# KNN = Theory

KNN = It's the [Simplest] yet one of the most (Commonly) used M.L. algo.

## Advantages of KNN

→ Easy to implement

→ KNN has [no-model] — No complex function like other ML algorithms such as Linear Regression, Logistic Regression, Decision Tree, Naïve Bayes et

→ It just [stores data], so [no learning] is required.

→ Training is [not required] hence all work is done during prediction.

→ No (prior knowledge) or any knowledge about distribution of data required.

→ Defer data processing untill receives request to classify unlabeled data — Lazy learner.

## Example of application of KNN

KNN is simple but powerful algorithm & many times [outperforms] other algorithms.

## Is my heart functioning well?

⟹ In medical research lots of data are generated.

⟹ KNN is capable of [handling huge amount of data] provided we've sufficient resources,

⟹ And it [doesn't get affected] by outliers.

⟹ So based on [similarity among cases of diagnosis] it can [predict] diseases.

✹ Since it works on basis of similarity b/w instances so it has wide application in [Recommender System.]

→ φ. Similarity b/w (users' ratings, product ratings)

Can efficiently be [Recognized] by KNN algo.

N.B.:- KNN algo ~~~~~ will be talked & its application is discussed. & in Recommender system.

## KNN in a nutshell

→ It's a form of [Supervised learning]

* Data processed by KNN has input features & output feature.

→ Used in both [Classification & Regression]

→ Uses [nearest majority] of neighbors to predict data points.

# How does KNN work?

* **Step 1:** Build a neighborhood with [K-members]
(K = 1 to any no., only [odd NO] because the algo has to choose the heighest No. of rankings)

* **Step 2:** Find [distance] from [query point] (Test data) to each point of neighborhood.

* **Step 3:** [Assign] query point to that class having [maximum members around (Majority voting)].
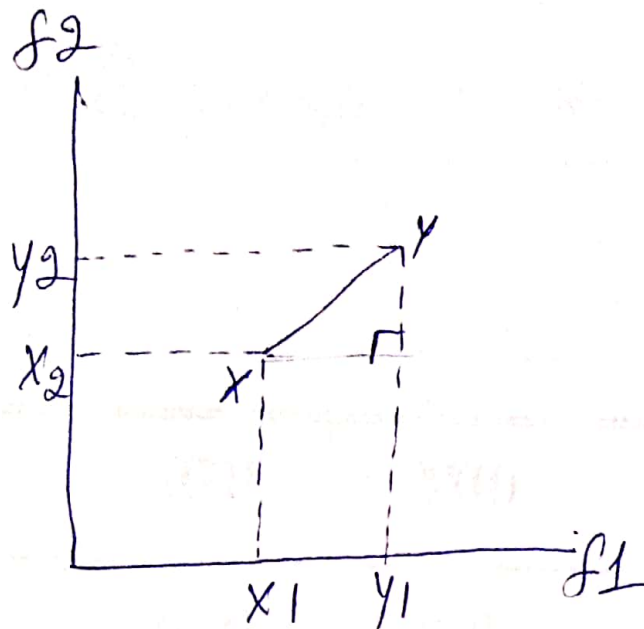
→ Classification is done in KNN based on (majority of neighbors.)

→ where [K] denotes number of neighbors.

---

There are a no. of distance measures for continuous data such as Euclidean, manhattan, Minkowski most common being [Euclidean Distance]

Euclidian distance measures shortest distance b/w two points (X & Y)

$f2$



$$\Rightarrow \sqrt{\sum_{i=1}^{k} (x_i - y_i)^2}$$

$X(X1, X2)$ & $Y(Y1, Y2)$

$f1 = X1, Y1.$ & $f2 = X2, Y2$

$$d = \sqrt{(Y1 - X1)^2 + (Y2 - X2)^2}$$

For Continuous Data (Regression)

( Euclidian distance ) is used [as mentioned above]

For Categorical Data (classification)

( Hamming distance ) is used.

# Hamming distance is measured as below

This is locations where binary vectors (2 classes) differ.

$$D_H = \sum_{i=1}^{k} |x_i - y_i|$$

$$x = y \implies D = 0$$
$$x \neq y \implies D = 1$$

| X | Y | Distance |
|------|--------|----------|
| Male | Male | 0 |
| Male | Female | 1 |

← Sequence of X & Y differs (Male & Female)

A  B  B  A A B B A B A
B  B  A  A A B A B B A

Hamming distance is __4__, because classes differ at 4 locations.

# KNN

① Import packages (Load the data)

②③ Clean data ⎯⎯ ⎡ missing values
⎣ outlier treatment

② EDA ⎯⎯ ⎡ Box plot
⎢ Correlation
⎢ Histogram
⎢ Pairplot
⎣ many others

---

## Step - 1

### Create features and Labels

$X = \ldots$

$y =$

X.shape

y.shap :

If the data is in np. array convert it into pd. Data Frame

$df = pd.DataFrame(*\text{file}*).apply(pd.to\_numeric)$

## Step 2

Splitting the data into Train & Test ~~set~~ set

(Need to verify)

from sklearn.model_selection.import train-test-split

## Step 3    Create instant of the model

Building the K NN Model

from sklearn.neighbors import KNeighborsClassifier

↳ clf_knn = ~~KNN~~ KNeighborsClassifier(n_neigh bors=10

## Step 4 => Fit the model

∅ clf_knn = clf_knn.fit(X_train, y_train)

## Step 5 => Predict the model

y_pred = clf_knn.predict(X_test)

## Step 6 => Evaluate the model

Cross validation is used to see how well a model performs in an independent dataset. Different samples ca. used for training & test & accuracy score is calculated

from sklearn-model-selection import Cross-Val_score

Scre_knn = Cross_Val_score (Clf_knn, X-test, y-test, Cv=5)

→ Print ("Cross Validation Score:" + str (Scre_knn))

Print ("Cross Validation mean Score:" + str (Scre_knn_mean()))

## Classification report

from sklearn.metrics → import Classification_report, Confusion matrix, accuracy_score.

Cr = Classification_report (X-test, y_pred)

pr (Cr)

## Confusion matrix

Cm = Confusion_matrix (y-test, y_pred).

Cm

## Optimization of Neighbours

↳ Use different Counts of neighbours of by setting value of k (5, 7, 9, 11, 13, 15) ⇒ (Any range of Values)

↳ Perform CV (5) on each value of k (k in KNN)

↳ Compare the CV score to deduce the best model.

↳ Use the best value of k.

```
list_k = [5,7,9,11,13,5,] => Always odd No/Even No
                                       is also fine
for -k in list_k:
     clf_knn = KNeighbors Classifier (n_neighbors = k)

     Score_knn = Cross_Val_Score (clf_knn, X_test,
                                 y_test, CV=5)

     Print ("k: ", k)

     Print ("Cross validation Score: " +Str(Score_knn))

     Print ("Cross validation mean Score : " +Str(Score_knn.mean
                                                  ())
     Print ("  ")
```

~~we can~~ we can use (Grid Search to)

get faster result.