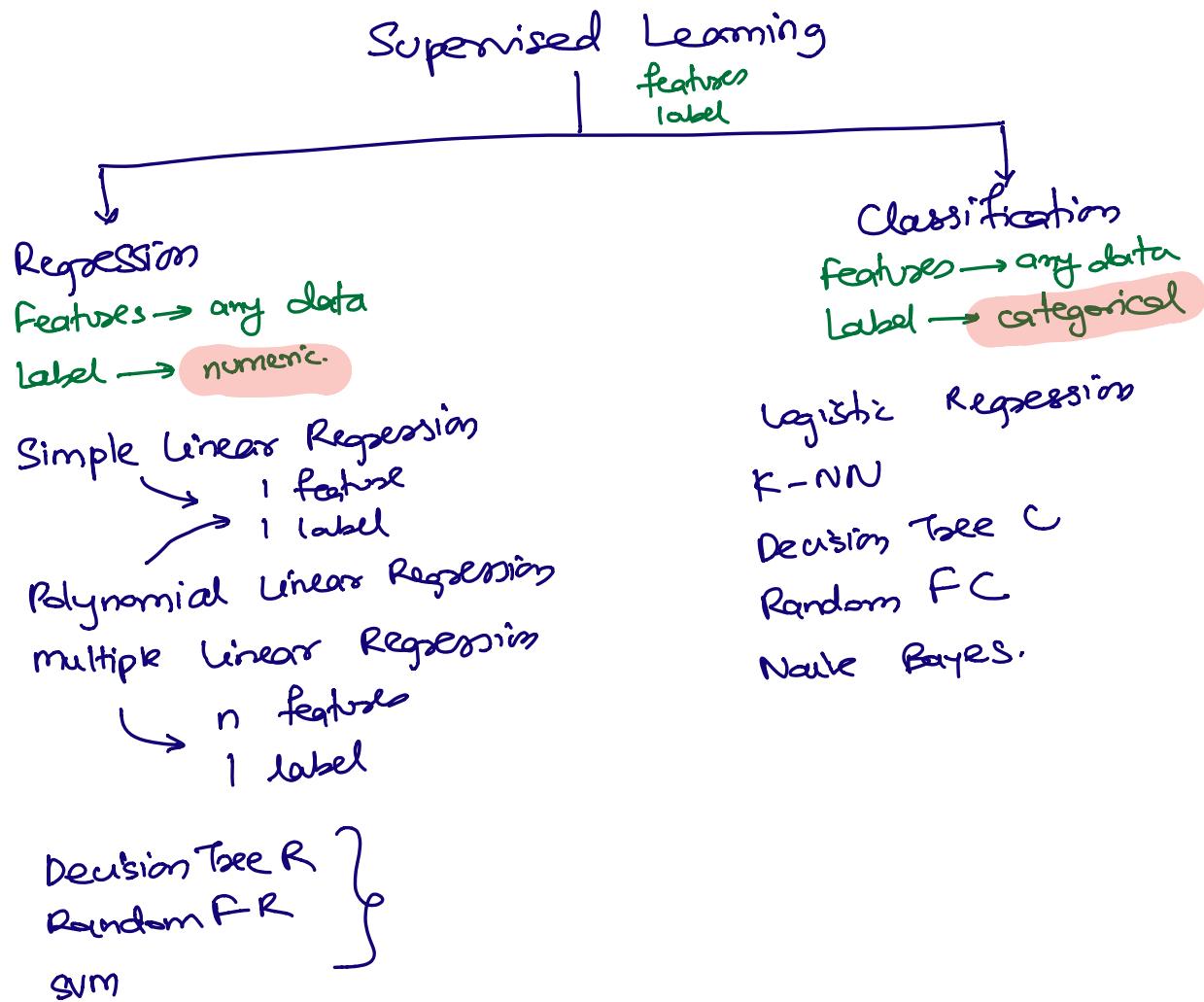


Classification Basics

Approaching a Classification Problem

Logistic Regression

K-Nearest Neighbourhood.



Classification
dealing with categorical labels
in terms of grouping the data

Binary classification
label will be having
two unique data points

Yes | No

0 | 1

spam | ham

multi-class classification.
label has more than
two unique data points

setosa | versicolor | virginica
Primary | Social | Updates | Promotions

K-NN

Naive Bayes

DTC

RFC

SVC

Logistic Regression

K-NN

Naive Bayes

DTC

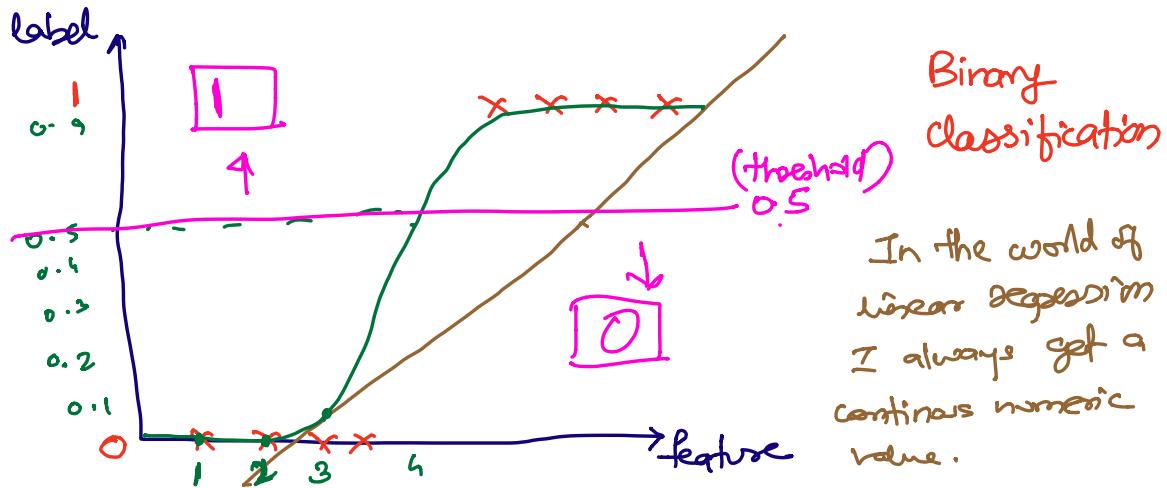
RFC

SVC

Logistic Regression

Linear Regression & Sigmoid

misleading term.



$f^n \rightarrow$ sigmoid function - or -
sigmoid curve

$$S = \frac{1}{1 + e^{-y}}$$

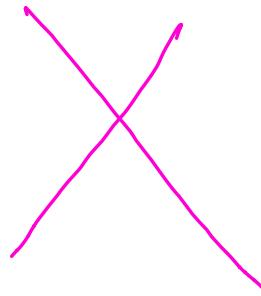
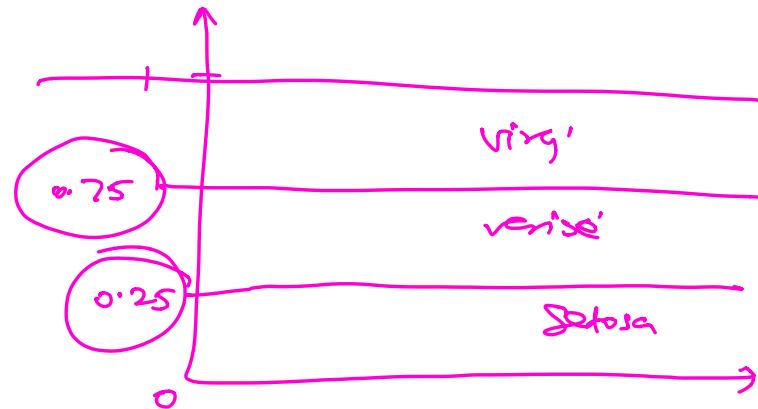
\therefore Let y be the linear regression equation.

$$S = \frac{1}{1 + e^{-(b_0 + b_1(\text{feature}))}}$$

\hookrightarrow Logistic Regression

The output of sigmoid curve is always a probability value ranging from 0 to 1

for multi-class LR

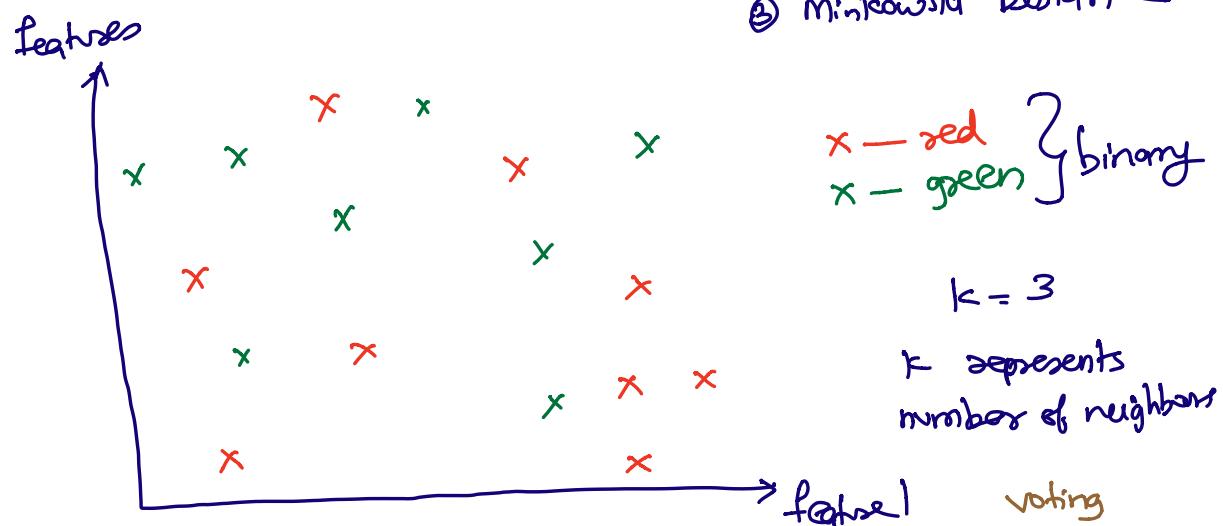


K - Nearest Neighbourhood

↓
uses Euclidean dist as
baseline for distance calc.

↳ Nearest Neighbourhood.

- ↳ Distance formula
- ① Euclidean Distance
 - ② Manhattan Distance
 - ③ Minkowski Distance



Algo:-

- ① Define the value of K. (Choose only odd numbers)
- ② Get K nearest data points from the unknown point.
- ③ Perform voting and select the majority class.
- ④ Reassign majority class as final class output for the given feature.

EVALUATION METRICS IN SUPERVISED LEARNING

How we approach a regression problem?

- ① Label is numeric
- ② features and label must be a 2D array
- ③ Ensure our model is a generalized model

$$\text{Score}_{(\text{train})} < \text{Score}_{(\text{test})}$$

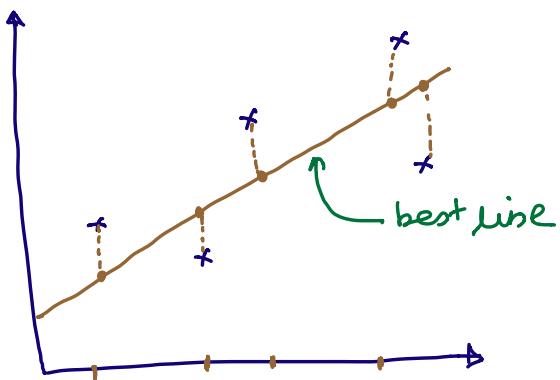
The model is generalized

- ④ Both features and label must be strictly numeric before initiating model training.

⑤ Evaluate your model.

① check MSE / MAE / RmSE (Root mean square error)

(mean squared error) (mean absolute error)

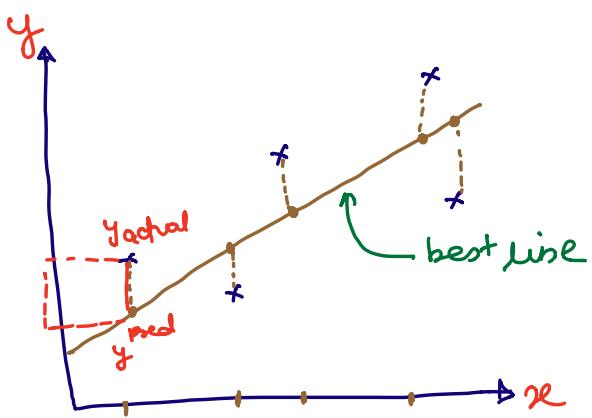


motto of linear regression model is to get the best line.

which line is the best line?
Any line that has less or no error.

Calculate the error of best line?

MSE (mean Squared Error) ↓



$$\text{error} = (y_{\text{actual}} - y_{\text{pred}})^2 \quad \text{for each pts}$$

$$MSE = \frac{\sum_{i=1}^n (y_{\text{actual}} - y_{\text{pred}})^2}{n}$$

n is no of data pts.

To get a line with minimal or no error.

⑥ Check score of generalized model. If convinced with score, perform deployment else improve.

How we approach a classification problem?

- ① Check whether the given problem is a **binary class** or **multi-class** problem.

eg. `data.Purchased.unique()`

- ② Check whether the given dataset is a balanced dataset or not. (which metric to follow for model evaluation)



by checking the value counts of your label column.

eg: `iris.csv`

↓	setosa	↓	50	}	Balanced dataset
	versicolor		50		
	virginica		50		

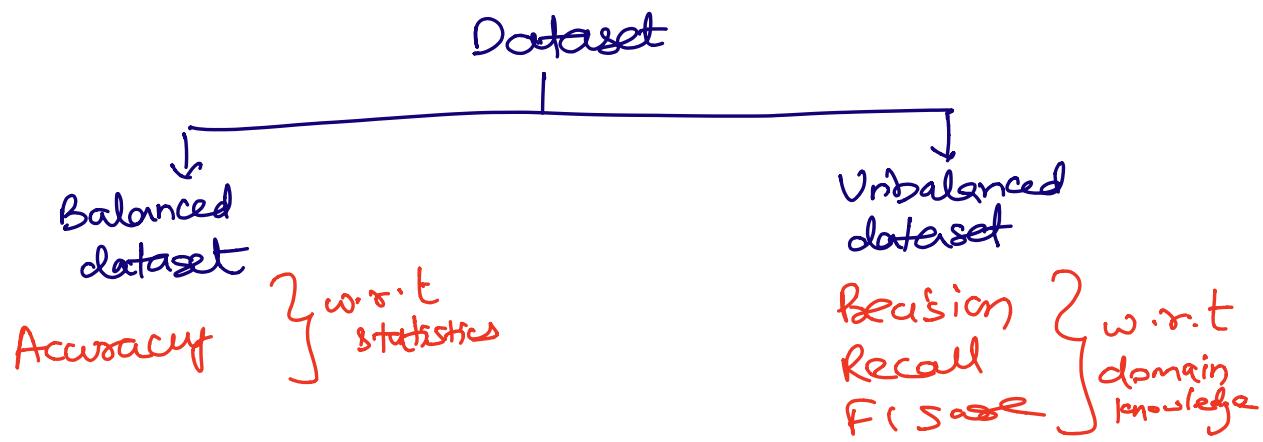
- ③ Ensure features are 2d array and label is 1d array

EVALUATION METRICS FOR CLASSIFICATION.

- ① Accuracy
- ② Precision
- ③ Recall
- ④ F1 score



Evaluate your classification model



CONFUSION MATRIX

Binary classification → 0 & 1

Predicted values

		Predicted values	
		0	1
Actual values	0	Correct 50 Predictions	wrong 20 Predictions
	1	wrong 40 Prediction	Correct 40 Predictions

Accuracy = $\frac{\text{Correct Pred}}{\text{Total Preds}}$
 $= \frac{50+40}{50+20+40+40}$
 $= 0.6$

Precision → horizontal
 $(0) = \frac{50}{50+20} = 0.71$
 $(1) = \frac{40}{40+40} = 0.5$

Recall \rightarrow vertical
 (each label)

$$\text{Recall}_{(0)} = \frac{50}{50+40} = 0.55$$

$$\text{Recall}_{(1)} = \frac{45}{20+40} = 0.6$$

F1 score \rightarrow

		predicted		
		Setosa	Versi	Vir
Setosa		50	0	0
Versicolor	0	45	5	
	0	1	49	

$$\begin{aligned} \text{Acc} &= \frac{50+45+49}{50+45+5+49+1} \\ &= 0.96 \\ &= \underline{\underline{ }} \end{aligned}$$

Decision

$$P_{\text{setosa}} = \frac{50}{50} = 1$$

$$P_{\text{versicolor}} = \frac{45}{45+5} = 0.9$$

$$P_{\text{virginica}} = \frac{49}{49+1} = 0.98$$

Recall

$$R_{\text{setosa}} = 1$$

$$R_{\text{versi}} = 0.9$$

$$R_{\text{vir}} = 0.9$$

Email Spam Classification

h - ham

s - spam

dataset is unbalanced!

		h	s
h	10	100	
s	5	20	

spam \rightarrow inbox ✓
 \rightarrow ham \rightarrow spam ✗

$$\text{accuracy} = 0.27$$

$$\text{Precision}_{(h)} = 0.09$$

$$\text{Precision}_{(s)} = 0.85$$

$$\text{Recall}_{(h)} = 0.66$$

$$\text{Recall}_{(s)} = 0.23$$