# HOMEWORK 1

## DATA EXPLORATION

Solution to question (a.i), (a.ii)
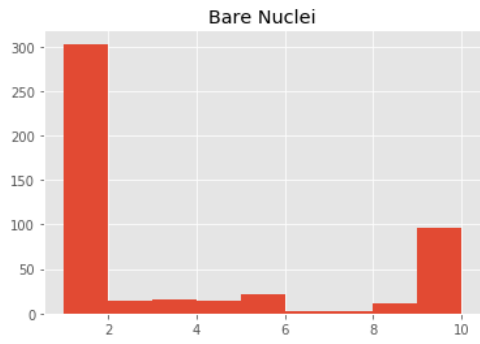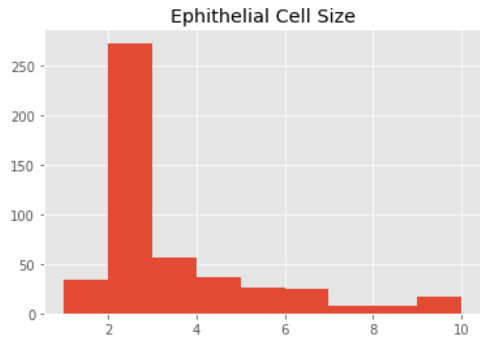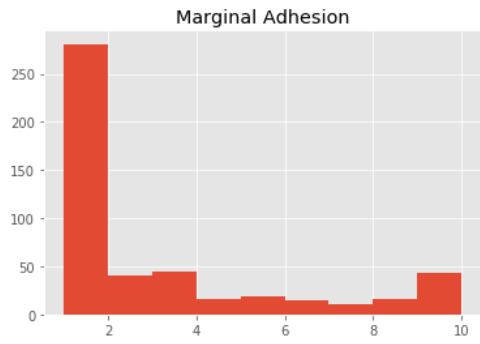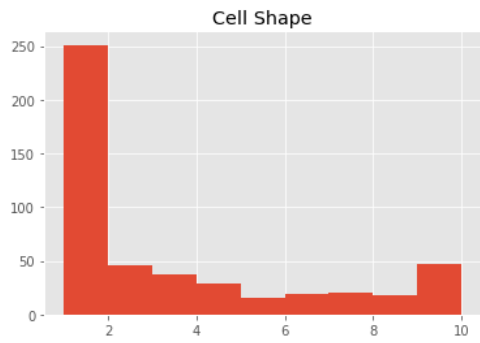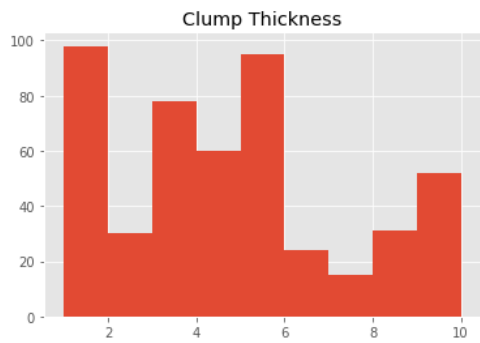
In this section, I have just parsed the dataset into lists to plot the histograms as required by question 2.(a) and 2.(b). It is visible that the data is not equally distributed. There are more data samples in the Benign case than in the Malignant case. Most features are more or less following a noraml distribution.
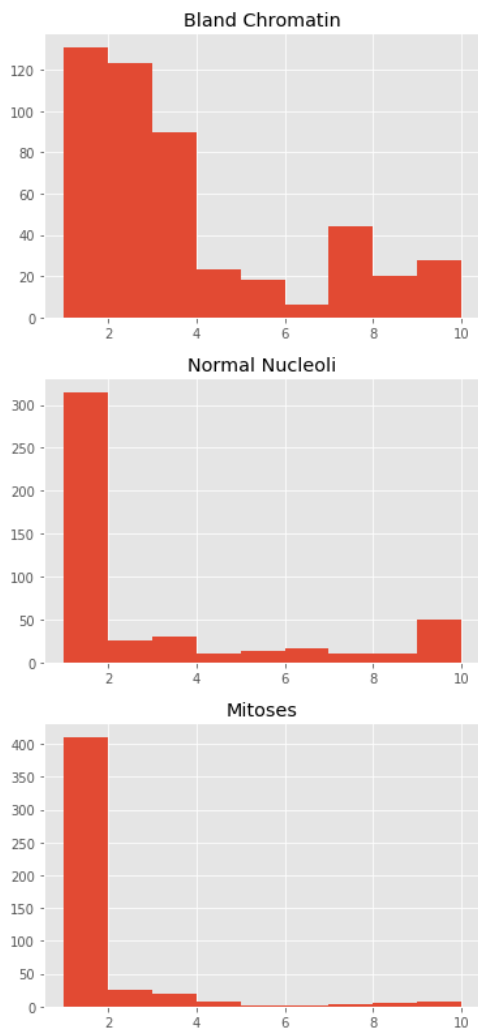
```python
1  import csv
2  import matplotlib.pyplot as plt
3  import numpy as np
4
5  malign = 0
6  benign = 0
7  Clump_Thick =  []
8  Cell_Size = []
9  Cell_Shape =  []
10 Marg_Adh =  []
11 E_Cell_Size = []
12 B_Nuclei =  []
13 Bland_Chromatin =  []
14 Norm_Nucleoli =  []
15 Mitoses =  []
16 type_cancer =  []
17
18
19 with open('hw1_question1_train.csv') as csvfile:
20   readCSV = csv.reader(csvfile)
21   for row in readCSV:
22     '''read each value from the row and place it in an appropriate list.
23     these lists contain each individual column of our data.
24     '''
25     Clump_Thick.append(int(row[0]))
26     Cell_Size.append(int(row [1]))
27     Cell_Shape.append(int(row [2]))
28     Marg_Adh.append(int(row [3]))
29     E_Cell_Size.append(int(row [4]))
30     B_Nuclei.append(int(row [5]))
31     Bland_Chromatin.append(int(row [6]))
32     Norm_Nucleoli.append(int(row [7]))
33     Mitoses.append(int(row [8]))
34     type_cancer.append(row [9])
35     '''
36     maintain two variables that hold the number of malignant cases and
37     the number of benign cases
38     '''
39     if (row [9]== '2'):
40        benign = benign + 1
41     else:
42        malign = malign + 1
43
44 #PLOTTING THE FEATURES
45 print ("Number of Benign cases are:", benign, "\nNumber of malignant cases are:",malign,"\nNo, the two classes are not equally distributed.\n")
46 plt.style.use('ggplot')
47 '''
48 Plotting the histogram of the various features.
49 '''
50 plt.title('Clump Thickness')
51 plt.hist(np.sort(Clump_Thick),bins=[1,2,3,4,5,6,7,8,9,10])
52 plt.show()
53 plt.title('Cell Size')
54 plt.hist(np.sort(Cell_Size),bins=[1,2,3,4,5,6,7,8,9,10])
55 plt.show()
56 plt.title('Cell Shape')
57 plt.hist(np.sort(Cell_Shape),bins=[1,2,3,4,5,6,7,8,9,10])
58 plt.show()
59 plt.title('Marginal Adhesion')
60 plt.hist(np.sort(Marg_Adh),bins=[1,2,3,4,5,6,7,8,9,10])
61 plt.show()
62 plt.title('Ephithelial Cell Size')
63 plt.hist(np.sort(E_Cell_Size),bins=[1,2,3,4,5,6,7,8,9,10])
64 plt.show()
65 plt.title('Bare Nuclei')
66 plt.hist(np.sort(B_Nuclei),bins=[1,2,3,4,5,6,7,8,9,10])
67 plt.show()
68 plt.title('Bland Chromatin')
69 plt.hist(np.sort(Bland_Chromatin),bins=[1,2,3,4,5,6,7,8,9,10])
70 plt.show()
71 plt.title('Normal Nucleoli')
```

```
71 plt.title('Normal Nucleoli')
72 plt.hist(np.sort(Norm_Nucleoli),bins=[1,2,3,4,5,6,7,8,9,10])
73 plt.show()
74 plt.title('Mitoses')
75 plt.hist(np.sort(Mitoses),bins=[1,2,3,4,5,6,7,8,9,10])
76 plt.show()
```

### Clump Thickness



### Cell Size



### Cell Shape



### Marginal Adhesion



### Ephithelial Cell Size



### Bare Nuclei

## IMPLEMENTATION OF k-NN

In my implementation of k-NN, I first created a class called 'Datapoint'. An object of this class holds all the feautres of our dataset. Also the class has a couple of helper functions to make things easier.

Solution to 2(a.iii)

> As part of further data exploration, I next randomly selected 5 pairs of datapoints. Then plotted these points. This helps to judge the possible correlation between features. From my observation I could not see too much correlation between any of the features that I had chosen. Also the scatter plot helped to show that the data was not linearly seperable and hence something simple like a perceptron would be unable to classify this data.

The dataset is parsed and the corresponding datapoints are then created. To compute the nearest neighbors I have used 3 different distance metrics; euclidean distance, hamming distance and cosine similarity. Having used the training and test data I finally compare the performance of the 3 distance metrics and for different values of 'k'.

Solution to 2(b.i)

> To perform nearest neighbor classification the following steps were followed.
>
> 1. Find the neighbors of the datapoint being classified
> 2. Of the neighbors find the neighbor neighbor closest to the datapoint in question. Here is where the value of 'k' comes into play. The number 'k' tells us how many of the nearest neighbors to a given point we consider. The nearness of a point in the space and the datapoint can be measured by different metrics.
> 3. Each of the 'k' nearest neighbors get to vote to steer the results toward the class they belong to. Once all the neighbors vote, the datapoint is said to belong to the class which got the most votes.

> In this question the euclidean distance metric was used. The code for this can be found from line 85 to 153 below.

Solution to 2(b.ii),2(b.iii) and 2(b.iv)

Now, that it is clear how a k-NN works, the next step was to train the k-NN to classify a given point. To ensure this is done accurately we use a development set to train a model while varying the values of k. Also, we ran 3 different models for 3 different distance metrics i.e. Euclidean Distance, Cosine Similarity and the Hamming distance. Using the Acc and B.Acc performance metrics we compared the trained models and chose the best model of this to use on our test data.

```
 1 import numpy as np
 2 import matplotlib.pyplot as plt
 3 from sklearn.metrics.pairwise import cosine_similarity
 4 from scipy.spatial import distance
 5 from scipy import spatial
 6 import csv
 7 import math as m
 8 import operator
 9
10 #Creating a class for representing the Datapoints.
11
12 class Datapoint(object):
13
14   def __init__ (self, feats):
15       self.feature_1 = feats['Clump_Thick']
16       self.feature_2 = feats['U_CSize']
17       self.feature_3 = feats['CShape']
18       self.feature_4 = feats['M_Adh']
19       self.feature_5 = feats['E_CSize']
20       self.feature_6 = feats['B_Nuclei']
21       self.feature_7 = feats['Bl_Chroma']
22       self.feature_8 = feats['Norm_Nucleoli']
23       self.feature_9 = feats['Mitoses']
24       self.type_of_tumor = feats['type_of_tumor']
25   #Returns the features as a Numpy Array.
26   def feature_vector (self):
27       return np.array([self.feature_1, self.feature_2, self.feature_3,\
28                        self.feature_4, self.feature_5, self.feature_6,\
29                        self.feature_7, self.feature_8, self.feature_9, self.type_of_tumor])
30
31   def __str__(self):
32     return "\nClump Thickness:{}, \nCell Size:{}, \nCell Shape:{}, \nMarginal Adhesion:{},\
33       \nEpithelial Cell Size:{}, \nBare Nuclei:{}, \nBland Chromatin:{}, \
34       \nNormal Nucleoli:{}, \nMitoses:{}, \nType of Tumor:{}".format(self.feature_1, self.feature_2, self.feature_3, self.feature_4, self.feature_5,\
35                                                 self.feature_6, self.feature_7, self.feature_8,  self.feature_9, self.type_of_tumor)
36 #Function that creates the datapoints and writes the corresponding feature value to the datapoints.
37 def parse_dataset(filename):
38     with open(filename) as csvfile:
39       dataset = []
40       readCSV = csv.reader(csvfile)
41       for row in readCSV:
42         a=Datapoint({'Clump_Thick':int(row[0]), 'U_CSize':int(row[1]), 'CShape':int(row[2]), 'M_Adh':int(row[3]), \
43                          'E_CSize':int(row[4]), 'B_Nuclei':int(row[5]), 'Bl_Chroma':int(row[6]), 'Norm_Nucleoli':int(row[7]),
44                          'Mitoses':int(row[8]),'type_of_tumor':int(row[9])})
45
46         dataset.append(a.feature_vector())
47
48
49     return dataset
50
51 #PARSE the training , development and testing  dataset.
52 dataset_train = parse_dataset('hw1_question1_train.csv')
53 dataset_test  = parse_dataset('hw1_question2_test.csv')
54 dataset_dev = parse_dataset('hw1_question2_dev.csv')
55 #Printing some statistics about the data.
56 print("Total Number of Data Points in Training set: {0}".format(len(dataset_train)))
57 print("Total Number of Data Points in Dev set: {0}".format(len(dataset_dev)))
58 print("Total Number of Data Points in Testing set: {0}".format(len(dataset_test)))
59
60 #Function to plot.
61 def plot_data(dataset,x,y):
62
63   benign_feat1 = [data[x] for data in dataset if data[9] == 2]
64   benign_feat2 = [data[y] for data in dataset if data[9] == 2]
65   malign_feat1 = [data[x] for data in dataset if data[9] == 4]
66   malign_feat2 = [data[y] for data in dataset if data[9] == 4]
67   plt.scatter(benign_feat1,benign_feat2, c='b',marker ='x', label = 'Benign')
68   plt.scatter(malign_feat1,malign_feat2, c='r',marker ='o', label = 'Malignant')
69   plt.legend()
70   plt.show()
71 #Plotting 5 random sets of features. It is not really random, but just different combinations
72 #features
73 plot_data(dataset_train, 6, 4)
74 plot_data(dataset_train, 0, 2)
75 plot_data(dataset_train, 1, 3)
76 plot_data(dataset_train, 3, 5)
77 plot_data(dataset_train, 4, 2)
78
79 print("The Datapoints are not linearly sperable.")
80
81 #--------------------------------------------------------------CLASSIFICATION--------------------------------------------------------------
82 '''
83 As a first step in classification, I define a function that can compute the eucledian distance between two datapoints
84 '''
85 #DISTANCE METRICS------------------------------------------------
86 def EuclidDistance (instance1, instance2, length):
87   distance = 0
```

```
 88    for x in range(length):
 89      distance += pow((instance1[x]-instance2[x]),2)
 90    return m.sqrt(distance)
 91
 92 def hamming_distance(instance1, instance2, length):
 93    return(distance.hamming(instance1, instance2))
 94
 95 def cosine_similarity(instance1, instance2, length):
 96    return(spatial.distance.cosine(instance1, instance2))
 97 #----------------------------------------------------------------
 98 # Function to compute the neighbors to the test point.
 99 def getNeighbors(dataset, testcase, k):
100    distances = []
101    length = len(testcase) - 1
102    for x in range((len(dataset))):
103      dist = EuclidDistance(testcase,dataset[x],length)
104      distances.append((dataset[x],dist))
105    distances.sort(key=operator.itemgetter(1))
106    neighbors = []
107    for x in range (k):
108        neighbors.append(distances[x][0])
109    return neighbors
110
111 def getNeighbors_hamming(dataset, testcase, k):
112    distances = []
113    length = len(testcase) - 1
114    for x in range((len(dataset))):
115      dist = hamming_distance(testcase,dataset[x],length)
116      distances.append((dataset[x],dist))
117    distances.sort(key=operator.itemgetter(1))
118    neighbors = []
119    for x in range (k):
120        neighbors.append(distances[x][0])
121    return neighbors
122
123 def getNeighbors_cosine(dataset, testcase, k):
124    distances = []
125    length = len(testcase) - 1
126    for x in range((len(dataset))):
127      dist = cosine_similarity(testcase,dataset[x],length)
128      distances.append((dataset[x],dist))
129    distances.sort(key=operator.itemgetter(1))
130    neighbors = []
131    for x in range (k):
132        neighbors.append(distances[x][0])
133    return neighbors
134 #----------------------- End of Distance Metrics-----------------------------
135 #Function that calculates the final class based on a vote.
136 def getResponse(neighbors):
137     votes = []
138     four = 0
139     two = 0
140     for x in range (len(neighbors)):
141       response = neighbors [x][-1]
142       votes.append(neighbors[x][-1])
143     for x in votes:
144       if x == 2:
145         two += 1
146       else:
147         four += 1
148     if (two>four):
149       result = 2
150     else:
151       result = 4
152     return (result)
153 #----------------------------------------------------------------------------
154 #Accuracy calculations
155 def Acc (classified,dataset):
156    count = 0
157    for x in range (len(dataset)):
158      if (classified[x]==dataset[x][-1]):
159        count += 1
160    acc = float(count/(len(dataset)))
161    return(acc)
162
163 def bAcc (classified,dataset):
164    one = 0
165    two = 0
166    count1 = 0
167    count2 = 0
168    for x in range (len(dataset)):
169      if (dataset[x][-1]== 2):
170        one += 1
171      else:
172        two += 1
173    for x in range (len(dataset)):
174      if (classified[x]==dataset[x][-1] & classified [x]==2):
175        count1 += 1
```

```
176      elif (classified[x]==dataset[x][-1] & classified [x]==4):
177        count2 += 1
178    bacc = 0.5*((count1/one) + (count2/two))
179    return(bacc)
180 #---------------------------------------------------------------------------
181 #KNN CLASSIFICATION (EUCLIDEAN Distance)
182 print("\nTRAINING MODEL WITH EUCLIDEAN DISTANCE AS THE DISTANCE METRIC")
183 k_values = [1,3,5,7,9,11,13,15,17,19]
184 acc_values = []
185 bacc_values = []
186 for k in k_values:
187    neighbors = []
188    classified = []
189    print("\nRunning for K =", k)
190    for data in dataset_dev:
191      neighbors= (getNeighbors(dataset_train, data, k))
192      classified.append(getResponse (neighbors))
193    print('Accuracy:',Acc(classified,dataset_dev),'\nBalanced Accuracy:',\
194        bAcc(classified,dataset_dev))
195    acc_values.append(Acc(classified, dataset_dev))
196    bacc_values.append(bAcc(classified,dataset_dev))
197 plt.title('Accuracy')
198 plt.bar(k_values,acc_values)
199 plt.show()
200
201 plt.title('Balanced Accuracy')
202 plt.bar(k_values,bacc_values)
203 plt.show()
204 euclidean_metric = [acc_values,bacc_values]
205 print('Best K1 is', acc_values.index(max(acc_values))+2)
206 print('Best K2 is', bacc_values.index(max(bacc_values))+2)
207 #---------------------------------------------------------------------------
208 #KNN using Hamming Distance
209 print("\nTRAINING MODEL WITH HAMMING DISTANCE AS THE DISTANCE METRIC")
210 k_values = [1,3,5,7,9,11,13,15,17,19]
211 acc_values = []
212 bacc_values = []
213 for k in k_values:
214    neighbors = []
215    classified = []
216    print("\nRunning for K =", k)
217    for data in dataset_dev:
218      neighbors= (getNeighbors_hamming(dataset_train, data, k))
219      classified.append(getResponse (neighbors))
220    print('Accuracy:',Acc(classified,dataset_dev),'\nBalanced Accuracy:',\
221        bAcc(classified,dataset_dev))
222    acc_values.append(Acc(classified, dataset_dev))
223    bacc_values.append(bAcc(classified,dataset_dev))
224 plt.title('Accuracy Hamming')
225 plt.bar(k_values,acc_values)
226 plt.show()
227
228 plt.title('Balanced Accuracy Hamming')
229 plt.bar(k_values,bacc_values)
230 plt.show()
231 hamming_metric = [acc_values,bacc_values]
232 print('Best K1 is', acc_values.index(max(acc_values))+2)
233 print('Best K2 is', bacc_values.index(max(bacc_values))+2)
234
235 #---------------------------------------------------------------------------
236 #KNN Using Cosine Similarity
237 k_values = [1,3,5,7,9,11,13,15,17,19]
238 acc_values = []
239 bacc_values = []
240 print("\nTRAINING MODEL WITH COSINE SIMILARITY AS THE DISTANCE METRIC AND k=3")
241 for k in k_values:
242    neighbors = []
243    classified = []
244    print("\nRunning for K =", k)
245    for data in dataset_dev:
246      neighbors= (getNeighbors_cosine(dataset_train, data, k))
247      classified.append(getResponse (neighbors))
248    print('Accuracy:',Acc(classified,dataset_dev),'\nBalanced Accuracy:',\
249        bAcc(classified,dataset_dev))
250    acc_values.append(Acc(classified, dataset_dev))
251    bacc_values.append(bAcc(classified,dataset_dev))
252 print('\n\n\n\n\n')
253 plt.title('Accuracy (Cosine)')
254 plt.bar(k_values,acc_values)
255 plt.show()
256
257 plt.title('Balanced Accuracy (Cosine)')
258 plt.bar(k_values,bacc_values)
259 plt.show()
260 cosine_metric = [acc_values,bacc_values]
261 print('Best K1 is', acc_values.index(max(acc_values))+2)
262 print('Best K2 is', bacc_values.index(max(bacc_values))+2)
263 #
```
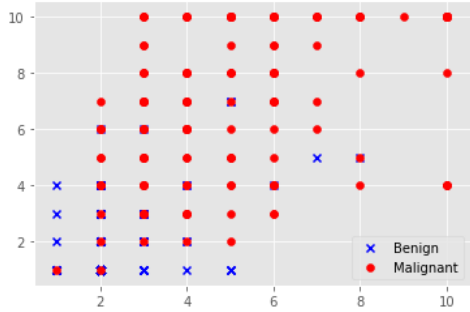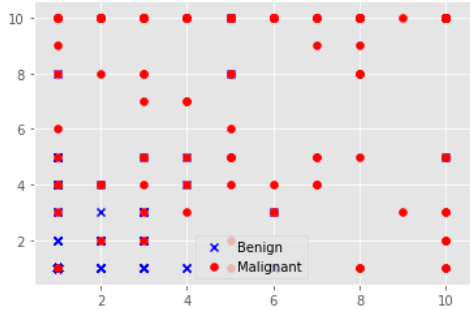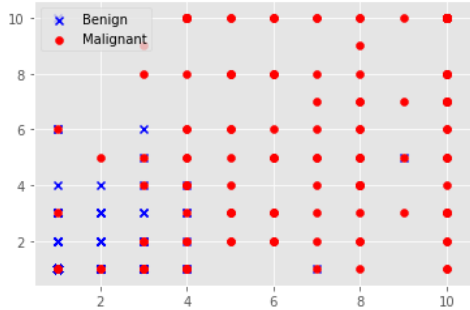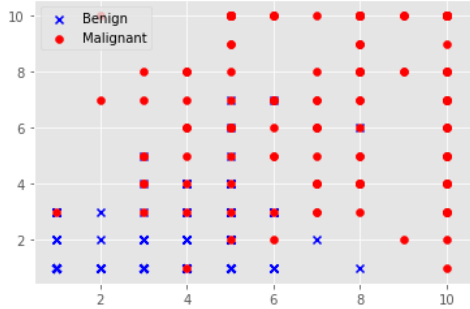
```python
263  #--------------------------------------------------------------------------------------------
264  print("\nCOMPARING THE 3 DISTANCE METRICS!")
265  print("ACC and bACC Values for Euclidean Distance: ",euclidean_metric[0],'\n',euclidean_metric[1])
266  print("ACC and bACC Values for Hamming Distance: ",hamming_metric[0],'\n',hamming_metric[1])
267  print("ACC and bACC Values for Cosine Similarity: ",cosine_metric[0], '\n',cosine_metric[1])
268  #--------------------------------TESTING--------------------------------------
269  #KNN using Hamming Distance
270  print("\nTRAINING MODEL WITH HAMMING DISTANCE AS THE DISTANCE METRIC")
271  k_values = [3]
272  acc_values = []
273  bacc_values = []
274  for k in k_values:
275    neighbors = []
276    classified = []
277    #print("\nRunning for K =", k)
278    for data in dataset_test:
279      neighbors= (getNeighbors_hamming(dataset_train, data, 3))
280      classified.append(getResponse (neighbors))
281    print('Accuracy:',Acc(classified,dataset_test),'\nBalanced Accuracy:',\
282          bAcc(classified,dataset_test))
283
```
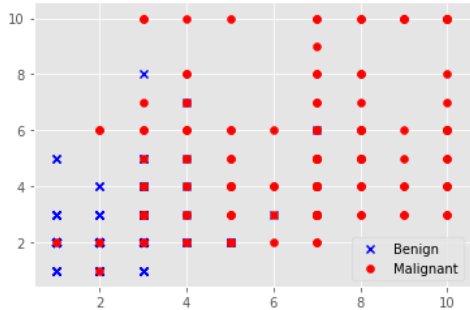
otal Number of Data Points in Training set: 483
otal Number of Data Points in Dev set: 100
otal Number of Data Points in Testing set: 100











he Datapoints are not linearly sperable.

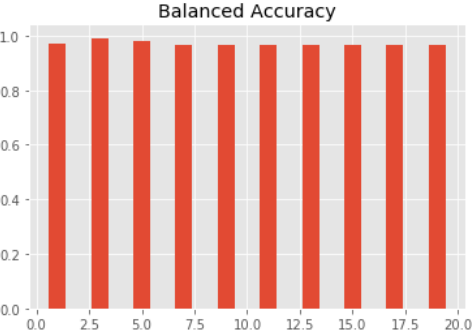RAINING MODEL WITH EUCLIDEAN DISTANCE AS THE DISTANCE METRIC

unning for K = 1
ccuracy: 0.97
alanced Accuracy: 0.9708557255064076

unning for K = 3
ccuracy: 0.99
alanced Accuracy: 0.9915254237288136

unning for K = 5
ccuracy: 0.98
alanced Accuracy: 0.9793303017775941

unning for K = 7
ccuracy: 0.97
alanced Accuracy: 0.9671351798263745

unning for K = 9
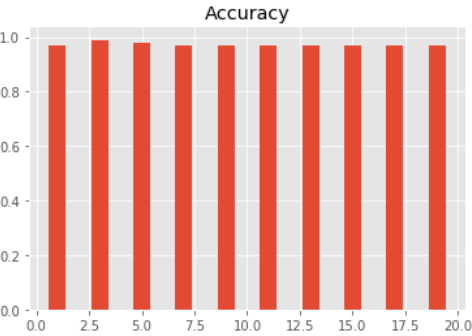ccuracy: 0.97

ccuracy: 0.97
alanced Accuracy: 0.9671351798263745

unning for K = 11
ccuracy: 0.97
alanced Accuracy: 0.9671351798263745

unning for K = 13
ccuracy: 0.97
alanced Accuracy: 0.9671351798263745

unning for K = 15
ccuracy: 0.97
alanced Accuracy: 0.9671351798263745

unning for K = 17
ccuracy: 0.97
alanced Accuracy: 0.9671351798263745

unning for K = 19
ccuracy: 0.97
alanced Accuracy: 0.9671351798263745

### Accuracy



### Balanced Accuracy



est K1 is 3
est K2 is 3

RAINING MODEL WITH HAMMING DISTANCE AS THE DISTANCE METRIC

unning for K = 1
ccuracy: 0.97
alanced Accuracy: 0.9671351798263745

unning for K = 3
ccuracy: 0.99
alanced Accuracy: 0.9915254237288136

unning for K = 5
ccuracy: 0.99
alanced Accuracy: 0.9915254237288136

unning for K = 7
ccuracy: 0.98
alanced Accuracy: 0.9793303017775941

unning for K = 9
ccuracy: 0.99
alanced Accuracy: 0.9915254237288136

unning for K = 11
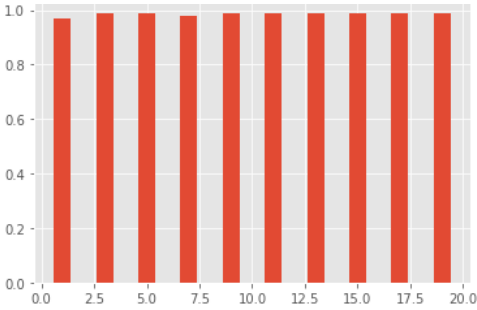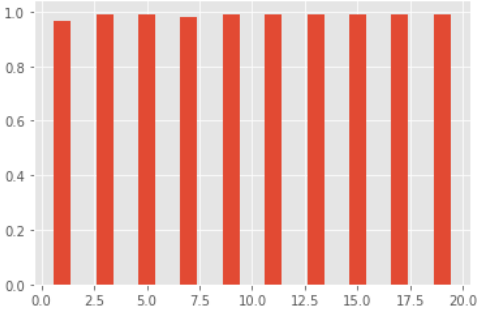ccuracy: 0.99
alanced Accuracy: 0.9915254237288136

unning for K = 13
ccuracy: 0.99
alanced Accuracy: 0.9915254237288136

unning for K = 15
ccuracy: 0.99
alanced Accuracy: 0.9915254237288136

unning for K = 17
ccuracy: 0.99
alanced Accuracy: 0.9915254237288136

unning for K = 19
ccuracy: 0.99
alanced Accuracy: 0.9915254237288136

### Accuracy Hamming

**Balanced Accuracy Hamming**



```
est K1 is 3
est K2 is 3

RAINING MODEL WITH COSINE SIMILARITY AS THE DISTANCE METRIC AND k=3

unning for K = 1
ccuracy: 0.91
alanced Accuracy: 0.9051260851591567

unning for K = 3
ccuracy: 0.89
alanced Accuracy: 0.8770152955766846

unning for K = 5
ccuracy: 0.91
alanced Accuracy: 0.9014055394791236

unning for K = 7
ccuracy: 0.9
alanced Accuracy: 0.8892104175279041

unning for K = 9
ccuracy: 0.9
alanced Accuracy: 0.8892104175279041

unning for K = 11
ccuracy: 0.9
alanced Accuracy: 0.8892104175279041

unning for K = 13
ccuracy: 0.89
alanced Accuracy: 0.8807358412567177

unning for K = 15
ccuracy: 0.91
alanced Accuracy: 0.9051260851591567

unning for K = 17
ccuracy: 0.93
alanced Accuracy: 0.9257957833815627

unning for K = 19
ccuracy: 0.94
alanced Accuracy: 0.9379909053327822
```
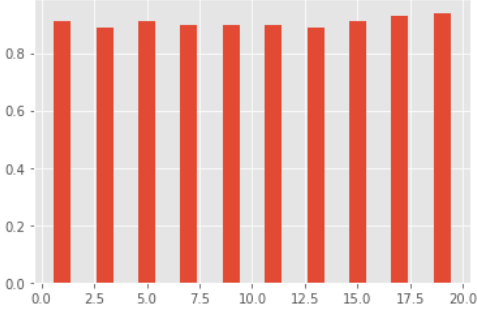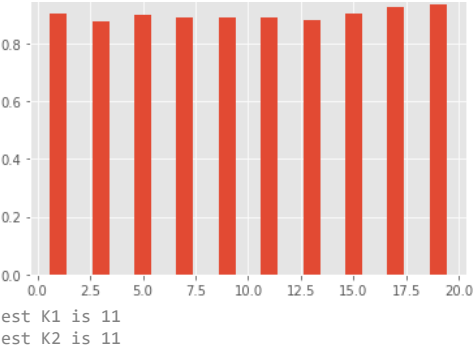
**Accuracy (Cosine)**



**Balanced Accuracy (Cosine)**

```
est K1 is 11
est K2 is 11

OMPARING THE 3 DISTANCE METRICS!
CC and bACC Values for Euclidean Distance:  [0.97, 0.99, 0.98, 0.97, 0.97, 0.97, 0.97, 0.97, 0.97, 0.97]
[0.9708557255064076, 0.9915254237288136, 0.9793303017775941, 0.9671351798263745, 0.9671351798263745, 0.9671351798263745, 0.9671351798263745, 0.96713517
CC and bACC Values for Hamming Distance:  [0.97, 0.99, 0.99, 0.98, 0.99, 0.99, 0.99, 0.99, 0.99, 0.99]
[0.9671351798263745, 0.9915254237288136, 0.9915254237288136, 0.9793303017775941, 0.9915254237288136, 0.9915254237288136, 0.9915254237288136, 0.99152542
CC and bACC Values for Cosine Similarity:  [0.91, 0.89, 0.91, 0.9, 0.9, 0.9, 0.89, 0.91, 0.93, 0.94]
[0.9051260851591567, 0.8770152955766846, 0.9014055394791236, 0.8892104175279041, 0.8892104175279041, 0.8892104175279041, 0.8807358412567177, 0.90512608

RAINING MODEL WITH HAMMING DISTANCE AS THE DISTANCE METRIC
ccuracy: 0.99
alanced Accuracy: 0.9864864864864865
```

# Results

At the end the best set of parameters is K=3 and using Hamming distance as the distance metric. We see an Accuracy if 99% and a Balanced Accuracy of 98.6%.