

# CSCE633 Homework 01

Lu Sun

September 2020

## Question 1

(i) Plot  $\vec{x}_1, \vec{x}_2, \vec{x}_3$  in 2D space. Plot the line of  $\omega(0)$  and the corresponding direction.

Answer:

$\omega(0) = [2, -1, 1]^T \rightarrow 2 - x + y = 0$  The figure is shown as follow:

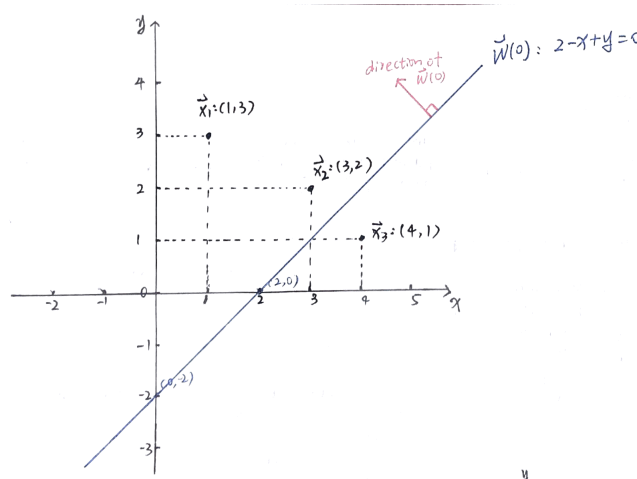


Figure 1: Q1.(i)

(ii) Indicate the class of  $\vec{x}_1, \vec{x}_2, \vec{x}_3$  by  $\omega(0)$ . Point out the incorrectly classified samples.

Answer:

- (a) The results are shown in the following table:
- (b) From the above table, the incorrectly classified sample is  $\vec{x}_2$ .

samples	$\text{sign}(\omega(0)^T x_n)$	predict class	true class ( $y_n$ )
$\vec{x}_1 : (1, 3)$	$\text{sign}(2 - 1 + 3) = \text{sign}(4) = +1$	Class 1	Class 1 (+1)
$\vec{x}_2 : (3, 2)$	$\text{sign}(2 - 3 + 2) = \text{sign}(1) = +1$	Class 1	Class 2 (-1)
$\vec{x}_3 : (4, 1)$	$\text{sign}(2 - 4 + 1) = \text{sign}(-1) = -1$	Class 2	Class 2 (-1)

Table 1: Q1.(ii)

(iii) Find the value of the new weight  $\omega(1)$  with LPA method and above misclassified sample. Find the new line corresponding to (1). Plot  $\omega(1)$  and its direction. Indicate which samples are correctly classified and which are not.

**Answer:**

(a) In (ii), only sample  $\vec{x}_2$  is misclassified. So, in process to get  $\omega(1)$  with PLA method, only  $\vec{x}_2 : (3, 2)$ ,  $y_2 : -1$  and  $\omega(0) : [2, -1, 1]^T$  are needed. Finally,  $\omega(1)$  is got through the following formula:

$$\omega(1) = \omega(0) + y_2 \cdot [1, \vec{x}_2]^T = [2, -1, 1]^T - 1 \cdot [1, 3, 2]^T = [1, -4, -1]^T$$

(b) The new line corresponding to  $\omega(1) : [1, -4, -1]^T$  is  $1 - 4x - y = 0$

(c) The plot is shown as follow:

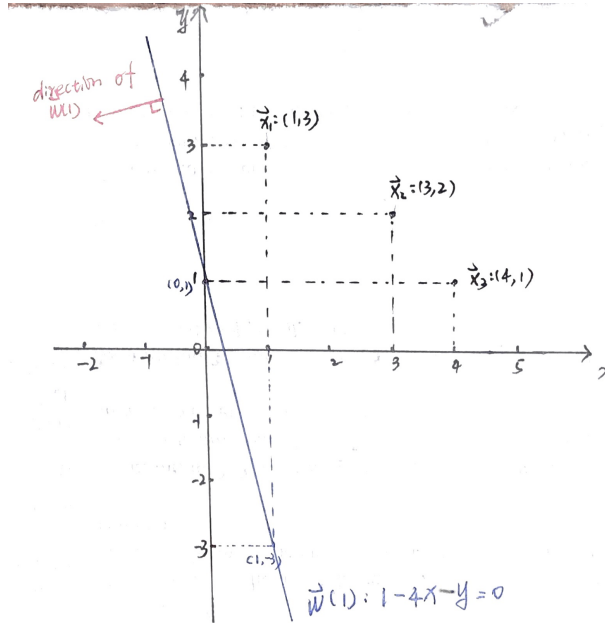


Figure 2: Q1.(iii)

(d) From the above plot, following results can be get:

samples	$\text{sign}(\omega(1)^T x_n)$	predict class	true class ( $y_n$ )	classified
$\vec{x}_1 : (1, 3)$	$\text{sign}(1 - 4 - 3) = \text{sign}(-6) = -1$	Class 2	Class 1 (+1)	Incorrect
$\vec{x}_2 : (3, 2)$	$\text{sign}(1 - 12 - 2) = \text{sign}(-13) = -1$	Class 2	Class 2 (-1)	Correct
$\vec{x}_3 : (4, 1)$	$\text{sign}(1 - 16 - 1) = \text{sign}(-16) = -1$	Class 2	Class 2 (-1)	Correct

Table 2: Q1.(iii)

Namely, sample  $\vec{x}_2, \vec{x}_3$  are correctly classified, while  $\vec{x}_1$  is not.

**(iv) Find the value of the new weight  $\omega(2)$  with LPA method and above misclassified sample. Find the new line corresponding to (2). Plot  $\omega(2)$  and its direction. Indicate which samples are correctly classified and which are not.**

**Answer:**

(a) In (iii), only sample  $\vec{x}_1$  is misclassified. So, in process to get  $\omega(2)$  with PLA method, only  $\vec{x}_1 : (1, 3)$ ,  $y_1 : +1$  and  $\omega(1) : [1, -4, -1]^T$  are needed. Finally,  $\omega(2)$  is got through the following formula:

$$\omega(2) = \omega(1) + y_1 \cdot [1, \vec{x}_1]^T = [1, -4, -1]^T + 1 \cdot [1, 1, 3]^T = [2, -3, 2]^T$$

(b) The new line corresponding to  $\omega(2) : [2, -3, 2]^T$  is  $2 - 3x + 2y = 0$

(c) The plot is shown as follow:

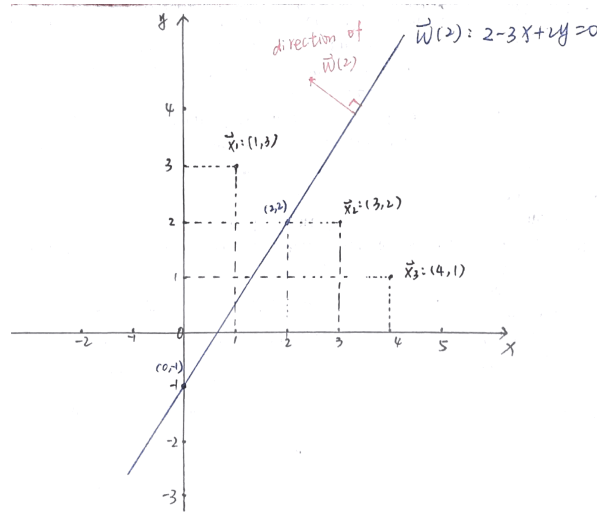


Figure 3: Q1.(iii)

(d) From the above plot, following results can be get:

samples	$sign(\omega(1)^T x_n)$	predict class	true class ( $y_n$ )	classified
$\vec{x}_1 : (1, 3)$	$sign(2 - 3 + 6) = sign(5) = +1$	Class 1	Class 1 (+1)	Correct
$\vec{x}_2 : (3, 2)$	$sign(2 - 9 + 4) = sign(-3) = -1$	Class 2	Class 2 (-1)	Correct
$\vec{x}_3 : (4, 1)$	$sign(2 - 12 + 2) = sign(-8) = -1$	Class 2	Class 2 (-1)	Correct

Table 3: Q1.(iv)

Namely, sample  $\vec{x}_1, \vec{x}_2, \vec{x}_3$  are all correctly classified and none of the samples is incorrectly classified.

## Question 2

(a.i) Compute the number of samples belong to each class.  
Are the classes equally distributed?

**Answer:**

The code for computing the number of samples is as follow: (in the code ,  
'Iris-setosa': 0, 'Iris-versicolor': 1,'Iris-virginica':2)

```
1 Lg = len(datasetTrain)
2 for i in range(3):
3     print(len([j for j in range(Lg) if datasetTrain[j].label==i]))
```

Finally, the following results are shown:

- The numbers of samples belong to 'Iris-setosa': 0, 'Iris-versicolor': 1,'Iris-virginica':2 are 30,30,30 respectively
- The classes equally distributed.

(a.ii) Plot the histogram of each feature. How are the features distributed?

**Answer:**

The code for plots is as follow:

```
1 import matplotlib.pyplot as plt
2 sl = [datasetTrain[i].sepal.length for i in range(Lg)]
3 sw = [datasetTrain[i].sepal.width for i in range(Lg)]
4 pl = [datasetTrain[i].petal.length for i in range(Lg)]
5 pw = [datasetTrain[i].petal.width for i in range(Lg)]
6 mapList = { 0:[sl,'sepal.length'], 1:[sw,'sepal.width'], ...
              2:[pl,'petal.length'], 3:[pw,'petal.width']}
7
8 plt.style.use('ggplot')
9 #plt.grid(True)
10 for i in range(4):
11     plt.hist(mapList[i][0], bins=20, color = 'w',ec = ...
               'darkblue',lw=1.5)
12     plt.xlabel(mapList[i][1])
13     plt.ylabel('Number of train data')
14     plt.show()
```

The following plots can be shown:

From the below figures, we can know that

- Features sepal length and sepal width are unimodal, uniform distributions
- Features petal length and petal width are bimodal distributions

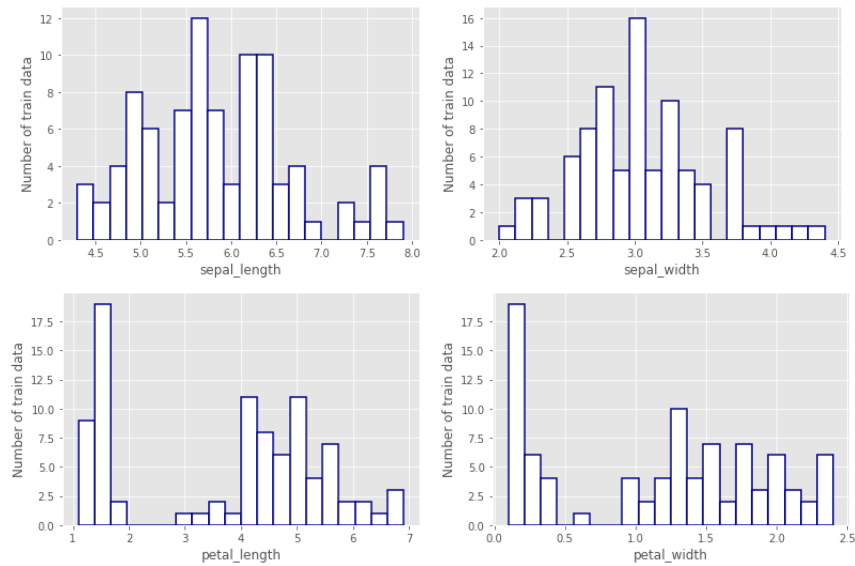


Figure 4: The histogram of each feature

(a.iii) Plot scatter plots. Indicate the class on the plots. What do you observe? How separable do the classes look? Are there feature combinations for which the three classes are more separable?

**Answer:**

The code for plots is as follow:

```

1 def plot.data(xdata,ydata,xlabl,ylabl):
2     #'Iris-setosa': 0, 'Iris-versicolor': 1,'Iris-virginica':2
3     SeVeVi_x = []
4     SeVeVi_y = []
5     for j in range(3):
6         SeVeVi_x.append([xdata[i] for i in range(Lg) if ...
7             datasetTrain[i].label == j]) # List: heights of cats
8         SeVeVi_y.append([ydata[i] for i in range(Lg) if ...
9             datasetTrain[i].label == j]) # List: weights of cats
10
11 plt.scatter(SeVeVi_x[0], SeVeVi_y[0], c='b', marker='x', ...
12     label='Iris-setosa') # put points corresponding to cats, ...
13     b: blue colored, x: marked with 'x'
14 plt.scatter(SeVeVi_x[1], SeVeVi_y[1], c='r', marker='+', ...
15     label='Iris-versicolor') # put points corresponding to ...
16     cats, b: blue colored, x: marked with 'x'
17 plt.scatter(SeVeVi_x[2], SeVeVi_y[2], c='g', marker='1', ...
18     label='Iris-virginica') # put points corresponding to ...
19     cats, b: blue colored, x: marked with 'x'

```

```

12     plt.xlabel(xlab1)
13     plt.ylabel(ylab1)
14     plt.legend()
15     plt.show()
16
17
18     #sl : "sepal.length", sw : "sepal.width", pl : "petal.length", ...
19     pw : "petal.width"
20     #plot_data(sl,sw,'sepal.length','sepal.width')
21
22     import itertools
23     for i, j in itertools.product(range(4),range(4)):
24         if i<j:
25             plot_data(mapList[i][0],mapList[j][0],mapList[i][1],mapList[j][1])

```

The following plots can be shown:

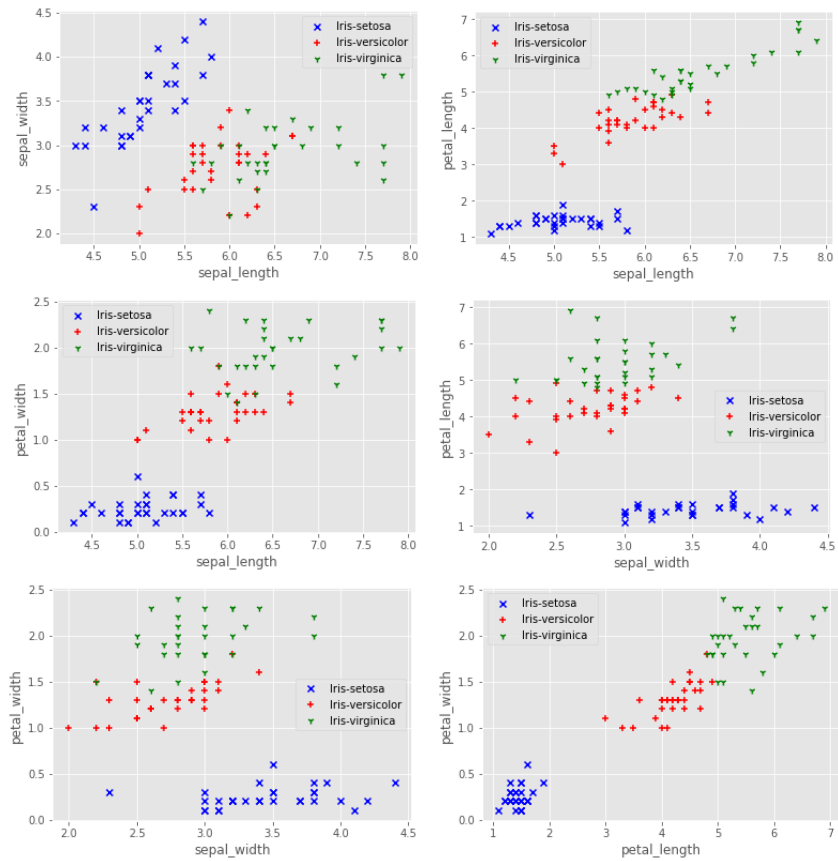


Figure 5: The scatter of combined two features

From the above figures, we can know that:

- Some combination of features can separate the three classes, while some can not.
- Combination of (sepal length, sepal width) can't separate class versicolor and virginica. While all the other five combinations look good in separating the three classes.
- Combination (petal length, petal width) is more separable for the three classes both in x-axis and y-axis.

### (b.i) Implement a KNN using l2-norm

**Answer:**

The pseudocode is shown as follow:

```

1 Knn( test_data, train_data,k):
2     distance = ||test_data - train_data||
3     indx_tep = k smallest points' order in distance list
4     countx = train_data.label[indx_tep]
5     number of label 0 = count(countx==0)
6     number of label 1 = count(countx==1)
7     number of label 2 = count(countx==2)
8     return the largest number's label

```

Here, we don't care about the order of k-th points. Namely, for example, [9,9,1,2,3]'s 1-th largest point is 9, at order 1. Although order 2 is also 9, the largest number, we don't care about that order. Then the code for KNN method is shown as follow:

Listing 1: KNN

```

1 #####Define Norm Function#####
2 def L0Fun(xtest,datatrain):
3     x_y = datatrain-np.array(xtest)
4     #return sum([1 for i in range(len(x_y)) if x_y[i]!=0])
5     return np.count_nonzero(x_y,axis = 1)
6 def L1Fun(xtest,datatrain):
7     x_y = datatrain-np.array(xtest)
8     #return sum(abs(x_y))
9     return np.sum(abs(x_y),axis=1)
10 def L2Fun(xtest,datatrain):
11     x_y = datatrain-np.array(xtest)
12     #return np.sqrt(x_y.dot(x_y))
13     return np.sqrt(np.sum(x_y * x_y, axis=1))
14 def CsFun(xtest,datatrain):
15     x = np.array(xtest)
16     #y = np.array(y)
17     #return x.dot(y) / (np.sqrt(x.dot(x))*np.sqrt(y.dot(y)))
18     xnorm = np.sqrt(x.dot(x))
19     ynorm = np.sqrt(np.sum(datatrain * datatrain, axis=1))
20     return np.sum(datatrain * x,axis = 1)/ (xnorm * ynorm)
21 mapNorm = {0:L0Fun,1:L1Fun,2:L2Fun,3:CsFun}

```



```

22
23
24 #####Define KNN Function#####
25 def KnnFun(xtest, intK, xtrain = datasetTrain, normch = 2):
26     datatrain = np.reshape([xtrain[i].full_vector() for i in
27         range(len(xtrain))], [-1,5])
28
29     #select norm
30     normFun = mapNorm[normch]
31
32     #compute all the distance
33     dist = normFun(xtest, datatrain[0::,0:4])           #total dist list
34
35     #select top k-th nearest point regardless the order
36     indx.tep = np.argpartition(-dist, -intK)[-intK:]    #index
37     countx = [int(datatrain[i,-1]) for i in indx.tep]    #K-th lable list
38
39     esti.tep = np.array([countx.count(i) for i in range(3)])
40     esti = np.argmax(esti.tep)
41
42     return esti

```

(b.ii) Set  $K = 1, 3, 5, 7, \dots, 19$ , train the model and compute the classification accuracy. Plot the ACC with  $K$ . Report the best  $K^*$

**Answer:**

First, the following code for ACC is implemented.

```

1 def AccFun(estiLabel, trueLabel):
2     x-y = np.array(estiLabel)-np.array(trueLabel)
3     Lg_xy = len(x-y)
4     #return sum([1 for i in range(Lg_xy) if x-y[i]==0]) / Lg_xy
5     return np.count.nonzero(x-y==0,axis=0) / Lg_xy, x-y

```

Then, the code for training the model and compute ACC on development data is shown as follow:

```

1 datasetDev = parse_dataset('/content/drive/My ...
    Drive/2020_FALL/CSCE633/HomeWork/HW01/iris_dev.csv')
2
3 intK = np.arange(1,21,2)
4 LgK = len(intK)
5 LgDev = len(datasetDev)
6 estiLabel = np.zeros([LgDev,LgK],dtype = int)
7 trueLabel = np.zeros([LgDev,LgK],dtype = int)
8 for i,k in itertools.product(range(LgDev), range(LgK)):
9     xdev = datasetDev[i].feature_vector()
10    lab = KnnFun(xdev,intK[k])
11    estiLabel[i][k] = lab
12    trueLabel[i][k] = datasetDev[i].label
13 print(accRes)

```

The results of ACC is shown in the following table:

K	1	3	5	7	9
ACC	0.96666667	0.93333333	0.93333333	0.93333333	0.93333333
K	11	13	15	17	19
ACC	0.93333333	0.93333333	0.93333333	0.93333333	0.93333333

Table 4: ACC

The code for plot is shown as follow:

```
1 plt.plot(intK, accRes, '-+')
2 plt.xlabel('choice of K')
3 plt.ylabel('Accuracy Rate')
```

The plot is shown as follow:

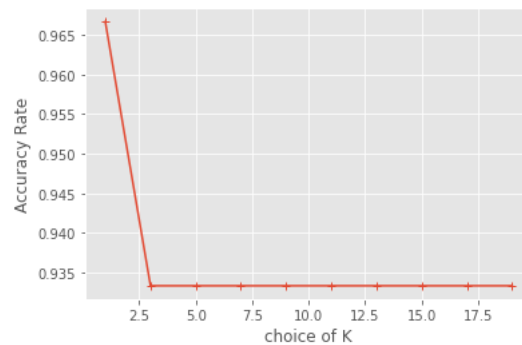


Figure 6: ACC

- From the above table and plot, the best  $K^* = 1$  whose ACC is the highest 0.96666667

**(b.iii) Report the ACC on test set with  $K^*$ . What do you observe?**

**Answer:**

The code is as follow:

```
1 datasetTest = parse_dataset('/content/drive/My ...
    Drive/2020_FALL/CSCE633/HomeWork/HW01/iris.test.csv') # ...
    Returns the list of datapoints
2
3 Kstar = 1
4 LgTest = len(datasetTest)
5 estiLabel = np.zeros([LgTest,1], dtype = int)
```

```

6 trueLabel = np.zeros([LgTest,1],dtype = int)
7 for i in range(LgTest):
8     xdev = datasetTest[i].feature_vector()
9     lab= KnnFun(xdev, Kstar)
10    estiLabel[i] = lab
11    trueLabel[i] = datasetTest[i].label
12 print(accRes)

```

Finally, we get:

- Acc = 0.96666667
- From the observation, the ACC for development set is equal to that of test set with  $K^* = 1$ .

**(b.iv) Experiment with different types of distances. Report your findings**

The plots are shown as follow:

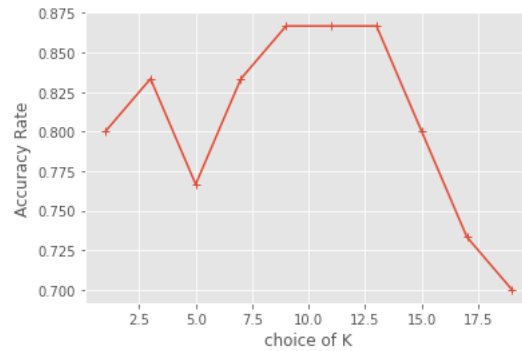


Figure 7: 10-norm

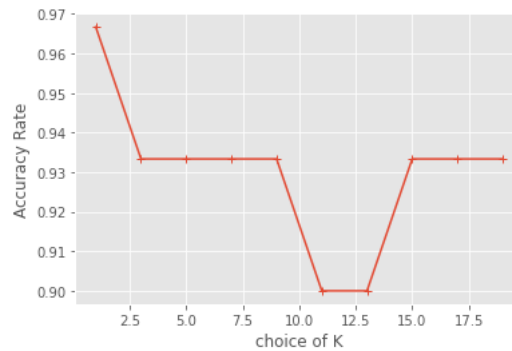


Figure 8: 11-norm

From the above plots, we can know:

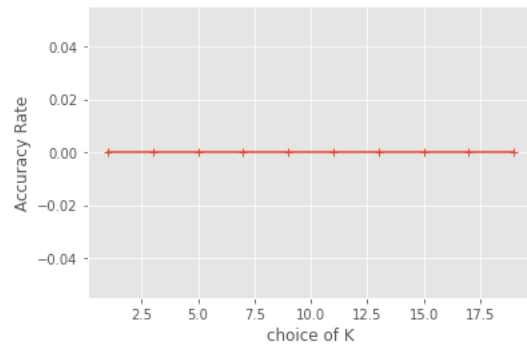


Figure 9: cosine similarity

- l0-norm's best  $K^*$  is 9,11,13 which is larger than that of l2-norm, while the accuracy is smaller.
- l1-norm has the same best  $K^*$  and ACC as that of l2-norm.
- Cosine similarity can't get good results in this example, whose ACC is always 0.

## Appendix

### (i) Code for upload data

```
1 import numpy as np
2
3 class DataPoint(object):
4
5     def __init__(self, feats):
6         """
7         Return a Datapoint object whose attributes are given by "feats"
8         """
9         self.sepal.length= feats['sepal.length']
10        self.sepal.width = feats['sepal.width']
11        self.petal.length= feats['petal.length']
12        self.petal.width = feats['petal.width']
13        self.label = feats['label']
14
15    def feature_vector(self):
16        """
17        Return feature vector as a numpy array
18        """
19        return np.array([self.sepal.length, self.sepal.width, ...
20                          self.petal.length, self.petal.width])
21
22    def feature_vector_with_bias(self):
23        """
24        Return features with 1 appended at the end
25        We 'absorb' bias (b) into the feature vector w by adding one ...
26        additional constant dimension
27        y = wx + b becomes y = w'x' where x' = [x; 1] and w' = [w; b]
28        """
29        return np.array([self.sepal.width, self.sepal.width, ...
30                          self.petal.length, self.petal.width, 1.])
31
32    def full_vector(self):
33        """
34        Return feature vector as a numpy array
35        """
36        return np.array([self.sepal.length, self.sepal.width, ...
37                          self.petal.length, self.petal.width, self.label])
38
39    def __str__(self):
40        """
41        print(object) uses this method
42        Define it in the way you want to see the object (e.g. ...
43        Height:27.0, Weight:50.0, Label:1)
44        """
45        return "sepal.length:{}, sepal.width:{}, petal.length:{}, ...
46                petal.width:{}, Label:{}".format(self.sepal.length, ...
47                self.sepal.width, self.petal.length, self.petal.width, ...
48                self.label)
49
50    def parse_dataset(filename):
51        datafile = open(filename, 'r') # Open File "to read"
52        dataset = [] # List to hold Datapoint objects
```

```

45     label_map = {'Iris-setosa': 0, 'Iris-versicolor': ...
46                 1, 'Iris-virginica': 2} # define labels as 0/1/2.
47     for index, line in enumerate(data_file):
48         if index == 0: # First line: "Height Weight Label" ...
49             describes the datapoint, it's not an actual datapoint,
50             continue # do nothing, it will skip all the following lines
51         #print(line)
52         sepal.length, sepal.width, petal.length, petal.width, label = ...
53         line.strip().split(',') # line.strip().split('\t') # ...
54         strip() removes '\n', and split('\t') splits the line at ...
55         tabs
56         dataset.append(DataPoint({'sepal.length': float(sepal.length), ...
57                                 'sepal.width': float(sepal.width),
58                                 'petal.length': float(petal.length), ...
59                                 'petal.width': float(petal.width),
60                                 'label': label_map[label]})) # ...
61         Create DataPoint object for ...
62         the given data
63
64     print("Total Number of Data Points: {0}".format(len(dataset)))
65     print("Examples: 1 - {0} , 2 - {1}".format(dataset[0], ...
66         dataset[-1]))
67     return dataset
68
69 datasetTrain = parse_dataset('/content/drive/My ...
70                             Drive/2020-FALL/CSCE633/HomeWork/HW01/iris.train.csv') # ...
71                             Returns the list of datapoints

```

## (ii) Code for bonus

```

1  #####10 norm#####
2  intK = np.arange(1, 21, 2)
3  LgK = len(intK)
4  LgDev = len(datasetDev)
5  estiLabel = np.zeros([LgDev, LgK], dtype = int)
6  trueLabel = np.zeros([LgDev, LgK], dtype = int)
7  for i, k in itertools.product(range(LgDev), range(LgK)):
8      xdev = datasetDev[i].feature_vector()
9      lab = KnnFun(xdev, intK[k], normch = 0)
10     estiLabel[i][k] = lab
11     trueLabel[i][k] = datasetDev[i].label
12 #estiLabel = np.reshape(estiLabel, [-1, 1])
13 #trueLabel = np.reshape(trueLabel, [-1, 1])
14 #print(estiLabel, trueLabel)
15 accRes, _ = AccFun(estiLabel, trueLabel)
16 print(accRes)
17
18 plt.plot(intK, accRes, '-+')
19 plt.xlabel('choice of K')
20 plt.ylabel('Accuracy Rate')
21 plt.show()
22 #####11 norm#####
23 intK = np.arange(1, 21, 2)
24 LgK = len(intK)
25 LgDev = len(datasetDev)

```

```

26 estiLabel = np.zeros([LgDev,LgK],dtype = int)
27 trueLabel = np.zeros([LgDev,LgK],dtype = int)
28 for i,k in itertools.product(range(LgDev),range(LgK)):
29     xdev = datasetDev[i].feature_vector()
30     lab = KnnFun(xdev, intK[k],normch=1)
31     estiLabel[i][k] = lab
32     trueLabel[i][k] = datasetDev[i].label
33 #estiLabel = np.reshape(estiLabel,[-1,1])
34 #trueLabel = np.reshape(trueLabel,[-1,1])
35 #print(estiLabel,trueLabel)
36 accRes,_ =AccFun(estiLabel,trueLabel)
37 print(accRes)
38
39 plt.plot(intK,accRes,'-+')
40 plt.xlabel('choice of K')
41 plt.ylabel('Accuracy Rate')
42 plt.show()
43 #####cosine similarity#####
44 intK = np.arange(1,21,2)
45 LgK = len(intK)
46 LgDev = len(datasetDev)
47 estiLabel = np.zeros([LgDev,LgK],dtype = int)
48 trueLabel = np.zeros([LgDev,LgK],dtype = int)
49 for i,k in itertools.product(range(LgDev),range(LgK)):
50     xdev = datasetDev[i].feature_vector()
51     lab = KnnFun(xdev, intK[k],normch=3)
52     estiLabel[i][k] = lab
53     trueLabel[i][k] = datasetDev[i].label
54 #estiLabel = np.reshape(estiLabel,[-1,1])
55 #trueLabel = np.reshape(trueLabel,[-1,1])
56 #print(estiLabel,trueLabel)
57 accRes,_ =AccFun(estiLabel,trueLabel)
58 print(accRes)
59
60 plt.plot(intK,accRes,'-+')
61 plt.xlabel('choice of K')
62 plt.ylabel('Accuracy Rate')
63 plt.show()

```

### (iii) Code for KNN with order considered

```

1 def KnnFun(xtest, intK, xtrain = datasetTrain, normch = 2):
2     datatrain = np.reshape([xtrain[i].full_vector() for i in ...
3         range(len(xtrain))],[-1,5])
4     normFun = mapNorm[normch]
5     #compute all the distance
6     dist = normFun(xtest,datatrain[0::,0:4]) ...
7     #total dist list
8     #select top k-th nearest point regardless equal nearest points' ...
9     order
10    indx.tep = np.argpartition(-dist,-intK)[-intK:] #index
11    #indx = dist[indx.tep] ...
12    #intK-th dist list
13    #print(dist,'\n',indx) ...

```

```

12     #check top intK-th
countx = [int(datatrain[i,-1]) for i in indx.tep] ...
    #intK-th lable list
13 #print(countx)
14
15 esti.tep = np.array([countx.count(i) for i in range(3)])
16 esti = np.argmax(esti.tep)
17 #print(esti.tep,esti) ...
    #check lable selection
18
19 """
20 #select top k-th nearest point and record possible equal ...
    nearest top k-th points
21 indx = []
22 indist = []
23 while (len(indx) <intK):
24     indx.tep = np.where(dist==min(dist))[0]
25     #if len(indx.tep) == 1:
26     #    indx.append(datatrain[indx.tep][0])
27     #else:
28     for i in range(len(indx.tep)):
29         indx.append(datatrain[indx.tep[i]])
30         indist.append(dist[indx.tep[i]])
31     datatrain = np.delete(datatrain,indx.tep,axis = 0)
32     dist = np.delete(dist,indx.tep,axis = 0)
33
34 if len(indx) > intK:
35     print("# nearest points is {0} larger than ...
        {1}".format(len(indx),intK))
36
37 countx = [indx[i][4] for i in range(len(indx))]
38 esti.tep = np.array([countx.count(i) for i in range(3)])
39 esti = np.where(esti.tep == max(esti.tep))[0]
40 if len(esti) > 1:
41     indistnew = np.where(indist!=indist[-1])[0]
42     indx = [indx[i] for i in indistnew]
43     countx = [indx[i][4] for i in range(len(indx))]
44     esti.tep = np.array([countx.count(i) for i in range(3)])
45     esti = np.where(esti.tep == max(esti.tep))[0]
46     print("Erro: More than one opt estimator")
47 return esti ,indx, countx
48 """
49 return esti

```