# 1   Nearest Neighbor

## 1.1   Describe Implementation

The general idea is to compute the L2 Norm to the nearest neighbor, K point, and return the most recurring label in those neighbors.

### 1.1.1   KNN_TEST

This function genrates a label based on the points given to it, computes the distance, to each of those points (training points) and puts them in a list. This list is then sorted based on the distance that was calculated before. The subset of the K nearest neighbor or point will have their labels summed and then we will output the sign of the summation as the label, so if our predicted label is less than 0 we return -1 and if it is greater than 0 we return 1. In case the predicted value is 0 we will choose the result based on the ordering of the training data and from the set of neighbors in the sorted list.

```
1    def KNN_test(X_train, Y_train, X_test, Y_test,K):
2      NumberOfCorrect = 0
3      for i, j in enumerate(X_test):
4        radius = []
5        for n, m in enumerate(X_train):
6          radius.append((((math.pow((m[0] - j[0]), 2) + math.pow((m[1] - j[1]), 2)), Y_train[n
     ][0]))
7      List_sorted = sorted(radius, key=lambda member: member[0], reverse=False)[:K]
8      predictedvalue = sum([member[1] for member in List_sorted])
9      if (not predictedvalue):
10       predictedvalue = List_sorted[0][1]
11     elif (predictedvalue < 0):
12       predictedvalue = -1
13     elif (predictedvalue > 0):
14       predictedvalue = 1
15     if predictedvalue == Y_test[i][0]:
16       NumberOfCorrect = NumberOfCorrect+1
17     return NumberOfCorrect / len(X_test)
18
```

### 1.1.2   Choose_K

This function is simply checking all possible accuracies on the given test data and returns the index with the greatest accuracy which will be our K.

```
1      OldAccuracy = 0
2      NewAccuracy = 0
3      Predicted_K = 1
4      K = len(X_train)
5      for i in range(1, K):
6        OldAccuracy = KNN_test(X_train,Y_train,X_val,Y_val, i)
7        if(OldAccuracy>NewAccuracy):
8          NewAccuracy = OldAccuracy
9          Predicted_K = i
10     return Predicted_K
11
```

## 1.2   Algorithm Classification

Using the testing data given to us from the project description we will get the following results

The accuracy on K=1, K=2, K=3 are 0.3, 0.4, 0.6 respectively. You can see more details like which sample was predicted incorrectly on the picture below

1

The best K for our Algorithm would be 9 with the accuracy of 0.7, for more details like which specific sample was predicted wrong refer to the picture below.



2

---

[1]The picture has a high resolution so you should be able to zoom in and see the detail, but the picure has also been included in the zip file so please refer to that if there is any issues

[2]The picture has a high resolution so you should be able to zoom in and see the detail, but the picure has also been included in the zip file so please refer to that if there is any issues

# 2   KMeans Clustering

## 2.1   Describe Implementation

### 2.1.1   K_Means

Our implementation for KMeans starts with selecting a random number generated using the random library. This number is used as an index to grab the points from the data; these points are selected as our center.

The datapoints are passed into a function along with the centers and enter a while loop with the condition being it runs until error equals 0. Error is defined as the centers from one run not being the same centers generated from the next run.

A list of lists is generated depending on the number of centers passed in, e.g. if three centers were generated, K=3, so three clusters are necessary. Each point in the data is tested against the center points, calculating the Euclidean distance and putting the point in the corresponding cluster that had the least amount of distance to the center.

The new centers are then calculated using the mean of each individual cluster and tested against the old centers to see if the loop should be exited. If they are not the same these steps are repeated.

```
1  randCenterIndexes = random.sample(range(len(X)), K)
2  randCenters = []
3  for k in randCenterIndexes:
4  randCenters.append(X[k])
5
6  output = FindClusters(X, randCenters)
7
8  return output
9
```

### 2.1.2   K_Means_better

This function takes X, our sample data-set and the K value and runs the K-Mean algorithm 1000 times (can be changed by the user) we will then count the number of occurrences of each cluster center and return the ones with the most repeated or occurred cluster centers, this value will be stored in a numpy array of cluster centers and then returned.
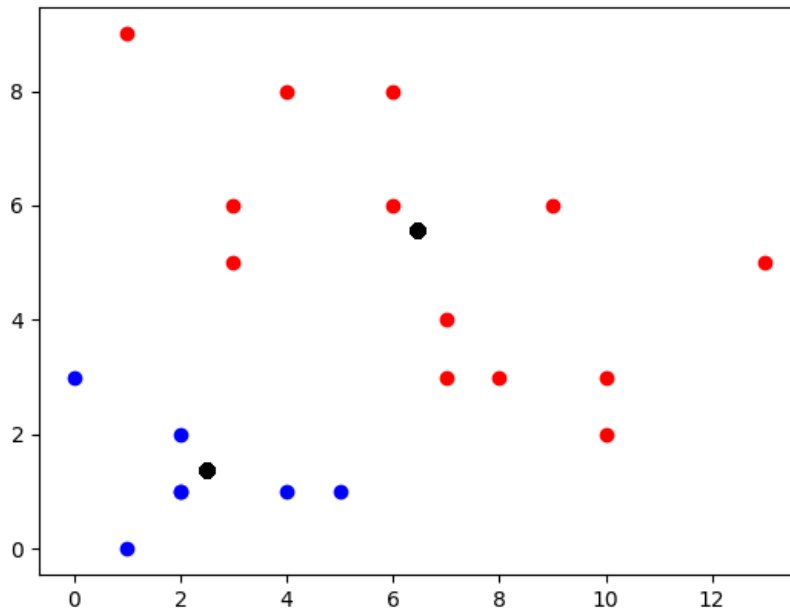
```
1    def K_Means_better(X, K):
2    run = 1000
3    bestCenterList = []
4    occurences = []
5    j = 0
6    for i in range(run):
7      try:
8        result = K_Means(X, K)
9        bestCenterList.append(result)
10       j += 1
11   except:
12     pass
13   for j in range(len(bestCenterList)):
14     count = bestCenterList.count(bestCenterList[j])
15     occurences.append(count)
16   maxNum = 0
17   for k in occurences:
18     if (k > maxNum):
19       maxNum = k
20   return bestCenterList[maxNum]
21
```
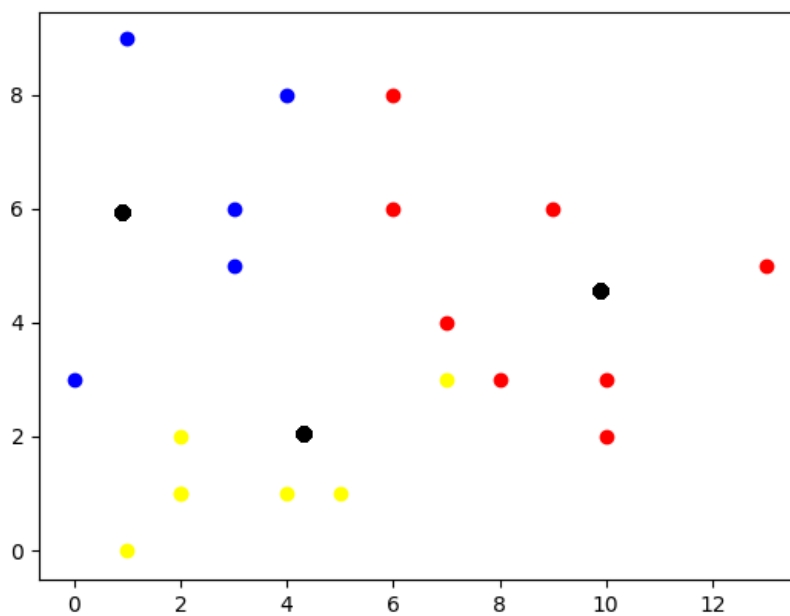
## 2.2    Algorithm Classification

### 2.2.1    KMeans on K = 2, K = 3

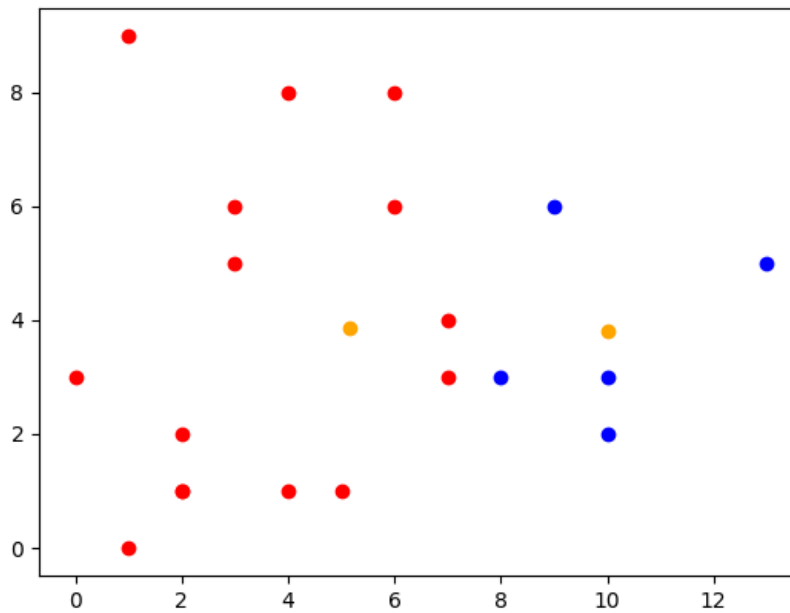K = 2 The cluster centers have been labeled as a black dot.



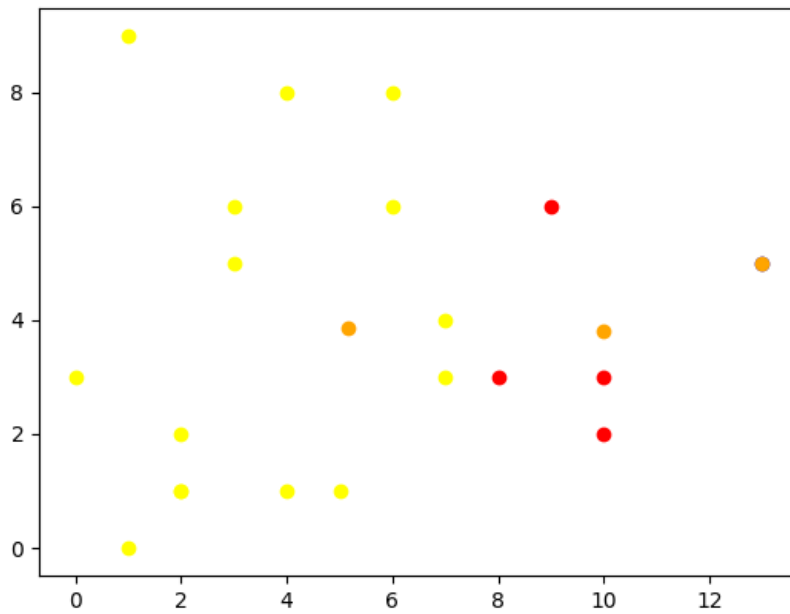K = 3 The cluster centers have been labeled as a black dot.

### 2.2.2   KMeans Better on K = 2, K = 3

K = 2 The centers are the two orange dots



K = 3 the center is the point at 13, 4.5 ( it is a point and a center)

# 3   Perceptron

## 3.1   Describe Implementation

### 3.1.1   perceptron_train

This function will be computing a weight (vector) and a bias using the data-set X and the labels for that set Y. We start with the initial weights and bias of 0 and we will kep track of epoch by a variable called epoch, we would also have to consider the fact that some data might not be linearly separable hence we will have another variable called maxEpoch which we would use to break our if the algorithm is not successful after maxEpoch tries. We will be using a while loop that checks if we are done yet, so we've gone an entire epoch without updating, or if we have reached the maxEpoch in which case the program gets out of the while loop and returns the weights and bias of the last iteration. Inside the while loop we are computing the activation and if the activation returns a number less than or equal to 0 we would ned to update hence we update our weight vector and the bias, they're both in the same list, the first element is weight vector, 2nd is bias(W_B_List) if data is linearly separable we would update our Done flag to true and get out of the loop and finally return the weight vector and bias (W_B_List)

```
1   Done = False
2   epoch = 0
3   maxEpoch = 199
4   activation = 0
5   W_B_List = [[],[]]
6   W_B_List[0] =[0]*len(X[0])
7   W_B_List[1] = [0]
8   while ( not Done or epoch < maxEpoch):
9     epoch+=1
10    for i in range (len(X)):
11      activation = 0
12      for m in range(len(X[i])):
13        activation += X[i][m] * W_B_List[0][m]
14      activation+= W_B_List[1][0]
15      if  (( activation * Y[i][0] ) <= 0):
16        for j in range(len(X[0])):
17          W_B_List[0][j] = W_B_List[0][j] + Y[i][0] * X[i][j]
18        W_B_List[1][0] = W_B_List[1][0] + Y[i][0]
19      else:
20        Done = True
21  return W_B_List
22
```

### 3.1.2   perceptron_test

This function computes the accuracy of the weights and bias given the data-set X and their label Y. This function computes the activation for each of the samples in the data-set X and compares it with their label. If the activation matches the label we would increment our NumberOfCorrect variable at the end we will divide this number by the length of our data-set X to find the accuracy and we would return that number.
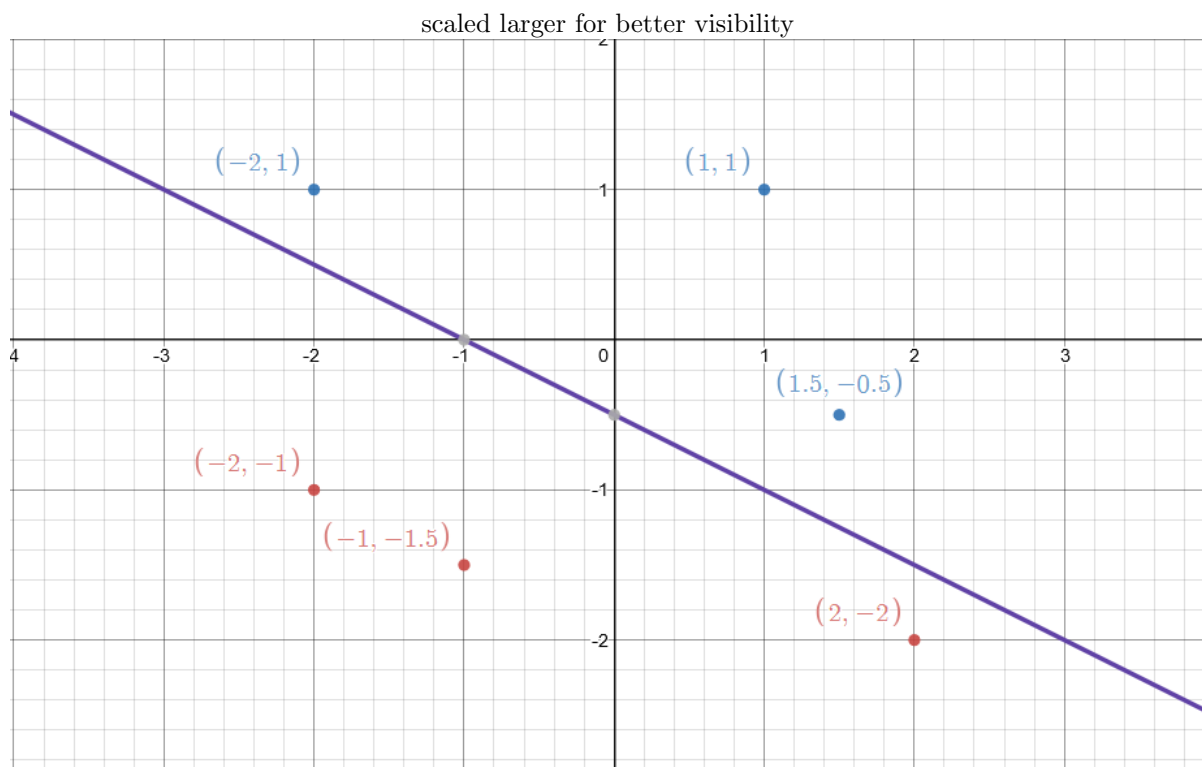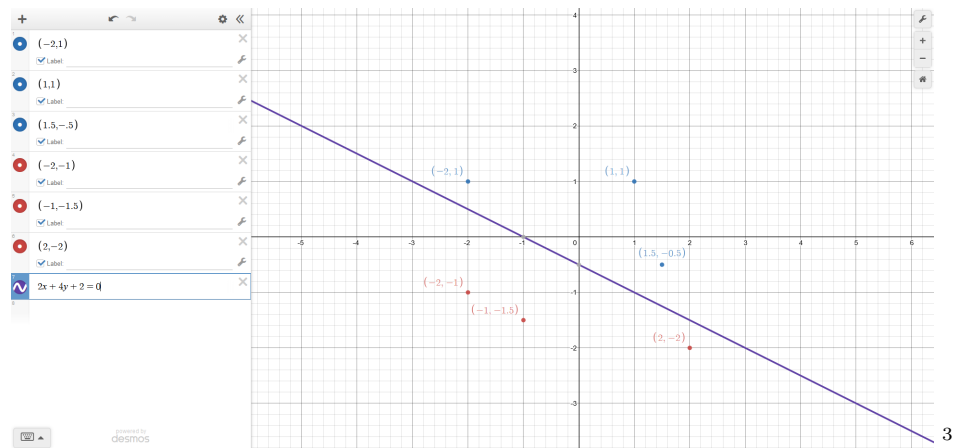
```
1   NumberOfCorrect = 0
2   activation = 0
3   for i, j in enumerate(X_test):
4     activation = np.dot(j,w) + b
5     if activation * Y_test[i] > 0:
6       NumberOfCorrect+=1
7   return NumberOfCorrect/len(X_test)
8
```

## 3.2   Decision Boundary

I have used the Desmos website to plot the points and the decision boundary for this part. Blue points are the ones with label 1 and red points have label of -1. The algorithms returns 2,4 for our weights and 2 for

our bias hence we will be using the function $2X_1 + 4X_2 + 2 = 0$ we could simplify this function but since we will be drawing it to show the decision boundary there is no need.



scaled larger for better visibility



---

[3]The picture has a high resolution so you should be able to zoom in and see the detail, but the picure has also been included in the zip file so please refer to that if there is any issues