



# RAJSHAHI UNIVERSITY OF ENGINEERING AND TECHNOLOGY

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

---

## CSE 4204

---

### Lab Report

**Submitted By:**

Riyad Morshed Shoeb

Roll: 1603013

Computer Science and Engineering  
Rajshahi University of Engineering  
and Technology

**Submitted To:**

Rizoan Toufiq

ASSISTANT PROFESSOR

Computer Science and Engineering  
Rajshahi University of Engineering  
and Technology

July 24, 2022

# K Nearest Neighbor

## Introduction

K-Nearest Neighbour, abbreviated as KNN, is one of the simplest Machine Learning algorithms based on Supervised Learning technique. This algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most similar to the available categories. It stores all the available data and classifies a new data point based on the similarity. This means when new data appears, it can be easily classified into a well suite category by using the algorithm.

There are various distance matrices.

### Euclidean Distance

$$d(X, Y)_{euc} = \sqrt{\left(\sum_{i=1}^n (X_i - Y_i)^2\right)}$$

### City block Distance (Manhattan)

$$D_{eb} = \sum_n |X_i - Y_i|$$

### Square Distance $D_{sq} = MAX|X_i - Y_i|$

## The Algorithm

1. Select the number neighbors ( $K$ ).
2. Calculate the distance between  $K$  number of neighbors and the rogue pattern.
3. Assign class to the rogue pattern using a discriminant function,  $f(X)$ . For a two-class problem:

$$f(X) = \text{closest}(class_1) - \text{closest}(class_2)$$

If  $f(X)$  is positive, the pattern belongs to  $class_2$ , and  $class_1$  otherwise.

## Source Code

```
1 import numpy as np
2 from sklearn import datasets
3 import matplotlib.pyplot as plt
4
5 class KNearestNeighbor():
6     def __init__(self, X, y):
7         self.X = X
8         self.y = y
9     def predict(self, x, measure="euclidean"):
10         closest_class_1 = self.distance(self.X[self.y==0], x, measure)
11         closest_class_2 = self.distance(self.X[self.y==1], x, measure)
12         f_x = closest_class_1 - closest_class_2
```

```

13         return 1 if f_x>0 else 0
14     def distance(self, X, x, measure="euclidean"):
15         d_X_x = None
16         if measure=="city-block":
17             d_X_x = self._city_block(X, x)
18         elif measure=="square":
19             d_X_x = self._square(X, x)
20         else:
21             d_X_x = self._euclidean(X, x)
22         return d_X_x
23     def _euclidean(self, X, x):
24         d = 0
25         for sample in X:
26             d += np.sqrt((sample[0]-x[0])**2 + (sample[1]-x[1])**2)
27         d /= len(X)
28         return d
29     def _city_block(self, X, x):
30         d = 0
31         for sample in X:
32             d += np.abs((sample[0]-x[0]) + (sample[1]-x[1]))
33         d /= len(X)
34         return d
35     def _square(self, X, x):
36         d = 0
37         for sample in X:
38             d += np.max([(sample[0]-x[0]), (sample[1]-x[1])])
39         d /= len(X)
40         return d
41
42 X, y = datasets.make_blobs(n_samples=1500, n_features=2, centers=2,
43     ↪ cluster_std=1.5)
44
45 sc = plt.scatter(X[:,0], X[:,1], c=y, marker='.')
46 plt.colorbar(sc)
47 plt.title('Classified Data')
48 plt.savefig('classified-input.png')
49 plt.show()
50
51 model = KNearestNeighbor(X, y)
52 x = np.random.normal(loc=0, size=2) * 5
53 predicted = model.predict(x)
54
55 sc = plt.scatter(X[:,0], X[:,1], c=y, marker='.')
56 plt.scatter(x[0], x[1], marker='D')
57 plt.colorbar(sc)
58 plt.title(f"Rogue input: {x} is classified as class {predicted}")
59 plt.savefig('rogue-input.png')
60 plt.show()

```

## Output

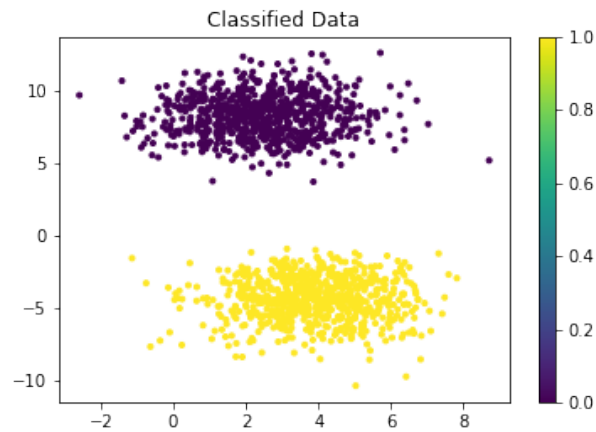


Figure 1: Classified Data

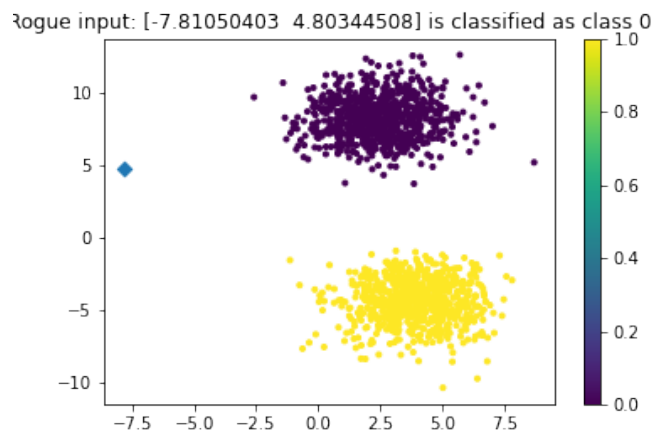


Figure 2: Classifying a rogue pattern

## Advantages

1. It is simple to implement.
2. It is robust to the noisy training data.
3. It can be more effective if the training data is large.

## Disadvantages

1. Always needs to determine the value of K which may be complex some time.
2. The computation cost is high because of calculating the distance between the data points for all the training samples.

# Single Layer Perceptron Learning Algorithm

## Introduction

The Perceptron algorithm is a two-class (binary) classification machine learning algorithm. It is a type of neural network model, perhaps the simplest type of neural network model. It consists of a single node or neuron that takes a row of data as input and predicts a class label. This is achieved by calculating the weighted sum of the inputs and a bias. The weighted sum of the input of the model is called the activation.

$$y = f_h \left( \sum_{i=1}^n w_i x_i - \theta \right)$$

If the activation is above 0, the model will output 1; otherwise, it will output 0.

$$\begin{aligned} f_h(x) &= 1 & x > 0 \\ f_h(x) &= 0 & x \leq 0 \end{aligned}$$

The Perceptron is a linear classification algorithm. This means that it learns a decision boundary that separates two classes using a line in the feature space.

## The Algorithm

1. Initialize weights,  $w$  and threshold,  $\theta$ .
2. Present input,  $x(t)$  and desired output,  $d(t)$  to the perceptron.
3. Calculate actual output

$$y(t) = f_h \left( \sum_{i=0}^n w_i(t) x_i(t) \right)$$

4. Adapt weights

$$\begin{aligned} \Delta &= d(t) - y(t) \\ w_i(t+1) &= w_i(t) + \eta \Delta x_i(t) \\ d(t) &= \begin{cases} 1, & \text{if input from class A} \\ 0 & \text{if input from class B} \end{cases} \end{aligned}$$

## Source Code

```
1  import numpy as np
2  from sklearn import datasets
3  import matplotlib.pyplot as plt
4
5  class Perceptron:
6      def __init__(self, learning_rate=0.1, n_iters=10):
7          self.lr = learning_rate
```

```

8         self.n_iters = n_iters
9         self.weights = None
10        self.bias = None
11    def learn(self, X, y):
12        n, m = X.shape
13        self.weights = np.random.uniform((m, 1))
14        self.bias = np.random.random_sample
15        for _ in range(self.n_iters):
16            f_x = np.dot(X, self.weights) + self.bias
17            y_pred = self.activation_function(f_x)
18            delta = self.lr * (y-y_pred)
19            self.weights += delta * X
20            self.bias += delta
21        return
22    def predict(self, X):
23        linear_output = np.dot(X, self.weights) + self.bias
24        y_predicted = self.activation_func(linear_output)
25        return y_predicted
26    def activation_function(self, x):
27        return np.where(x>=0, 1, 0)
28
29    X, y = datasets.make_blobs(n_samples=1500, n_features=2, centers=2,
    ↪ cluster_std=1.5)
30    plt.scatter(X[:,0], X[:,1], c=y, marker='.')
31    plt.plot(X[np.argmin(X[:,0])], X[np.argmax(X[:,0])])
32    plt.title('Before learning')
33    plt.savefig('before_learning.png')
34    plt.show()
35
36    model = Perceptron()
37    model.learn(X, y)
38
39    x0_1 = np.amin(X[:, 0])
40    x0_2 = np.amax(X[:, 0])
41    x1_1 = (-model.weights[0]*x0_1 - model.bias) / model.weights[1]
42    x1_2 = (-model.weights[0]*x0_2 - model.bias) / model.weights[1]
43    ymin = np.amin(X[:, 1])
44    ymax = np.amax(X[:, 1])
45
46    plt.scatter(X[:,0], X[:,1], marker=".", c=y)
47    plt.ylim([ymin-3, ymax+3])
48    plt.plot([x0_1, x0_2], [x1_1, x1_2])
49    plt.title('After learning')
50    plt.savefig('after_learning.png')
51    plt.show()

```

## Output

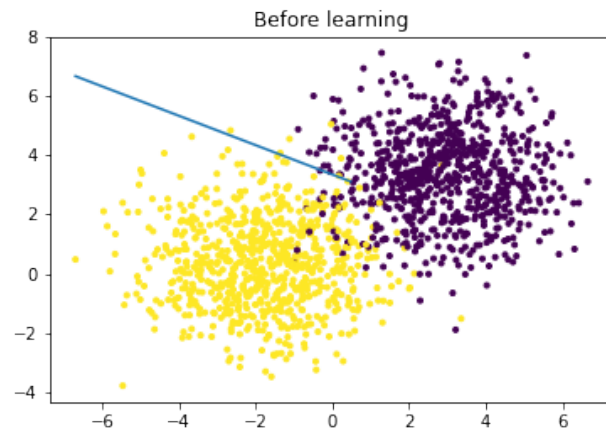


Figure 3: Before Learning Algorithm applied

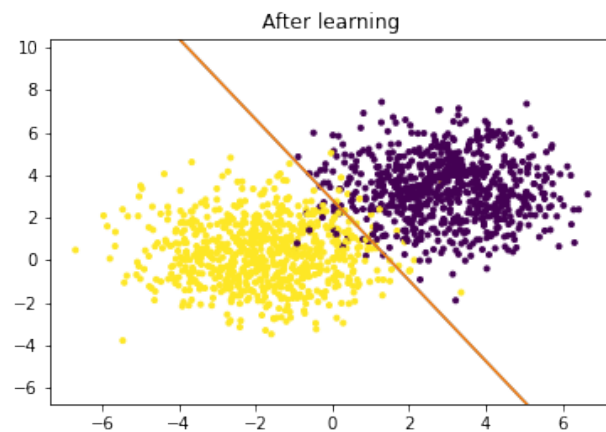


Figure 4: After Learning Algorithm applied

## Advantages

1. Single Layer Perceptron is quite easy to set up and train.
2. If the data is linearly separable, then perceptron will always reach a solution in finite time.

## Disadvantages

1. Can represent only a limited set of functions.
2. Only works for the linearly separable data.

# Naive Bayes Algorithm

## Introduction

Naive Bayes is a probabilistic machine learning algorithm based on the Bayes Theorem, used in a wide variety of classification tasks. Bayesian classification relies on the basic statistical theory of probabilities and conditional probabilities. If  $G_i, i = 1, 2, \dots, n$  be our possible list of groups, or classes, then we can define the probability of a pattern belonging to a class as  $P(G_i)$ , (where  $0 \leq P(G_i) \leq 1$ ). Given a set of measurements,  $X$ , the Bayes's rule assigns a likelihood, or probability, of it belonging to a class  $G_i$ , i.e.  $P(G_i|X)$ .  $X$  belongs to class  $i$  for

$$P(G_i|X) > P(G_j|X) \quad \text{for } i = 1, 2, \dots, n \quad i \neq j$$

i.e. we assign a pattern to the class that has the highest conditional probability of the vector  $X$  belonging to it. According to Bayes's law,

$$P(G_i|X) = \frac{P(X|G_i)P(G_i)}{\sum_j P(X|G_j)P(G_j)}$$

## Source Code

```
1  import numpy as np
2  import pandas as pd
3  import pprint
4
5  class NaiveBayes():
6      def __init__(self, X, y):
7          self.X = X
8          self.y = y
9          self.lookup_table = {}
10     def learn(self):
11         self.lookup_table['class'] = {}
12         for cls in self.y.unique():
13             self.lookup_table['class'][cls] = len(y[y==cls])/len(y)
14         for column in self.X.columns:
15             self.lookup_table[column] = {}
16             for category in self.X[column].unique():
17                 self.lookup_table[column][category] = {}
18                 for cls in self.y.unique():
19                     self.lookup_table[column][category][cls] =
20                     ↪ len(X[y==cls][X[column]==category])/len(y[y==cls])
21     def display_lookup_table(self):
22         pp.pprint(self.lookup_table)
23     def predict(self, x):
24         prediction = {}
25         p_X = 0
26         for cls in self.y.unique():
27             prediction[cls] = self.lookup_table['class'][cls]
28             for column in self.X.columns:
```



```

28         prediction[cls] *=
           ↪ self.lookup_table[column][x[column]][cls]
29     p_X += prediction[cls]
30     for cls in prediction:
31         prediction[cls] /= p_X
32     return prediction
33
34 pp = pprint.PrettyPrinter()
35 data = pd.read_csv("https://raw.githubusercontent.com/datasciencedojo/
           ↪ datasets/master/titanic.csv")
36 X, y = data[['Pclass', 'Sex', 'Embarked']], data['Survived']
37 nb = NaiveBayes(X, y)
38 nb.learn()
39 nb.display_lookup_table()
40 x = {'Embarked':'Q', 'Pclass':3, 'Sex':'female'}
41 nb.predict(x)

```

## Output

```

{'Embarked': {nan: {0: 0.0, 1: 0.0},
              'C': {0: 0.1366120218579235, 1: 0.2719298245614035},
              'Q': {0: 0.08561020036429873, 1: 0.08771929824561403},
              'S': {0: 0.7777777777777778, 1: 0.6345029239766082}},
 'Pclass': {1: {0: 0.14571948998178508, 1: 0.39766081871345027},
            2: {0: 0.1766848816029144, 1: 0.2543859649122807},
            3: {0: 0.6775956284153005, 1: 0.347953216374269}},
 'Sex': {'female': {0: 0.14754098360655737, 1: 0.6812865497076024},
         'male': {0: 0.8524590163934426, 1: 0.31871345029239767}},
 'class': {0: 0.6161616161616161, 1: 0.3838383838383838}}

# calculated probabilities for the given input
{0: 0.3978469832302337, 1: 0.6021530167697664}

```

## Advantages

1. It is easy and fast to predict class of test data set. It also perform well in multi class prediction.
2. It needs less training data compared to other algorithms.
3. It performs well in case of categorical features compared to numerical features.

## Disadvantages

1. If categorical variable has a category in test data that was not observed in training data, then model assigns a 0 probability.