

## PART - 3 [Classification]

continuous number

- Unlike regression, we predict a category here.

• These include

- linear models like Logistic Regression, SVM
- non-linear models like K-NN, Kernel SVM & Random Forests.

Classification Models:

(1) Logistic Regression

(2) K-Nearest Neighbors (K-NN)

(3) Support Vector Machine (SVM)

(4) Kernel SVM

(5) Naive Bayes

(6) Decision Tree Classification

(7) Random Forest Classification



## ① Logistic Regression.

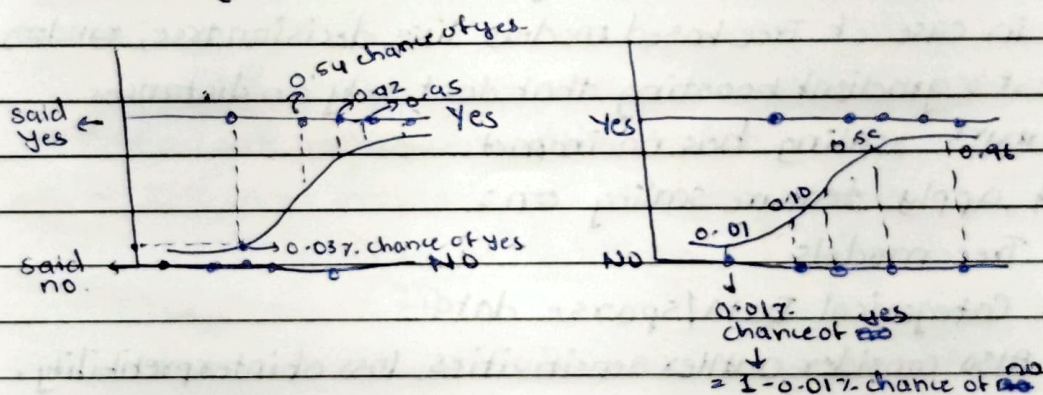
- Predict a categorical dependent variable from a no. of independent variables.
- Need to scale data here.

$$\ln \frac{P}{1-P} = b_0 + b_1 x_1 + \dots + b_n x_n$$

- a statistical method used for binary classification.
- Predicts probability of binary outcome.
- outputs probabilities mapped to discrete classes.

### \* Maximum Likelihood Estimation:

- Logistic R. uses MLE to estimate coeff. of model.



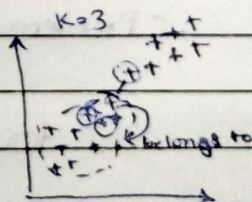
$$M.L. = 0.03 \times 0.54 \times \dots \times (1 - 0.01) \times (1 - 0.04) \times \dots$$

(multiply all these numbers).

and the curve with the max. likelihood is best curve.

## ② K-NN.

- Find euclidean distance from all points.
- Select the K nearest ones (min. distance)
- Take the max vote to identify class.



- Need to rescale data in KNN too, or a single feature will weigh too much on output.



Max margin Hyperplane  
(Max margin classifier)

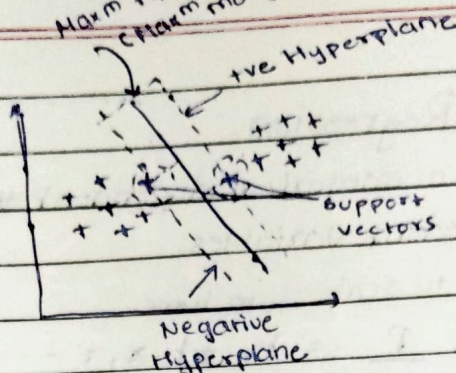
DATE \_\_\_\_\_

PAGE NO. \_\_\_\_\_

### ③ SVM.

(Linear).

Coz of linearly separable data



Working Principle

- The sum of the 2 distances between closest pts. i.e. support vectors should be maxm.
- Both <sup>support</sup> vectors are equidistant from margin.

Again we apply feature scaling. It only improves the performance. So never bad to apply it.

But in case of Tree based models like decision tree, random forest & gradient boosting that don't rely on distances feature scaling has no impact.

Don't apply feature scaling on:

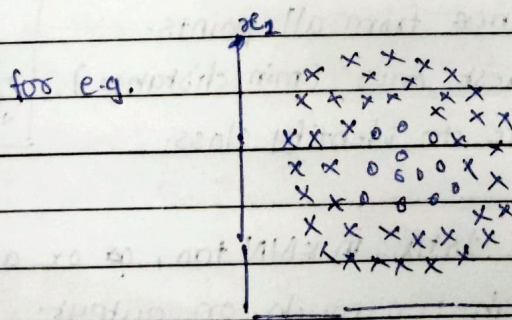
① Tree models

② Categorical Data/sparse data.

Also consider outlier sensitivities, loss of interpretability, & algo tuned for original data like naive bayes.

### ④ Kernel SVM

- Problem when data is not linearly separable.



non linear data

- Way to get this around is to map points in higher dimension & get linearly separable data set, then invoke svm algo, build a decision



boundary for dataset & project all back to original dimensions.

## \* Mapping To Higher Dimension.

### • 1D example.

not separable. Dot is separator in 1D.

Line in 2D

Hyperplane in 3D.

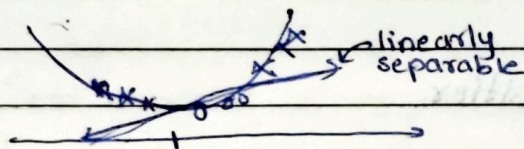
xxx ooo xxx  $\rightarrow x_1$

### - Mapping function example

$$f_1 = x - 5$$

$$f_2 = (x - 5)^2$$

xxx ooo xxx  $\rightarrow x_1$



### - Same with higher dimensional data.

e.g. - 2D  $\rightarrow$  project in 3D applying some fn.

$$\varphi(x_1, x_2) \Rightarrow (x_1, x_2, z) \begin{matrix} \text{new} \\ \text{dimension} \end{matrix}$$

2D 3D space

Project it back to 2D

# Problem: ① Computationally ~~Expensive~~ Intensive  
going to higher dimension then back to lower dimension.

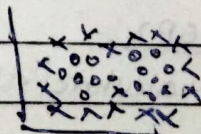
So we use Kernel Trick.

### - Gaussian RBF Kernel

$$K(\vec{x}, \vec{x}_i) = e^{-\frac{\|\vec{x} - \vec{x}_i\|^2}{2\sigma^2}}$$

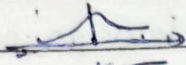
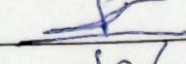
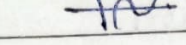
- increase  $\sigma$ , the circumference in 2D increases & vice versa.

Can be added too.  $K(\vec{x}, \vec{x}^1) + K(\vec{x}, \vec{x}^2) \Rightarrow$





## Types of Kernel functions

- (1) RBF 
- (2) Sigmoid. 
- (3) Polynomial 

## (5) Naive Bayes.

### \* Bayes Theorem:

Probability of event A on B  $\swarrow$  ie. B event has been occurred now P(A) needs to happen

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

### \* Naive Bayes classifier.

- It is a probabilistic classifier. algo.
- It predicts probability of data point belonging to a class

- Calculation when adding new data point.

- STEPS 1:
- (1) Prior Probability:  $P(\text{class}) \leftarrow$  for each class No. of class / Total data points
  - (2) Marginal likelihood:  $P(x) \leftarrow x$  is feature set.
  - (3) Likelihood:  $P(\text{class} | x)$
  - (4) Posterior Probability:  $P(\text{class} | x) \leftarrow$  To calculate

Bayes Theorem:  $P(\text{class} | x) = \frac{P(x | \text{class}) \times P(\text{class})}{P(x)}$

where  $P(\text{class}) = \frac{\text{Prob Total class members}}{\text{Total data points}}$

$P(x) = \frac{\text{No. of observations around datapoint (radius is chosen)}}{\text{Total observations}}$

$P(x | \text{class}) = \frac{\text{No. of similar obs. of class around datapoint}}{\text{Total no. of class observations}}$

### STEP 2:

Calculate Probability that feature  $x$  belongs to a class by Bayes theorem.



### STEP 3:

Similarly calculate probabilities for all classes.

### STEP 4:

compare all probabilities. The maximum probability is where new <sup>data pt</sup> ~~class~~ belongs by default.

- It works on assumptions (independent assumptions <sup>from</sup> some underlying data), thus its called Naïve.

## ⑥ Decision Tree

- Looks for optimal splits to maximize pts.
- Same as regression but we predict class here. rather than averages.

## ⑦ Random Forest

- ① pick random  $K$  data pts from train set.
- ② Build DTree on it
- ③ Build  $N$  trees following step ① & ② repeatedly
- ④ During prediction, predict with each of decision tree & take majority vote.

## # Model Selections.

### ① Confusion Matrix & Accuracy.

		Prediction	
		Post	NEG-
Actual	Post	TP	FN
	NEG-	FP	TN

$$\text{Accuracy} = \frac{TP + TN}{TP + FN + FP + TN}$$

To remember, Imagine for disease

Type I error - False +ve  
Type II error - False -ve  
[Just for purpose of remember, Doesn't actually mean anything]

No problem, doctor checks.  
Deadly as doctor won't check.



## # Evaluating Classification Models Performance:

### \* Accuracy paradox:

- More Large data.

- Less no. of +ve data but -  $TP < FP$
- More no. of -ve data and -  $TN \gg \gg FN$ .
- According to accuracy metric model will be good, accurate but actually for +ve's its performing worse

### \* Better Metrics:

\* CAP Curve (Cumulative Accuracy Profile).

\* ROC (Receiver Operating Curve)

Which model to use?

- Identify problem  $\begin{cases} \text{If} \rightarrow \text{linear [Logistic, SVM, etc]} \\ \text{or} \downarrow \text{non-linear [KNN, Naive Bayes, Tree, Forest]} \end{cases}$

- Logistic or Naive Bayes to rank <sup>Prediction</sup> ~~prob~~ by probability
- SVM to predict segment/group.
- Tree for clear interpretation.
- Forest for high performance on cost of less interpretation



## SUMMARY

### ① Logistic Regression

- Binary classifier with output as probability, belonging to discrete class
- Uses MLE to estimate coeff
- Need to scale data before training model.

### ② K-NN

- Works on euclidean distance with K nearest neighbors
- Need to rescale data

### ③ SVM (kernel = 'linear')

- Closest pts are support vectors that are equidistant from Maxm margin hyperplane. & sum of distances b/w is maxm.
- Need to rescale data.

### ④ Kernel SVM

- For nonlinear data, Map to higher dimension & then project back to original dimension.
- Kernel example = 'RBF'
- Need to rescale data

### ⑤ Naive Bayes

- Probabilistic classifier predicting probability of pt. belonging <sup>to particular</sup> class
- Rescale not really needed unless theres extreme difference of range in data. (Avoid scaling for Bernoulli & Multinomial).

### ⑥ Decision Trees & Random Forest

- rescaling data has little to no impact.