# K-Nearest Neighbors (KNN) Algorithm

**INTRODUCTION**

The K-Nearest Neighbors (KNN) algorithm is a fundamental supervised machine learning method used for classification and regression tasks. Developed by **Evelyn Fix and Joseph Hodges** in **1951** and later expanded by **Thomas Cover**, KNN is known for its simplicity and effectiveness in pattern recognition, data mining, and various other applications.

**What is the K-Nearest Neighbors Algorithm?**

KNN is a non-parametric algorithm that does not make any assumptions about the distribution of data. It classifies or predicts the value of a data point based on the majority class or average value of its K nearest neighbors in the training dataset. The key characteristics of KNN include:

- **Non-Parametric**: No assumptions about the data distribution.

- **Simple to Implement**: Easy to understand and apply.

- **Versatile**: Applicable to both classification and regression problems.

**Intuition Behind KNN Algorithm**

KNN operates on the principle of similarity. When a new data point needs to be classified or predicted, the algorithm examines the K nearest neighbors from the training set and uses their labels or values to make a prediction. The idea is that points close to each other are likely to share the same class or value.

**Example**

Given a set of data points with known classifications, an unclassified point is assigned a label based on the majority class of its K nearest neighbors. For instance, if a point is surrounded predominantly by points labeled 'Red', it is likely to be classified as 'Red'.

**Why Do We Need the KNN Algorithm?**

KNN is popular for several reasons:

- **Simplicity**: Easy to understand and implement.

- **No Assumptions**: Does not assume any specific data distribution.

- **Flexibility**: Handles both numerical and categorical data.

- **Adaptability**: Updates dynamically as new data is added.

**Distance Metrics Used in KNN Algorithm**

To determine the nearest neighbors, KNN uses distance metrics to measure the similarity between data points. Common distance metrics include:

**Euclidean Distance**

Euclidean distance is the straight-line distance between two points in a plane or space:

$$\text{distance}(x, X_i) = \left(\sum j=1 d(x_j − X_i,j)\right)1/2$$

where $x_j$ and $X_i,j$ are the coordinates of the points.

**Manhattan Distance**

Manhattan distance is the sum of absolute differences between coordinates:

$$\text{distance}(x,y) = \left(\sum i=1\text{-}n(|x_i − y_i|)\right)$$

This metric measures the total distance traveled along the axes.

**Minkowski Distance**

Minkowski distance generalizes both Euclidean and Manhattan distances:

$$\text{distance}(x,y) = \left(\sum i=1\text{-}n(|x_i − y_i|\char94 p)\right)\char94(1/p)$$

- When $p=2p = 2p=2$, it is equivalent to Euclidean distance.

- When $p=1p = 1p=1$, it is equivalent to Manhattan distance.

Other distance metrics, such as Hamming Distance, can also be used depending on the nature of the data.

**How to Choose the Value of K for KNN Algorithm**

Choosing the right value of K is crucial for KNN performance:

- **Higher K**: Reduces sensitivity to noise but may smooth out boundaries.

- **Lower K**: More sensitive to noise but can capture finer details.

- **Odd Values**: Prefer odd values for K to avoid ties in classification.

- **Cross-Validation**: Use cross-validation to find the optimal K value.


**Workings of KNN Algorithm**

KNN operates based on similarity:

1. **Selecting K**: Determine the number of neighbors.

2. **Calculating Distance**: Measure distances using a chosen metric.

3. **Finding Nearest Neighbors**: Identify the K nearest points.

4. **Voting/Averaging**:

    o **Classification**: Perform a majority vote among the K neighbors.

    o **Regression**: Calculate the average of the K neighbors' values.


**Advantages of the KNN Algorithm**

- **Ease of Implementation**: Simple and straightforward.

- **Adaptability**: Easily adjusts to new data.

- **Few Hyperparameters**: Only K and the distance metric need tuning.


**Disadvantages of the KNN Algorithm**

- **Scalability Issues**: Computationally expensive with large datasets.

- **Curse of Dimensionality**: Performance degrades with high-dimensional data.

- **Prone to Overfitting**: Sensitive to noisy data and irrelevant features.

**Applications of the KNN Algorithm**

- **Data Preprocessing**: Used in KNN Imputer to handle missing values.

- **Pattern Recognition**: Effective in classification tasks, such as with MNIST dataset.

- **Recommendation Engines**: Useful for grouping users and providing personalized recommendations based on similarity.