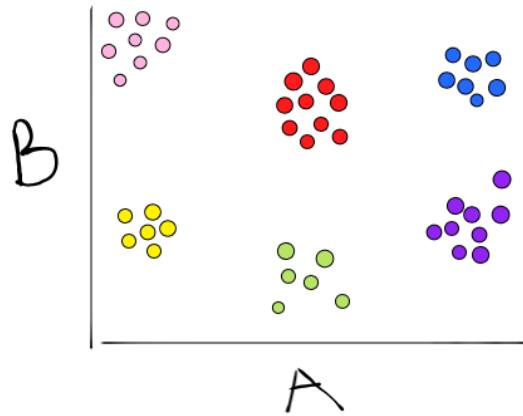# K-Nearest Neighbors Algorithm

- Nayan Man Singh Pradhan

The K-Nearest Neighbors (KNN) Algorithm is a popular **supervised** machine learning algorithm that can be used for both **classification** and **regression** problems (although it is used more for classification problems).
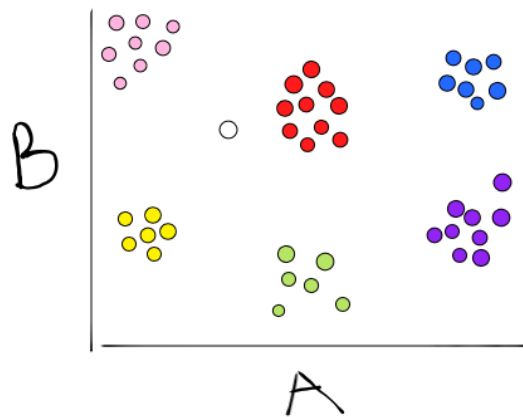
**This blog aims to:**

- Introduce the KNN algorithm

- Demonstrate how the KNN algorithm works through visual pictures

- Answer what the best value for K is

- Discuss the advantages and disadvantages of the KNN algorithm

- Show an example of the algorithm working
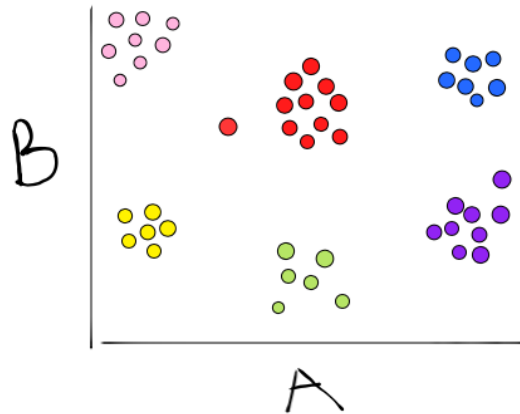
**How the KNN algorithm works:**

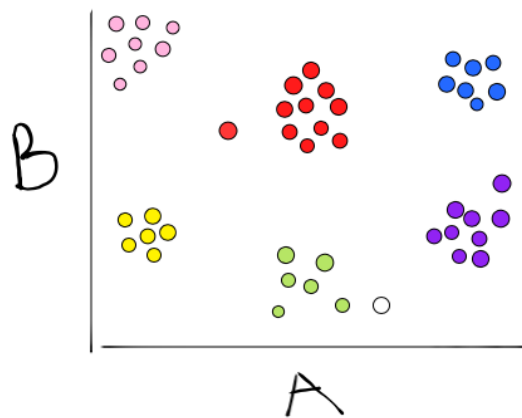- Cluster the given dataset into the known categories.



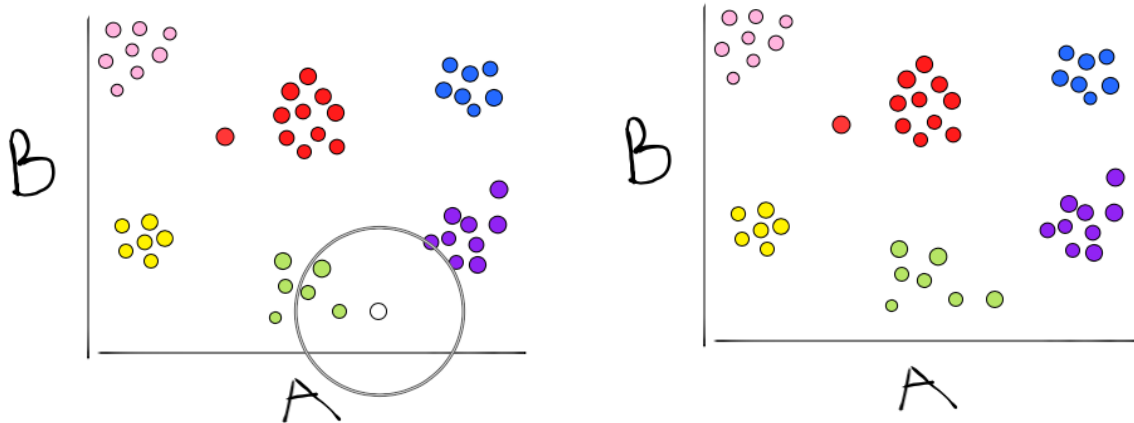- Add the unknown cell that we want to categorize.

- Classify the new data point by looking at the nearest neighbors. The K in KNN signifies how many neighbors the algorithm considers while classifying the new data point.



- The K plays a significant role in classification. Let us assume a new data point in the figure.

- If K = 3, the data point is classified as green. This is because, the 3 closest given data points are classified as green.



- But, if K = 14, the data point is classified as purple. This is because, there are 6 green nearest data points and 8 purple nearest data points. The KNN algorithm classifies the point to the data set with the largest number of nearest neighbors. Hence, one must be very careful while choosing K.

**How to pick the optimal K value:**

- There is no formula the gives the correct/optimal value for K.

- The value of K differs according the the number of data sets.

- Having K as an odd number is usually recommended as the extra data point acts as
  a tie breaker (if the data sets are evenly spread out between two classifications).

- Personally, I like to choose K as the integer equal to the square root of the length of
  the data set. If K is an even integer, I add 1 to K. This is illustrated by the code
  snippet below:

```
47  K = math.sqrt(len(data)) # adjust value of K accordingly
48  K = int(K)
49
50  if (K % 2 == 0):
51      K = K+1
52
53  print("K=", K)
```

**Advantages of the KNN algorithm:**

- Understanding and implementing the algorithm is simple.

- The algorithm can be used for classification and regression problems, hence, it is
  versatile.

- There is no need for building models and adjusting various parameters.

**Disadvantages of the KNN algorithm:**

- The algorithm gets slow as the number of independent variables increases.

**Example of the KNN algorithm:**

The algorithm below is predicting whether an employer gets accepted or rejected while

applying for a job.

```
1  # USING K NEAREST NEIGHBOURS (KNN)
2
3  import os
4  import numpy as np
5  import pandas as pd
6  from pandas import read_csv
7  import sklearn
8  from sklearn import linear_model
9  from sklearn.utils import shuffle
10 import matplotlib.pyplot as pyplot
11 import pickle
12 from matplotlib import style
13 from sklearn import preprocessing
14 from sklearn.neighbors import KNeighborsClassifier
15 import pickle
16
17 data = pd.read_csv("Placement.csv") # import data
18 data = data[["gender", "ssc_p", "ssc_b", "hsc_p", "hsc_b", "hsc_s", "degree_p", "degree_t", "workex", "etest_p"
19 data = data[["gender", "ssc_p", "ssc_b", "hsc_p", "hsc_b", "hsc_s", "degree_p", "degree_t", "workex", "etest_p"
20
21 # changing the strings into integers for prediction
22 le = preprocessing.LabelEncoder()
23 data.gender = le.fit_transform(list(data["gender"]))
24 data.ssc_p = le.fit_transform(list(data["ssc_p"]))
25 data.ssc_b = le.fit_transform(list(data["ssc_b"]))
26 data.hsc_p = le.fit_transform(list(data["hsc_p"]))
27 data.hsc_b = le.fit_transform(list(data["hsc_b"]))
28 data.hsc_s = le.fit_transform(list(data["hsc_s"]))
29 data.degree_p = le.fit_transform(list(data["degree_p"]))
30 data.degree_t = le.fit_transform(list(data["degree_t"]))
31 data.workex = le.fit_transform(list(data["workex"]))
32 data.etest_p = le.fit_transform(list(data["etest_p"]))
33 data.specialisation = le.fit_transform(list(data["specialisation"]))
34 data.mba_p = le.fit_transform(list(data["mba_p"]))
35 data.status = le.fit_transform(list(data["status"]))
36 # data.salary = le.fit_transform(list(data["salary"]))
37
38 # print(data.head())
39 predict = "status" # what to predict
40
41 X = np.array(data.drop([predict], 1))
42 y = np.array(data[predict])
43
44 x_train, x_test, y_train, y_test = sklearn.model_selection.train_test_split(X, y, test_size = 0.2)
45
46 # K = 5 # adjust value of K accordingly
47 K = math.sqrt(len(data)) # usually, best K is given by square root of n
48 model = KNeighborsClassifier(n_neighbors = int(K))
49
50 model.fit(x_train, y_train)
51 acc = model.score(x_test, y_test)
52 print("Accuracy using KNN:", acc)
```

Accuracy using KNN: 0.7906976744186046

**\*The full code can be found in my GitHub:**

https://github.com/nayan-pradhan/MachineLearning_Summer2020/blob/master/

jobAcceptance/predictionUsingKNN.ipynb