

* KNN =

- K-Nearest Neighbour
- supervised learning approach.
- assumes similarity b/w new case/data & available test case & put new case into category that is most similar to categories available.
- Regression & Classification (more preferred)
- lazy learner algo → not learning from training set immediately
→ stores dataset & time of classification → action.

Adv → simple to implement
robust → noisy training data
more effective → training data → large

Disadv → • need to determine K value
• computation cost high.

Working →

$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

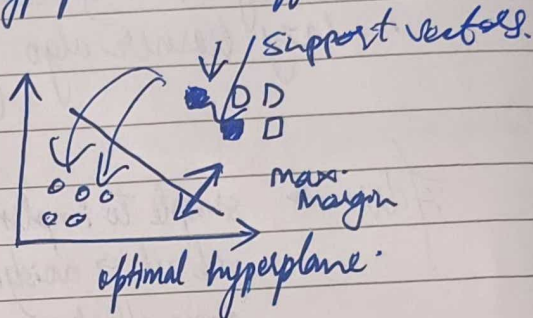
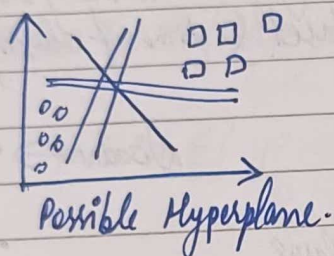
- select no. of K of neighbors
- Calculate Euclidean distance of K no. of neighbors
- Take K nearest neighbors as per calculated Euclidean distance
- K-neighbors, count no. of data point in each category
- Assign new data point to category which no. neighbor is max.
- Our model is ready.

Why KNN → discovery category or class of selected dataset without difficulty

SVM

- Support Vector Machine
- supervised learning
- classification & regression problem
- (more)

- boundary data points to create a hyperplane that differentiate data points.



Types → Linear SVM → • linearly separable info.
• data → labelled two classes → unmarked directly line.

Non Linear SVM → • non linearly separated facts.

- dataset → can't be labeled through usage of directly line
- fact → non-linear info
- classifier → non-linear SVM classifier.

Hyperplanes → decision boundary that help in classifying datapoints

Support vectors → datapoints that are closest to hyperplane & influence posⁿ & orientation of hyperplane.

* applⁿ SVM → Image classificⁿ - Handwritten character recogn

Limitation → • sensitive to noise.

- lies in choice of kernel
- optimal design for multiclass SVM classifier

Bagging

Goal → Reduce Variance

Learning → Parallel

more memory use

Each receives equal weight.

solves overfitting problem.

• used if classifier is unstable

eg → Random Forest

- combine homogeneous weak learners.
- aggregating predictions from weak learners

Boosting

Reduce Bias.

Sequential.

less memory usage

models are weighted acc to their performance.

only reduce bias

• Used if classifier is stable & simple

eg → AdaBoost, XGBoost

- weak learners aren't independent.
- weights of misclassified are increased.

* Random Forest-

- made up of multiple decision trees-
- classification & regⁿ problem-
- output of multiple decision trees to reach single result

Adv

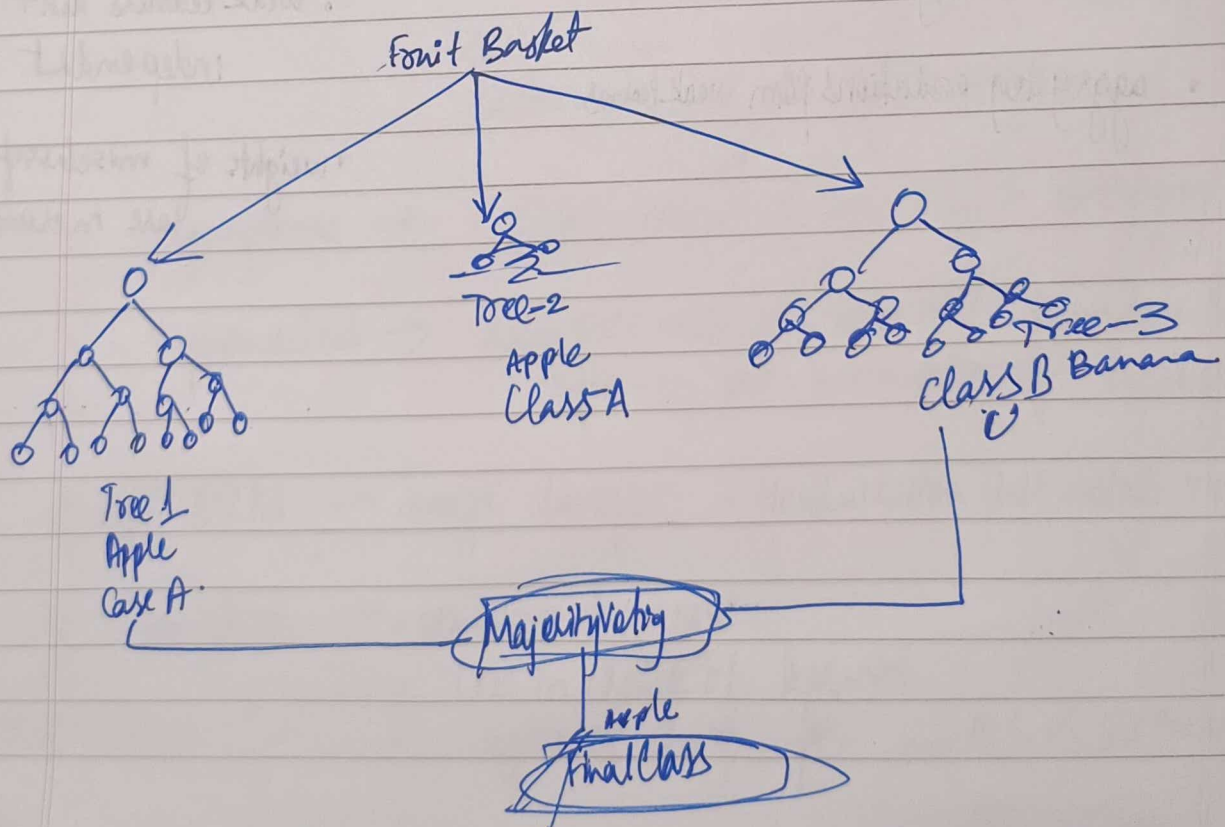
- less training time
- output with high accuracy → large dataset → efficient
- maintain accuracy → large portion of missing data

Disadv

- not more suitable for regression task

Appln →

- Banking
- Medicine
- Land use
- marketing



Binary Classifⁿ

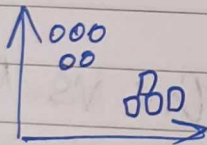
- supervised learning
- It is two types of groups. \rightarrow max two classes
- eg \rightarrow Medical Test \rightarrow disease yes or no

Algo \rightarrow

- "Logistic regⁿ"
- K-nearest neighbours
- SVM
- Naive Bayes

eg \rightarrow

- Spam detection \rightarrow spam yes or no
- Conversion predⁿ \rightarrow buy or no



Multiclass Classifⁿ

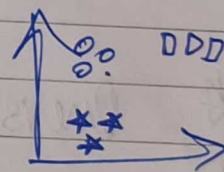
- supervised learning
- May be any variety of classes - many
- eg \rightarrow classify email categories

Algo \rightarrow

- K-Nearest nex
- Decision trees
- Random forest
- Gradient

eg \rightarrow

- plant species classification
- optical character recogⁿ



Balanced classification

one class.

- equal to other class
- ML models work best when no. of samples in each class are about equal.

Imbalanced classifⁿ.

- If observations in one class is way higher than other class
- may not be good ml model

Techniques

Cause \rightarrow 1) Sample Error
2) Problem Domain.

Technique to handle class imbalance -

- Random Resampling -
removing samples from majority class & add more in minority
- Tomek Links -
 - Majority class of each pair removed.
 - Link pairs of very close instance but opp classes
- SMOTE (Synthetic Minority Oversampling Technique)
 - Generate Synthetic data \rightarrow minority class.
 - Uses KNN.
- Class Weights \rightarrow Prioritizing any particular class

* One vs One

- Heuristic method for binary classification.
- splits multiclassification into binary problem.

eg \rightarrow classify colors \rightarrow R B G.

classifn \rightarrow Binary Problem - 1 \rightarrow red vs blue
 " " 2 \rightarrow red vs green
 " " 3 \rightarrow blue vs green

no. of binary dataset will be created one vs one.

$$\frac{C \times (C-1)}{2} \rightarrow C \rightarrow \text{no. of classes.}$$

argmax sum of score predicted \rightarrow class label

One vs All.

- Heuristic method for using binary classⁿ for multiclassification.

eg \rightarrow R G B V classification

B class Problem 1 \rightarrow R vs [G, B, V]
 2 \rightarrow G vs [R, B, V]
 3 \rightarrow B vs [R, G, V]
 4 \rightarrow V vs [R, G, B]

- one model for one class
- each model predicts \rightarrow class membership probability

argmax score \rightarrow predict class

$$\text{Precision} \rightarrow \frac{TP}{TP+FP}$$

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

* Macro average &

- when all classes needs to be treated equally

micro average

- need to weigh each instance of prediction equally.

* Cross Validation -

- given dataset \rightarrow randomly create training & testing dataset & test out various model to see which works best.

- evaluate model \rightarrow usage of subset of facts get after which compare use of complementary subset of facts set

Two ways \rightarrow 1) Non exhaustive \rightarrow don't split dataset in all ways leave subset of data for valn

2) exhaustive \rightarrow split data into all possible ways into training & validation sets

Methods

1) Holdout Method

- Non exhaustive
- Divide data \rightarrow 70:30 or 80:20
- training/testing
- data is shuffled randomly before splitting

2) K-fold cross valⁿ

- Split data in k parts
- Train data on (k-1) parts & test data on 1 part.
- perform steps k times
- total is $n \times k$

3) Leave P out

- exhaustive approach
- take p parts from total n no. of parts
- Train model on (n-p) points & test on p remaining
- Repeat all possible attributes
- Avg accuracy $\rightarrow \frac{\text{Average}}{P}$