

Classification

dealing with categorical labels in terms of grouping the data

PREDICT → The given features belongs to which group

features : anything
Label : categorical

BINARY CLASSIFICATION

label will have only two unique data points

Yes | No

0 | 1

True | False

spam | ham

MULTI-CLASS CLASSIFICATION

label will have more than two unique points.

setosa | versicolor | virginica

Primary | Social | Updates | Random

Logistic Regression ←

K-NN

Naive Bayes

Decision Tree Classifier

Random Forest Classifier

Support Vector Classifier

K-NN

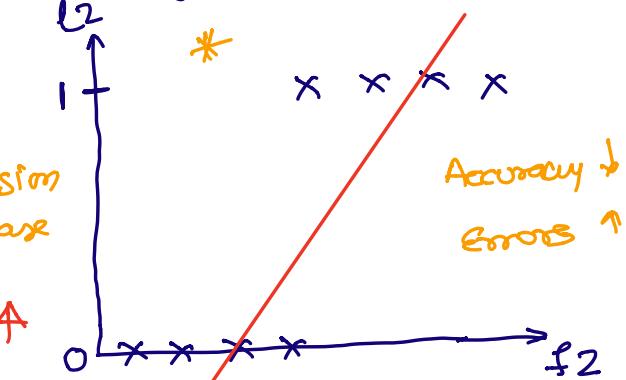
Naive Bayes

Decision Tree Classifier

Random Forest Classifier

Support Vector Classifier

Logistic Regression Classification algo. (Binary classification)



- * f^n that will return probability
- * ① Sigmoid f^n
- ② Bayes Theorem

The motto of the experiment was to modify the Linear Regression formula such that the outcome is a $\overset{(0-1)}{\text{probability value}}$ and based on probability value we can define the threshold for 0 & 1

Logistic Regression = Linear Regression applied over Sigmoid function.

$$S = \frac{1}{1 + e^{-y}} \quad \xrightarrow{\text{Sigmoid } f^n}$$



$$y = b_0 + b_1 f_2 \quad \text{—— Linear Regression formula} \quad \text{—— ①}$$

Substituting y in eqn ①

$$S = \frac{1}{1 + e^{-(b_0 + b_1 f_2)}} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Logistic Regression equation}$$

Applied ML \rightarrow Logistic Regression

- ① `sklearn.linear_model ... LogisticRegression`
- ② Applicable only for binary classification

Applied ML \rightarrow Classification.

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

e.g. `data.label.value_counts()`

- ② Check whether the given dataset is a balanced dataset or unbalanced dataset.

classification dataset

Balanced
your label data points
will have equal no. of data

Unbalanced
your label data points
will have an unequal amt of data.

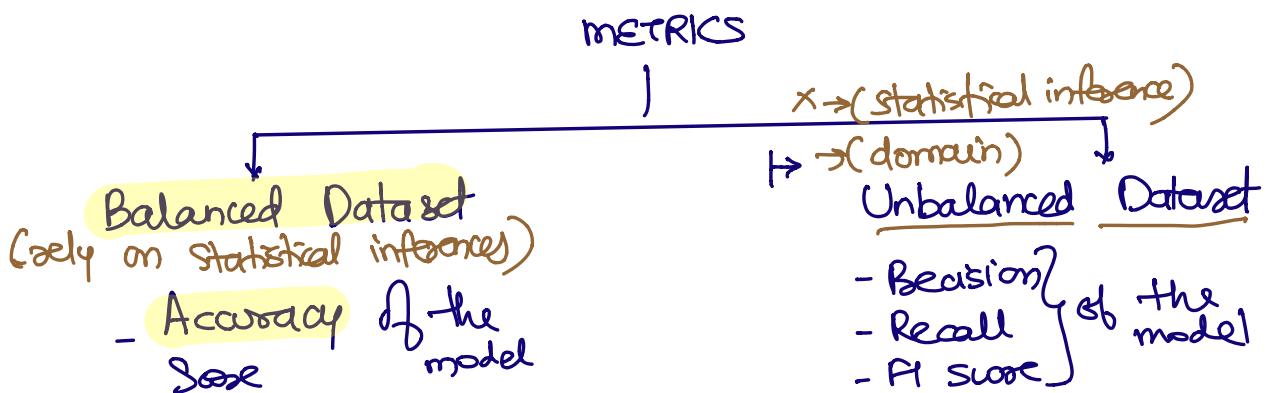
`iris.csv` → 150 observation / data points

Setosa → 50	}	→ <u>BALANCED DATASET.</u>
Versicolor → 50		
Virginica → 50		

Setosa → 51	}	→ <u>UNBALANCED DATASET</u>
Versicolor → 49		
Virginica → 50		

The reason we need to know whether the dataset is balanced or not is,

- ⇒ ① To check the quality of the model
 - * ② To decide the metric for feature optimization
- MODEL **QUALITY**



Note:- The rule for **generalization** remains same in classification. ($\frac{\text{test accuracy score}}{\text{train accuracy score}} > 1$)

EVALUATION METRICS FOR CLASSIFICATION

Confusion matrix \longrightarrow ① Accuracy
 CM is a table used to ② Precision
 describe the performance ③ Recall
 of the model. ④ F1 - Score

example of confusion matrix (iris.csv)

Balanced
 $S : 50\%$
 $Ve : 50\%$
 $Vi : 50\%$

Dataset used: iris.csv
 features: sepal.length, sepal.width, petal.length, petalwidth
 label: species

		Predicted			$y_{pred} \leftarrow$ predicted by model
		Setosa	Versicolor	Virginica	
Actual values y_{true}	Setosa	50	0	0	Support $\Rightarrow 50\downarrow$ $\Rightarrow 50\downarrow$ $\Rightarrow 50\downarrow$
	Versicolor	0	45	5	
	Virginica	0	1	49	

Accuracy = $\frac{50 + 45 + 49}{50+0+0+45+5+0+1+49} = \frac{144}{150} = 0.96$

(applicable for the entire model)

[Go Diagonal] $\underline{\underline{96\%}}$ \leftarrow

Precision → Precision is applicable for each label data point

{ Go VERTICAL }

	S	Ve	Vi
→ S	50	0	0
→ Ve	0	45	5
→ Vi	0	1	49

$$\text{Precision}_{(\text{setosa})} = \frac{50}{50+0+0} = 1 \approx 100\%$$

$$\text{Precision}_{(\text{versicolor})} = \frac{45}{45+1+0} = 0.97 \approx 97\%$$

$$\text{Precision}_{(\text{virginica})} = \frac{49}{0+5+49} = 0.9 = 90\%$$

Recall → is applicable for each data point

{ Go HORIZONTAL }

$$\text{Recall}_{(\text{Setosa})} = \frac{50}{50+0+0} = 100\%$$

$$\text{Recall}_{(\text{versicolor})} = \frac{45}{0+45+5} = 90\%$$

$$\text{Recall}_{(\text{virginica})} = \frac{49}{0+1+49} \approx 98\%$$

	S	Ve	Vi
S	50	0	0
Ve	0	45	5
Vi	0	1	49

example use - Case

Email spam classifier.

		ham	spam
ham	ham	10	2
	spam	5	30

① Is the given data set balanced or unbalanced?

Ans: Support of each label

$$\text{ham} = 10 + 100 = 110$$

$$\text{spam} = 5 + 30 = 35$$

$$\text{Total data pts} = 145 //$$

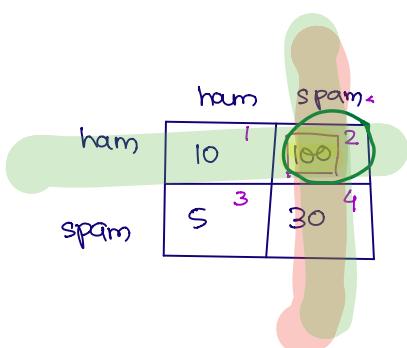
∴ given dataset is unbalanced
since $n(\text{ham}) \neq n(\text{spam})$

② What is the expectation of client in terms of error tolerance in domain perspective?

Spam → Inbox [mark as spam]

⇒ ham → Spam [☹]

↑ not tolerable (domain perspective)



$$P_{\text{ham}} \rightarrow \frac{10}{10+100} = 0.09 \approx 9\%$$

$$P_{\text{spam}} \rightarrow \frac{30}{100+30} = 0.23 \approx 23\%$$

use-case 2 → Medical diagnostics use-case

Role of this model is to predict whether the patient is sick or healthy.

		sick	healthy
Sick	True	50	1050
	False	0	30000

$$\text{Recall}_{\text{Sick}} \rightarrow \frac{50}{1050} = 0.05 = 5\%$$

$$\text{Precision}_{\text{Healthy}} \rightarrow \frac{30000}{31000} = 0.96 \approx 96\%$$

- ① Is the data balanced or not? → Unbalanced
- ② Binary / multiclass
- ③ To check generalization what need to be performed?
→ Raise a query of generalization check to DS.

The model is generalized!

- ④ Please check the quality of the model and send the same for approval.
Ask for error tolerance
- ⑤ Hospital is saying following point:

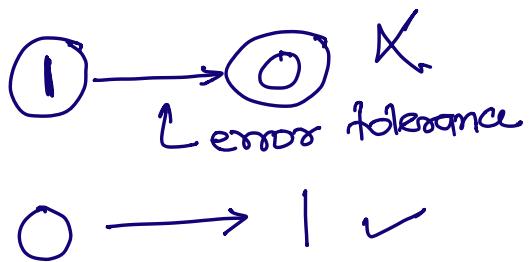
Healthy → Sick 😊 Business

Sick → Healthy 😕

↑ not tolerable

Shopping mall dataset

	0	1
0	238	19
1	43	100



$$\text{Precision}_{(0)} = 0.84$$

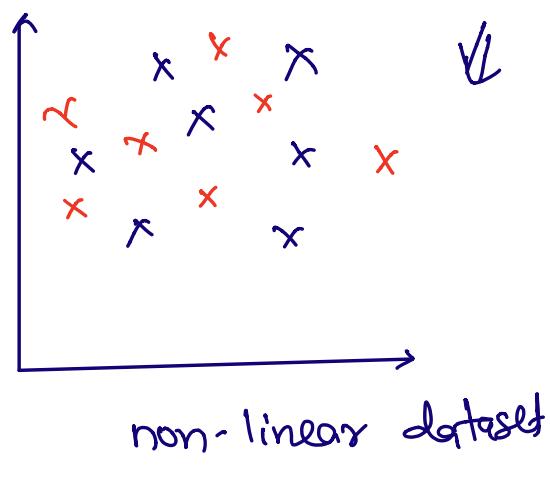
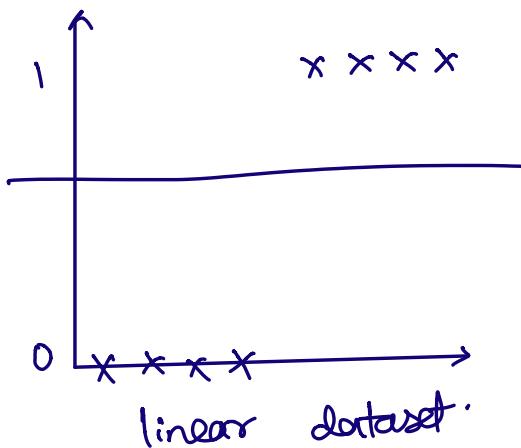
$$\text{Recall}_{(1)} = \frac{100}{143} = 0.69$$

K-Nearest Neighbourhood algorithm

Expectations of most of the classifn algs:

① Your data must be linearly separable

② Your use-case must be binary classifn (Logistic Regression)



K-NN is BEST when it comes to dealing non-linear separable datasets!

Algo:-

- ① Define the value of k. ($k \rightarrow$ no of neighbour)
- ② Get k nearest data points from unknown data pts.
- ③ Perform voting and select the majority class.
- ④ Reassign majority class as final class opp for given feature.

$k=5$
 $12 \text{ data} \rightarrow \text{sort} \rightarrow \text{head}(5)$

