

DEPARTMENT OF COMPUTER SCIENCE
INSTITUTE OF MANAGEMENT AND RESEARCH, JALGAON

Name Toke Pratiksha Raju
Expt. Title Implement simple KNN using Euclidean distance
Class _____ Batch B-4 Performed on _____
Roll No. 130 Expt. No. 05 Submitted on _____
Remarks _____ Returned on _____

- * k- Nearest Neighbour (KNN) algorithm
k- nearest Neighbour is one of the simplest machine learning algorithm based on supervised learning technique. K-NN algorithm stores all the available data and classified new data point based on the similarity. This means a when new data appear then it can easily classified into a well suite category by using k-NN algorithm.

k-NN is a non-parametric algorithm which means it does not make any assumption on underlying data. It is lazy learning algorithm where all assumption is deferred until classification.

* Algorithm.

The KNN algorithm classification is performed using the following four step.

- compute the distance metric between the test data point and all the labeled data points.
- ordered the labeled data points in the increasing order of this distance metric
- select the top k labelled data points and look at the class labels.
- Find the class label that the majority of these k labeled data points have and assign it to the test data point.

Listing

* Distance calculation formula:-

1) Euclidean Distance:-

It is generally used to find the distance between two real-valued vectors.

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

2) Manhattan Distance:-

This is the simplest and one technique to calculate the distance between two points often called Taxicab distance or city Block distance.

$$\text{Manhattan Distance} = \text{sum from } i=1 \text{ to } N \text{ of } |x_{1i} - x_{2i}|$$

3) Hamming distance:-

The hamming distance is mostly used in text processing or having boolean vector. Boolean vector means the data is in the form of binary digits 0 and 1.

$$\text{Hamming distance } (x_1, x_2) = \dots$$

4) Minkowski Distance:-

Minkowski distance is the generalization form of the euclidean and Manhattan Distance.

$$\|x_1 - x_2\| = \left(\sum_{i=1}^d |x_{1i} - x_{2i}|^p \right)^{1/p}$$

DEPARTMENT OF COMPUTER SCIENCE
INSTITUTE OF MANAGEMENT AND RESEARCH, JALGAON

Name Toke pratiksha Raju

Expt. Title _____

Class FYMCA

Batch B-4

Roll No. 170

Expt. No. 05

Performed on _____

Submitted on _____

Returned on _____

Remarks _____

* Formula of information gain

① Entropy:- $-p \log_2 p - q \log_2 q$

To build a decision tree, we need to calculate two type of entropy using Frequency table.

① Entropy using the Frequency table one attribute

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

② Entropy using the Frequency table of two attribute

$$E(T, X) = \sum_{c \in X} p(c) E(c)$$

$$\text{② Gini} = 1 - \sum_{i=1}^n (p_i)^2$$

③ Information Gain = Entropy.