# * Classification *

⊙ Classification:-

> Categorize given data set into classes, can operate on both structured & unstructured data.

↳ Goal – Assign i/p data points to predefined classes.

↳ It is under Supervised Learning.

↳ classes – target, label / categories.

⊙ Learners in Classification Problems:-

① Lazy Learner –
↳ Firstly, store training dataset & wait until receive test dataset.
↳ Classification done on the basis of most related data stored in training dataset.
↳ Less time training & more time in prediction.
↳ Algorithms: KNN algo, case-based reasoning.

② Eager Learner –
↳ Develop classification model based on training dataset, before receiving testing dataset.
↳ More time learning & less time in prediction.
↳ Algorithms – Decision Tree, Naive Bays, ANN.

# Types of Classification Algorithms :-

↳ Supervised Learning Algos, used to categorize data based on i/p provided.

↳ Most common problems — speech recognition, Image recognition, text from handwriting.

↳ Classification algos divided in 2 categories:

- Linear Modes —

  ① Logistic Regression
  ② Support Vector Machine

- Non-linear Models —

  ① Artificial Neural Network (ANN)
  ② Random Forest
  ③ Decision Tree
  ④ K-Nearest Neighbor (KNN)
  ⑤ Naive Bayes
  ⑥ Stochastic Gradient Descent

- Terminologies used in Classification in ML —

① Classifier — Model used to map i/p data to specific category.

② Classification model — Model used to predict/ draw conclusion to i/p data given for training, predict class/category for data.

③ Feature — Individual property of the dataset being observed.

④ **Initialize** — Assign classifier to be used for classification.

⑤ **Train** — Train the model using fit method. Using train_x & train_y dataset.

⑥ **Predict** — Use model to predict o/p for new set of i/p data.

⑦ **Evaluate** — Evaluation of performance of model.

● **Types of Classification :-**

1) **Binary classification** —

↳ Categorize i/p data to 1 of 2 possible classes or categories. Ex : True/False, Yes/No, 0/1.

↳ **Example** — 
- Email spam detection (spam or not).
- Churn prediction (churn / not)
- Conversion prediction (buy / not)
- Rain forecast (Yes/No).

↳ 2 classes — 1 is **normal** state & other **abnormal**.

           ↗                              ↓

    Assigned 0                          Assigned 1.

↳ **Ex :** 
- Email spam — "spam" is abnormal state.
- Medical diagnosis — "cancer detected" abnormal.

↳ **Algos used** — 
① Logistic Regression
② k-Nearest Neighbors
③ Decision Tree
④ Support Vector Machine
⑤ Naive Bayes

## 2) Multiclass Classification —

↳ Has more than 2 classes/categories.

↳ Examples —
- Face Classification.
- Plant species classification.
- Optical Character recognition.

↳ Has range of classes to predict.

↳ Algos used —
① k-Nearest Neighbors.
② Decision Tree
③ Naive Bayes
④ Random Forest
⑤ Gradient Boosting

↳ Each sample is assigned to only 1 label/target.

## 3) Multi-Label Classification —

↳ O/p have 2/more class labels, where one/more labels may be assigned to each example.

↳ Example —
- Photo classification- Have many objects like - man, bicycle, apple, tree, banana.

↳ Multi-label versions of Algos Used —
① Multi-label Decision Tree
② Multi-label Random Forests
③ Multi-label Gradient Boosting

# 4) Imbalanced Classification —

↳ Num of examples in each class is unequally distributed. i.e. one type has ~~some~~ more examples than others.

↳ Handled using Random Forests.

↳ Typically are of binary classification, where normal class is in majority & abnormal minority.

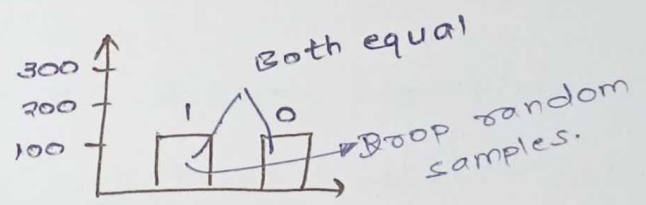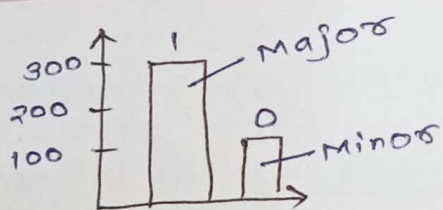- Techniques used for Balancing —
  - ① Under Sampling
  - ② Over Sampling
  - ③ Bagging / Boosting

① Under Sampling —
↳ Minority class remain as it is.
↳ Reduce majority to match minority.

```
from imblearn.under_sampling import NearMiss
```
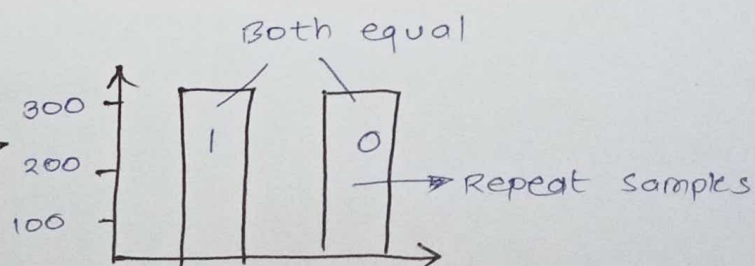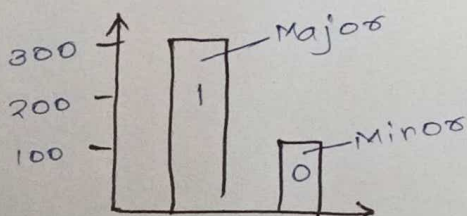


Both equal

Drop random samples.

↳ Disadvantage — Possible to lose imp. data.

② Over-Sampling — SMOT library used.
↳ Repeat features to ↑ minority class to match with majority class.

```
from imblearn.over_sampling import SMOT
```



Both equal

Repeat samples

## 5) Hierarchical classification —

↳ Classes ordered in hierarchy / tree-like structure.

↳ Suitable for hierarchical relationship bet^n classes.

↳ Example: species classification in taxanomy —
Kingdom, Phylum, Class, Order, Family, Genus & Species.

## 6) Ordinal Classification —

↳ Data where classes have natural ordering / ranking.

↳ Example — Movies rating: 1-star, 2-star etc.

## 7) Multiout Classification —

↳ Also k/a multiclass - multiout classification.

↳ Predict multiple o/p variables for each i/p data point.

↳ Each o/p variable can be binary / multiclass.

↳ It is extension of multilabel classification.

# ⊙ Model Evaluation Techniques for Classification :-
Used to find performance & efficiency of models.

## ① Confusion Matrix :

↳ Detailed view of model's performance by showing counts of true +ve (TP), True -ve (TN), False +ve (FP) & False -ve (FN) predictions.

↳ Used to derive other metrics, like - Precision, Recall, F1-score & specificity.

↳ Useful for multi-class variables.

Example -

| Y_actual | Y-pred |
|----------|--------|
| 0 | 1 |
| 1 | 1 |
| 0 | 0 |
| 1 | 1 |
| 1 | 1 |
| 0 | 1 |
| 1 | 0 |

| | | Actual | |
|---|---|---|---|
| | | True | False |
| | | 3 | 2 |
| Predicted | True | True +ve (TP) | False +ve (FP) |
| | False | False -ve (FN) | True -ve (TN) |

## ② Accuracy — Measure performance of the model.
↳ Measure how often model is correct.

↳ Ratio of total correct pred$^n$ to total Pred$^n$.

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

TP + TN → correct Predictions

TP + TN + FP + FN → Total Predictions

$$= \frac{4}{7} = \underline{0.57} \qquad \therefore \ 57\% \ Accuracy.$$

③ **Classification Error** — Also k/a <u>Misclassification Rate.</u>

↳ Measure how often classifier is incorrect.

$$Error = 1 - accuracy$$

④ **Precision** — How accurate model's +ve predictions are.

↳ How many predicted True are correct.

↳ Ratio of True +ve to total no. of +ve.

$$Precision = \frac{TP}{TP+FP} = \frac{3}{3+2} = \frac{3}{5} = 0.60$$

∴ 60% Precision.

↳ Goal is to reduce FP.

⑤ **Recall** — How many actual True correctly predicted.

↳ Ratio of TP to sum (TP + FN).

$$Recall = \frac{TP}{TP+FN} = \frac{3}{4} = 0.75$$

∴ 75% Recall.

↳ Goal — Reduce FN. (Health related Problems).

• <u>Mail Classifier</u> —

Good { Actual - Spam
       Pred - Spam    [TP]

May happen { Actual - Spam    [FN]
OK.        { Pred - Not

worst { Actual - ~~Spam~~ Not    [FP]
case  { Pred - ~~No~~ Spam

| 1 | 0 |
|---|---|
| [FP↓] | |
| | |

∴ <u>Precision</u>

• <u>Diabetes classifier</u> —

Both { Actual - Yes    Actual - No
Good { Pred - Yes      Pred - No
      [TP]             [TN]

[FP] Actual - No  } May be okay.
     Pred - Yes

[FN] Actual - Yes } worst case.
     Pred - No

| 1 | 0 |
|---|---|
| | |
| [FN↓] | |

∴ Recall.

**(6) F1- Score** — Evaluate overall performance of model. Harmonic mean of Precision & Recall.

$$F1. score = \frac{2 * Precision * Recall}{Precision + Recall}$$

$$F1 = \frac{2 * 60 * 75}{60 + 75} = 66.67$$

**(7) AUR- ROC Curve —**

- **ROC Curve** — Receiver Operating Characteristics.
- **AUC** — Area Under Curve of the ROC Plot.

↳ Probability graph to show performance of classification model at different threshold levels.

↳ Curve plotted bet$^n$ 2 params —
① TPR (True +ve Rate)
② FPR (False +ve Rate)

① **TPR** — Same as Recall. $\dfrac{TP}{TP+FN}$

② **FPR** — $\dfrac{FP}{FP+TN}$