

K-Nearest Neighbors (KNN)

Overview



Core idea of KNN



Distance Scores



KNN for Classification and Prediction

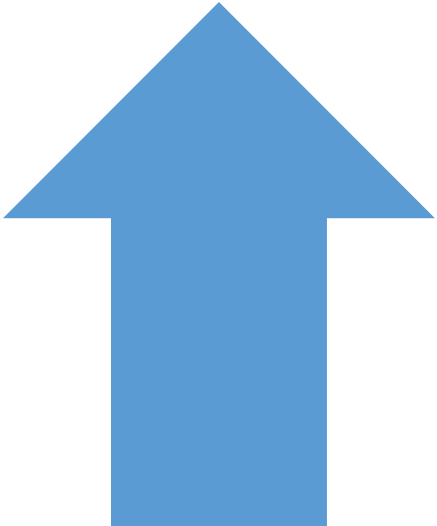


Advantages and Disadvantages of KNN



Summary

Categorical Outcome

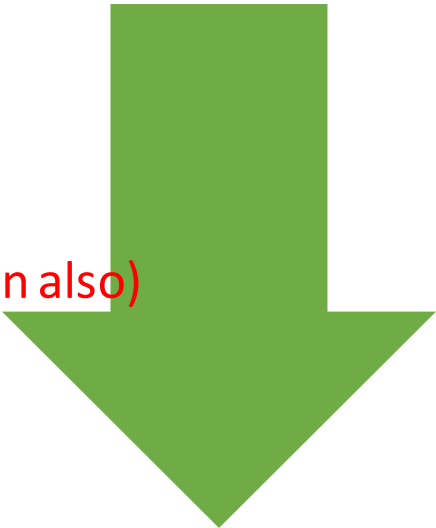


Classification

Pass / Fail
Purchase / No purchase,
Fraudulent / Genuine
Creditworthy / Defaulter
Owner / Non-owner

Prediction

(referred to as regression also)



Total Marks scored
Total Amount spent on purchase
Amount of Fraud
Amount of Default
Net worth

Numerical Outcome

Characteristics of KNN Algorithm

- KNN can be used for **classification** (categorical outcome) as well as **prediction** (numerical outcome).
- Data Driven
 - Little or no prior knowledge about the distribution of the data, i.e. there are no assumptions about the data such as normality
- Non-parametric
 - Unlike parametric models which have a fixed number of parameters, in a non-parametric model, the complexity of the model grows with the number of training data
- Instance-based learning/Case-based learning/Memory-based learning/Lazy learning
 - Instead of performing explicit generalization, new instances are compared with instances stored in the memory during training, wherein the computation is delayed until classification

Core Idea of KNN

- For a new datapoint to be classified, identify the **nearby** records
- “Nearby” means the labeled datapoints with similar predictor values x_1, x_2, \dots, x_p
- Classify the datapoint as belonging to the majority class among the near by datapoints (the “neighbors”)

Distance Scores

- Euclidean Distance
- Mahalanobis Distance
- Manhattan Distance
- Chebychev Distance
- Minkowski
- Hamming

$$ED(x, y) = \sqrt{\sum_{i=1}^n |x_i - y_i|^2}$$

$$d(\vec{x}, \vec{y}) = \sqrt{(\vec{x} - \vec{y})^T S^{-1} (\vec{x} - \vec{y})}.$$

$$MD(x, y) = \sum_{i=1}^n |x_i - y_i|$$

$$CD(x, y) = \max_i |x_i - y_i|$$

$$D_{Mink}(x, y) = \sqrt[p]{\sum_{i=1}^n |x_i - y_i|^p},$$

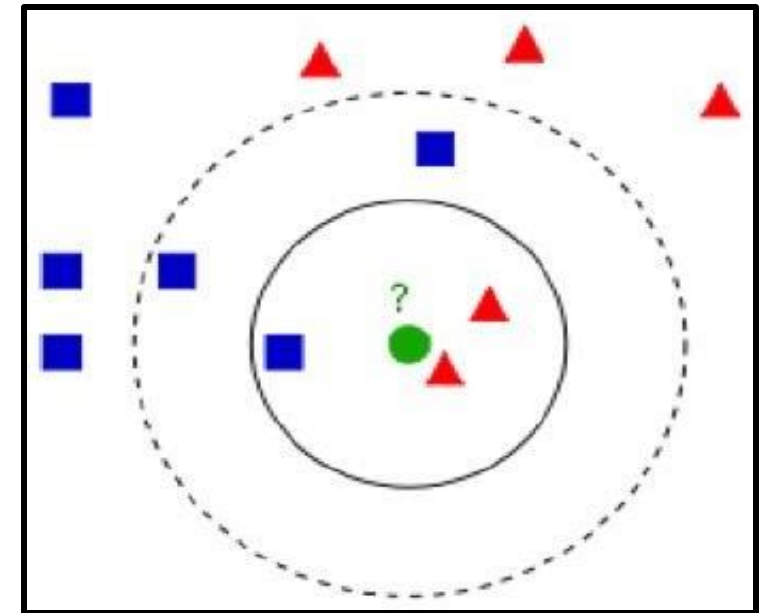
$$HamD(x, y) = \sum_{i=1}^n 1_{x_i \neq y_i}$$

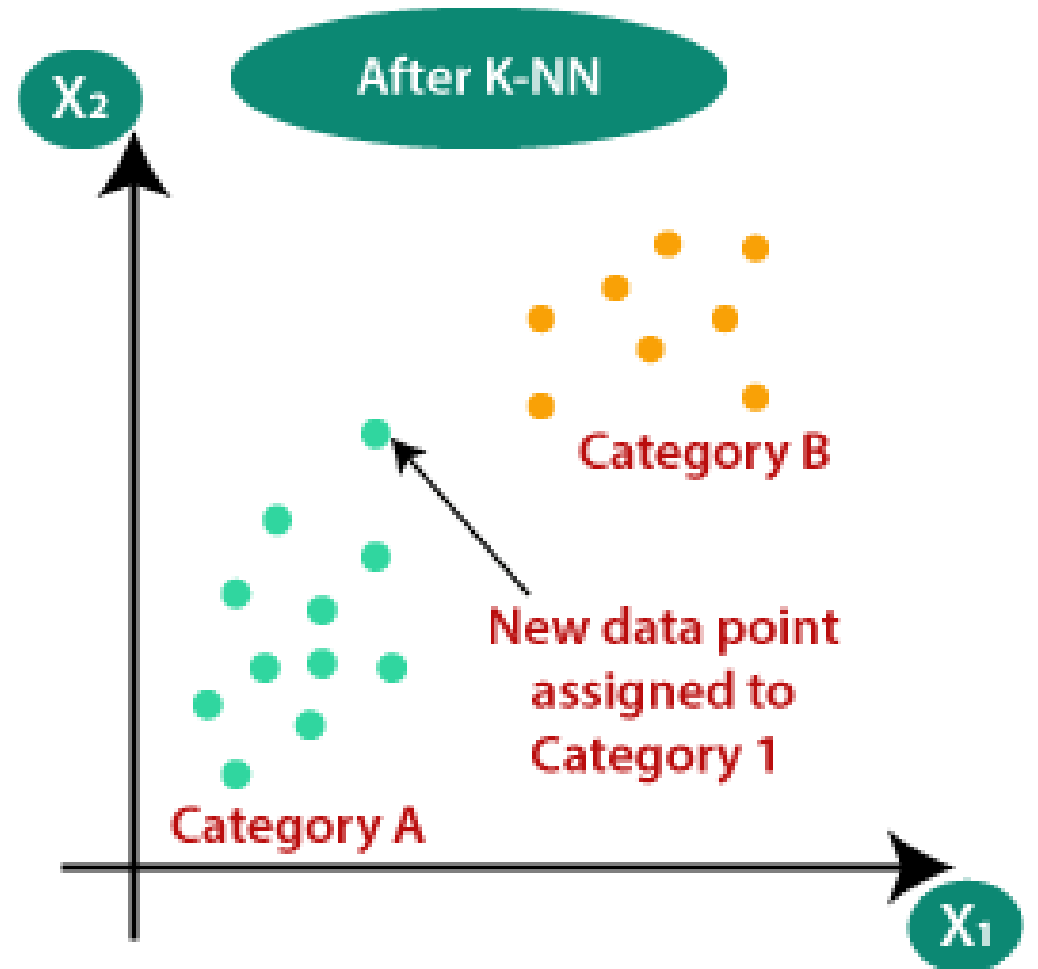
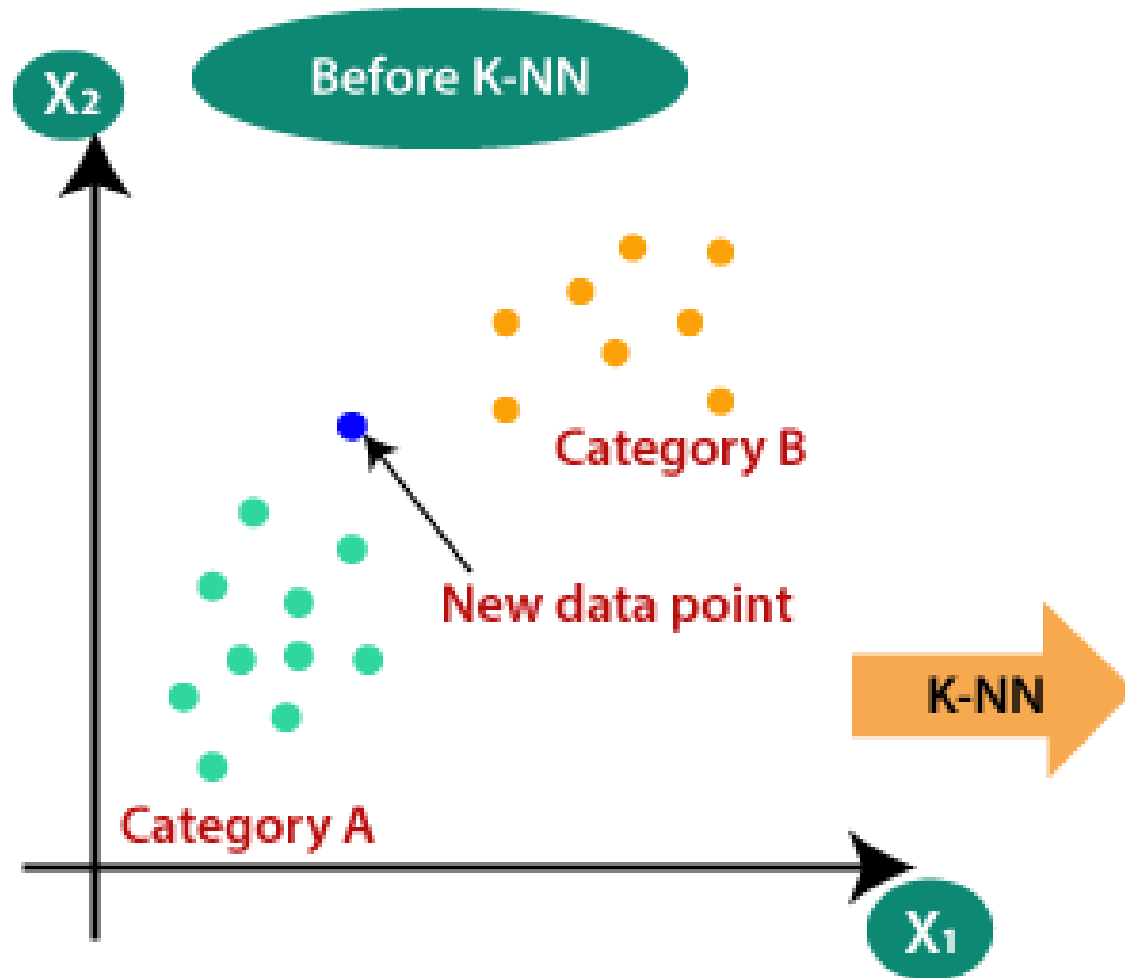
Basic KNN Algorithm - Classification

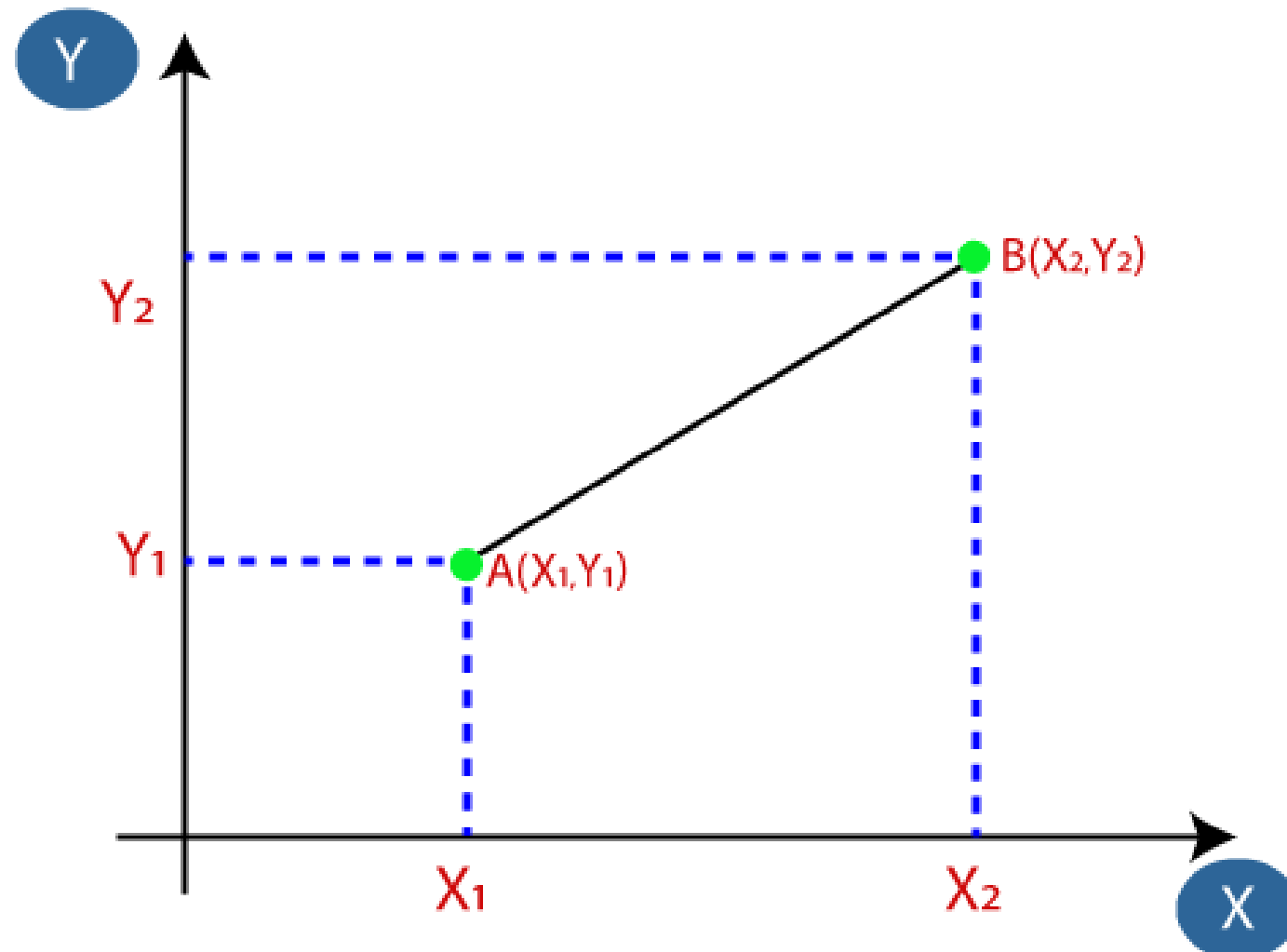
Input : Training Set D , Test datapoint d , Parameter k

Output : Class label of Test datapoint d

1. Compute the distance between test datapoint d and every datapoint in the training set D
2. Choose the k datapoints in training set D that are nearest to test sample d ; denote this set by $P (\in D)$
3. Assign the test datapoint d to the most frequent class (majority class)







Euclidean Distance between A_1 and $B_2 = \sqrt{(X_2 - X_1)^2 + (Y_2 - Y_1)^2}$

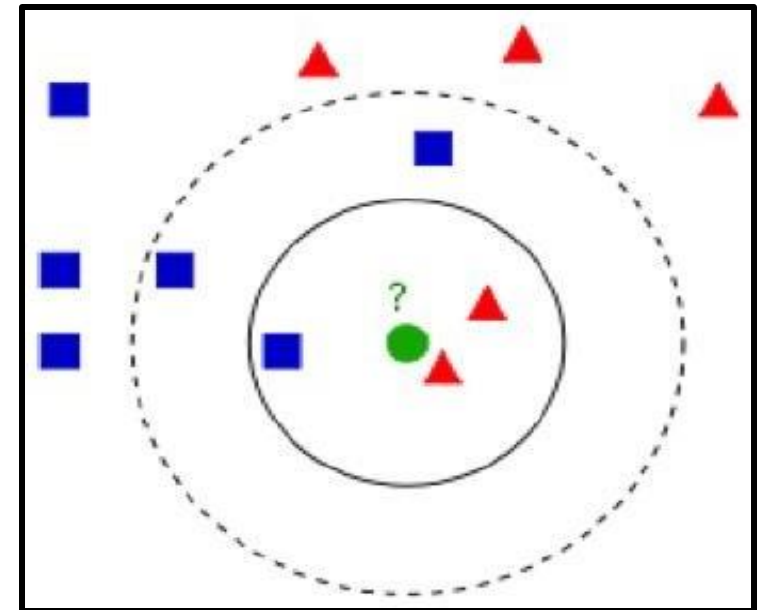


How to choose k

- k is the number of nearby neighbors to be used to classify the new record
 - $K=1$ means use the single nearest record
 - $K=3$ means use the 3 nearest records
- Low values of k (1, 3, ...) capture local structure in data (but also noise)
- High values of k provide more smoothing, less noise, but may miss local structure
 - extreme case of $k = n$ (i.e., the entire data set) is the same as the “naïve rule” (classify all records according to majority class. Here the relevance of predictors vanish and we can avoid the computation of distance scores altogether. Enough to just look at the proportion of outcome labels and choose the majority)
- Odd values of k preferred, to avoid tie caused by even values of k during majority voting

KNN for Prediction/Regression (for Numerical Outcome)

- Use average of outcome values, instead of “majority vote”
- Weighted average, weight decreasing with distance, etc



Advantages of KNN

- Simple approach and robust to noise in training data
- Learns complex structures in data easily (relatively) without having to define any statistical model
- Little or no prior knowledge about the distribution of the data, i.e. there are no assumptions about the data
- Can be used for 'imputation' of missing data

Disadvantages of KNN

- Lazy learning where computation is delayed until classification and hence need to compute distance scores each time
- Time consuming when dealing with large datasets
- Computational costs when applied to high dimensional data ('curse-of-dimensionality')

Summary

- KNN is a data driven, non-parametric approach
- Simple but powerful technique that selects k-nearest neighbors based on distance scores
- Can be used for both Classification and Prediction(Regression)
- Performance impacted by high dimensional data as well as volume of data.

If we have data from a survey which have an objective of testing with two attributes (acid durability and strength) to classify whether a special paper tissue is good or not. We have given 4 data samples.

Acid durability(x1)	Strength(x2)	Classification(y)	distance
7	7	Bad	4
7	4	Bad	5
3	4	Good	3
1	4	Good	3.6

Using this data, Classify the new data that pass the laboratory test with $x_1=3$ and $x_2=7$. Assume $k=3$.

Suppose we have height, weight and T-shirt size of some customers and we need to predict the T-shirt size of a new customer given only height and weight information we have. Data including height, weight and T-shirt size information is shown below -

Height (in cms)	Weight (in kgs)	T Shirt Size
158	58	M
160	60	M
163	61	M
163	64	L
168	62	L
170	64	L

Using this data, classify a new data point with height 161cm and weight 61kg. Assume $k=3$.

Use KNN classifier to classify new IRIS sample

Sepal Length	Sepal Width	Species
5.3	3.7	Setosa
5.1	3.8	Setosa
7.2	3.0	Virginica
5.4	3.4	Setosa
5.1	3.3	Setosa
5.4	3.9	Setosa
7.4	2.8	Virginica
6.1	2.8	Versicolor
7.3	2.9	Virginica
6.0	2.7	Versicolor
5.8	2.8	Virginica
6.3	2.3	Versicolor
5.1	2.5	Versicolor
6.3	2.5	Versicolor
5.5	2.4	Versicolor

Using this data, Classify the new sample, K=5

Sepal Length	Sepal Width	Species
5.2	3.1	?

Readings

- <https://www.analyticsvidhya.com/blog/2018/08/k-nearest-neighbor-introduction-regression-python/>
- <https://machinelearningmastery.com/k-nearest-neighbors-for-machine-learning/>
- <https://arxiv.org/pdf/1708.04321.pdf>