FloorPlanSize
localCrimeRate

But makes us loose information; Rather than dropping features to reduce the dimensionality of dataset, if we were able to combine the features in some way to maximize the amount of information we retain; → This is the goal of PCA.

Now,

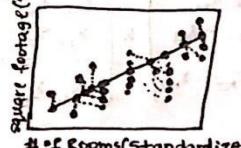
$$b_1(\text{lot size}) + b_2(\#\text{rooms}) + \dots$$

size feature

Neighborhood Quality feature

$$b_3(\text{crime rate}) + b_4(\text{schools}) + \dots$$

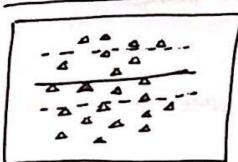
square footage (standardized)



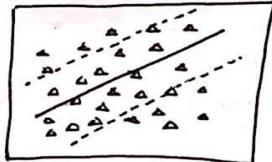
Dimensionality Reduction.

- Properties of Principal Components:
- Each component captures the largest amount of variance left in the data.

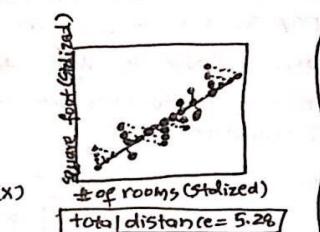
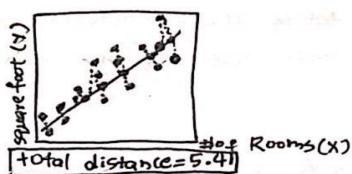
Not the largest variance



The largest variance

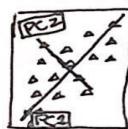


- The Information lost is the sum of dops to the components.



This will loose less info

- Components are orthogonal to one another

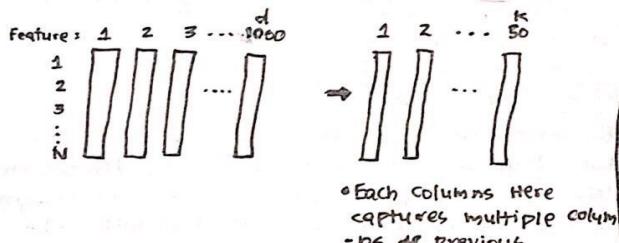


PC1 & PC2 are orthogonal if Angle b/w them is 90°

Principal components are also known as "eigenvectors"

RANDOM PROJECTION:

- Is a powerful Dimensionality Reduction Method



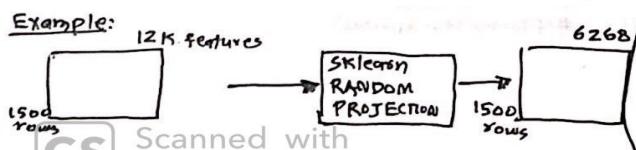
od is dimension \rightarrow into 'K' (either user given or computed) dimension

$$X_{K \times N}^{RP} = R_{N \times K} X_{d \times N} \quad [\text{original feature Matrix}]$$

(New Matrix)
Matrix of
Random
Numbers

$$\begin{bmatrix} X_{11}^{RP} & X_{12}^{RP} & \dots & X_{1n}^{RP} \\ X_{21}^{RP} & X_{22}^{RP} & \dots & X_{2n}^{RP} \\ \vdots & \vdots & \ddots & \vdots \\ X_{d1}^{RP} & X_{d2}^{RP} & \dots & X_{dn}^{RP} \end{bmatrix} = \begin{bmatrix} Y_{11} & Y_{12} & \dots & Y_{1d} \\ Y_{21} & Y_{22} & \dots & Y_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ Y_{K1} & Y_{K2} & \dots & Y_{Kd} \end{bmatrix} \times \begin{bmatrix} X_{11} & X_{12} & \dots & X_{1n} \\ X_{21} & X_{22} & \dots & X_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{d1} & X_{d2} & \dots & X_{dn} \end{bmatrix}$$

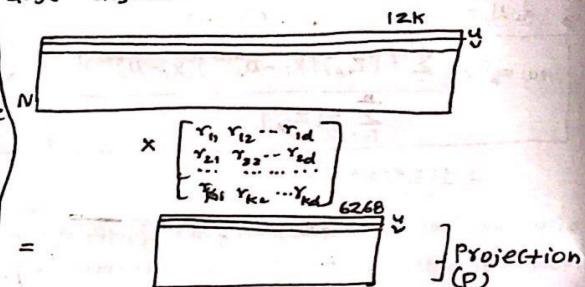
Example:



What is the theory behind RANDOM PROJECTION?

- Johnson-Lindenstrauss lemma:

A dataset of N Points in High Dimensional Euclidean Space can be mapped down to a Space in much lower dimension in a way that preserves the distance between the points to a large degree.



The distance b/w 2 points u, v in new DS is

$$\|P(u) - P(v)\|^2 > (1 - \epsilon_{PS}) \|u - v\|^2 < (1 + \epsilon_{PS}) \|u - v\|^2$$

ϵ_{PS} : level of error that we allow in Random Projection;

\therefore Distance b/w DP u, v in Projection will be larger than Distance b/w u, v in Original dataset

In Sklearn default value of epsilon (ϵ_{PS}) is 0.1

USE CASES:

- Random Projection before Applying K-means clustering to 400 face images of 64×64 dimensions.

- Application of RP in image & text data.

INDEPENDENT COMPONENT ANALYSIS (ICA):

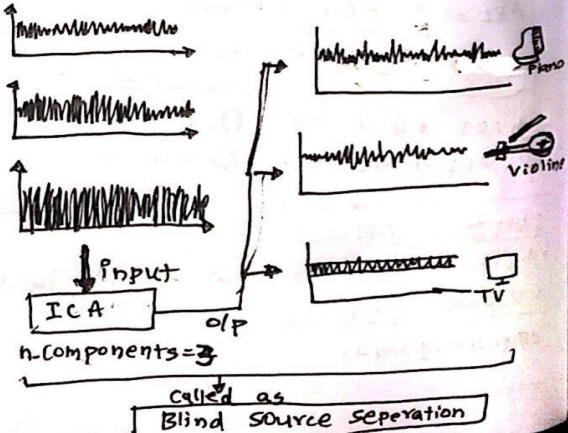
USE CASE SCENARIO:

- 3 persons go to an Auditorium with 3 curtains. One from TV, one from Piano & one from Violin.



- All of them record their favorite show sounds by using their phone recorder.

Now is it possible to separate individual of sound from all of these recordings?



fast ICA algorithm:

$$w = E\{xg(w^T x)\} - E\{g'(w^T x)\}w$$

$$g_1(u) = \tanh(a_1 u).$$

Algorithm:

#1: center, whiten x.

#2: choose initial Random weight Matrix $w_1, w_2 \dots w_n$

#3: Estimate W_g containing vectors.

#4: decorrelate W

#5: Repeat from step #3 until converged.

• NOTE: There should be n-number of recordings in order to get n-number of components out of it.

ICA Applications:

Medical fields:

• In Brain scanner \rightarrow EEG or MEG
(like the one in CHAPPIE movie); where each of these sensors sends the electric signals from the brain.

Electroencephalographic data \rightarrow EEG

↓
Apply ICA to separate the mixture.

• In MRI scan Images.

• In financial Data (stock signals). What factors that causes stock to go up & down.

TERM 1 OVER

