# Supervised Learning
- Nearest neighbour
- Naive Bayes
- Decision Trees
- Support Vector Machines (SVM)
- Neural Networks.
- Random Forest
- Similarity Learning
- Linear Regression
- Linear Discriminant analysis

## K - Nearest Neighbour
It is a non-parametric & Lazy algo

Non parametric means it doesnot depend on the data rather depends upon the proximity to other data points ~~regards~~
⇒ It doesnot make assumptions on the underlying data.

Parametric ~

Lazy learning algo implies there is little to no training phase.
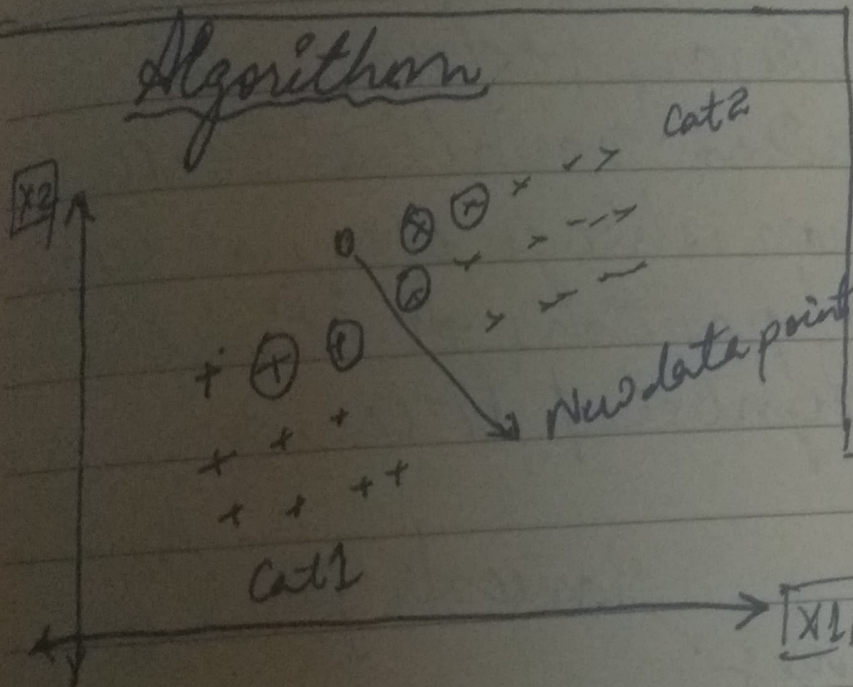⇒ We can immediately ~~train~~ classify new data points as they present themselves

| Pros: | Cons: |
|---|---|
| → No assumption about data date | High memory req? all the training data must be present to in memory in order to calc the closest k-neighbor |
| → Simple & easy to understand | |
| → can be used for classification & regression | → sensitive to irrelevant features |
| | → sensitive to the scale of the data since while computing the distances to the closest k-point. |

## Algorithm



Ⓧ → 3 neighbors
⊕ → 2 neighbors

(i) Pick a value for K. ( say K=5 )

(ii) Take K' nearest neighbors from the data point according to its Euclidean distance

(iii) Among these neighbors count the number of data points in each category & assign the new data point to that category.

**Date:** **CODE :- TUNING OF KNN:-**

→ Trying larger 'k' values to see if we could improve the performance of the algorithm.

→ Try different distance measures like ter : Hamming distance

Hamming distance: calculates the distance between two binary vectors.

In case of categorical variables, it is the difference b/n two strings of equal length i.e. the number of positions at which the corresponding symbols / letters are different.

Manhatten Distance :- → Tanicab geometry

$$d_1(p,q) = \| p - q \| = \sum_{i=1}^{n} |p_i - q_i|$$

↓ the sum of the absolute differences of their cartesian coordinates.

Minkowski Distance → Generalization of both the Euclidean & Manhatten distance

$$d(x,y) = \left( \sum_{i=1}^{n} |x_i - y_i|^p \right)^{\frac{1}{p}}$$

# Curse of Dimensionality

KNN works well with a small number of input variables ($p$), but struggles when the number of inputs is very large.

Each input variable can be considered as each dimension for instance ~~x1 & x2~~ or two $x1$ & $x2$ are two variables.

⟹ the input space would be 2-dimensional

As the number of dimensions increases

⟹ the volume of input space increases

⟹ points may be similar may have large distance

⟹ Curse of Dimensionality.