

# MACHINE LEARNING HOMEWORK SHEET-1

K-NN & DT.

Name: MAHALAKSHMI  
SABANAYAGAM

TUM ID: ge73yww  
DUE DATE: 29.10.2018.

Problem 1:

Build DT 7 for data  $\{x, y\}$ , using Gini Index, depth-2.

Arranging the given data  $\{x, y\}$  in ascending order with respect to  $x_1$  and then  $x_2$  and  $x_3$  (if  $x_2$  values are equal)  $\rightarrow$  (if  $x_3$  values are equal)

i	$x_1$	$x_2$	$x_3$	y
N	0.4	0.1	4.3	1
J	0.5	0	2.3	1
M	1	0.1	2.8	1
H	1.3	-0.2	1.8	1
O	2.7	-0.5	4.2	1
G	4.1	0.3	5.1	1
I	4.5	0.4	2	0
A	5.5	0.5	4.5	2
K	5.9	-0.1	4.4	0
C	5.9	0.2	3.4	2
F	6.8	-0.3	5.1	2
E	6.9	-0.1	0.6	2
B	7.4	1.1	3.6	0
L	9.3	-0.2	3.2	0
D	9.9	0.1	0.8	0

(possible 1st level cut)

(possible 1st level cut)

(possible 2nd level cut)

$$\text{Gini Index}_{i_g(t)} = 1 - \sum_{a \in C} \pi_a^2$$

@ depth 0 : (5, 6, 4)

$$\text{Gini}(5, 6, 4) = 1 - \left[ \left( \frac{5}{15} \right)^2 + \left( \frac{6}{15} \right)^2 + \left( \frac{4}{15} \right)^2 \right] = 0.6578$$

The trees generated by splitting at  $x_1 \leq 4.1$  and  $x_1 \leq 4.5$  are as follows,

$T_1$ : split @  $x_1 \leq 4.1$

$(5, 6, 4)$   $i_q(t) = 0.6578$

$x_1 \leq 4.1$

$(0, 6, 0)$

$(5, 0, 4)$

$i_q(t_L) \text{ gini}(0, 6, 0)$

$$= 1 - 1 = 0$$

$i_q(t_R) \text{ gini}(5, 0, 4)$

$$= 1 - \left( \frac{25+16}{81} \right)$$

$$= 0.4938$$

$T_2$ : split @  $x_1 \leq 4.5$

$(5, 6, 4)$   $i_q(t) = 0.6578$

$x_1 \leq 4.5$

$(1, 6, 0)$

$(4, 0, 4)$

$i_q(t_L) \text{ gini}(1, 6, 0)$

$$= 1 - \left( \frac{1+36}{49} \right)$$

$$= 0.2449$$

$i_q(t_R) \text{ gini}(4, 0, 4)$

$$= 1 - \left( \frac{16+16}{64} \right)$$

$$= 0.5$$

Calculating improvement at level 1 to decide between  $T_1$  &  $T_2$ ,

Improvement for split  $s$  of  $t$  into  $t_L$  &  $t_R$ ,

$$\Delta i(s; t) = i(t) - p_L i(t_L) - p_R i(t_R)$$

Improvement of  $T_1$

$$\Delta i(x_1 \leq 4.1, T_1) = 0.6578 - \frac{6}{15}(0) - \frac{9}{15}(0.4938) = 0.36152$$

Improvement of  $T_2$

$$\Delta i(x_1 \leq 4.5, T_2) = 0.6578 - \frac{7}{15}(0.2449) - \frac{8}{15}(0.5) = 0.27685$$

$$\Delta i(x_1 \leq 4.1, T_1) > \Delta i(x_1 \leq 4.5, T_2)$$

So  $T_1$  is better than  $T_2$ .

For 2nd level,

splitting @  $x_1 \leq 6.9$

$(5, 6, 4)$   $i_q(t) = 0.6578$

$x_1 \leq 4.1$

$(0, 6, 0)$   $i_q(t_L) = 0$

$(5, 0, 4)$   $i_q(t_R) = 0.4938$

$x_1 \leq 6.9$

$(2, 0, 4)$

$(3, 0, 0)$

$$i_q(t_{RL}) = 1 - \frac{4+16}{36} = 0.4444$$

$$i_q(t_{RR}) = 1 - 1 = 0$$

(Final Tree)

Problem 2:

$$\vec{x}_a = (4.1, -0.1, 2.2)^T$$

Using the generated  $T$  in problem 1,  $x_a$  will be classified as  $\hat{y}_a = 1$

$$p(C = \hat{y}_a | \vec{x}_a, T) = \frac{6}{6} = 1$$

$$\vec{x}_b = (6.1, 0.4, 1.3)^T$$

$$\hat{y}_b = 2$$

$$p(C = \hat{y}_b | \vec{x}_b, T) = \frac{4}{6} = 0.6667$$

K-NEAREST NEIGHBORS,

Problem 4, and 5:

Please check the ipynb. Both the problems are solved programmatically at the end of the given notebook.

Problem 6:

The scales of the features  $(x_1, x_2, x_3)$  of  $X$  are different.

Mean and variance of,

$$x_1 \quad \begin{matrix} \mu \\ \approx 5 \end{matrix} \quad \begin{matrix} \sigma^2 \\ \approx 9 \end{matrix}$$

$$x_2 \quad \approx 0.1 \quad \approx 0.14$$

$$x_3 \quad \approx 3 \quad \approx 2$$

As the scale of  $x_2$  is very less compared to  $x_1$  and  $x_3$ ,  $x_2$  is insignificant when  $kNN$  is used. So, the useful information held by  $x_2$  will be lost.

To avoid this, we have to scale each feature to zero mean and unit variance. So, standardising / normalising the data will help.

We can also use other distance metric like Mahalanobis.

This problem will not come up when using Decision Tree as it is not dependent on scales of the features.



Programming assignment 1: k-Nearest Neighbors classification

In [1]: import numpy as np
import math
from sklearn import datasets, model\_selection
import matplotlib.pyplot as plt
%matplotlib inline
import pandas as pd

Introduction

For those of you new to Python, there are lots of tutorials online, just pick whichever you like best :)

If you never worked with Numpy or Jupyter before, you can check out these guides

- https://docs.scipy.org/doc/numpy-dev/user/quickstart.html
- http://jupyter.readthedocs.io/en/latest/

Your task

In this notebook code to perform k-NN classification is provided. However, some functions are incomplete. Your task is to fill in the missing code and run the entire notebook.

In the beginning of every function there is docstring, which specifies the format of input and output. Write your code in a way that adheres to it. You may only use plain python and numpy functions (i.e. no scikit-learn classifiers).

Exporting the results to PDF

Once you complete the assignments, export the entire notebook as PDF and attach it to your homework solutions. The best way of doing that is

- Run all the cells of the notebook.
- Download the notebook in HTML (click File > Download as > .html)
- Convert the HTML to PDF using e.g. https://www.sejda.com/html-to-pdf or wkhtmltopdf for Linux (tutorial)
- Concatenate your solutions for other tasks with the output of Step 3. On a Linux machine you can simply use pdffunite. there are similar tools for other platforms too. You can only upload a single PDF file to Moodle.

This way is preferred to using nbconvert, since nbconvert clips lines that exceed page width and makes your code harder to grade.

Load dataset

The iris data set (https://en.wikipedia.org/wiki/Iris\_flower\_data\_set) is loaded and split into train and test parts by the function load\_dataset.

In [2]: def load\_dataset(split):
"""Load and split the dataset into training and test parts.

Parameters
-----
split : float in range (0, 1)
 Fraction of the data used for training.

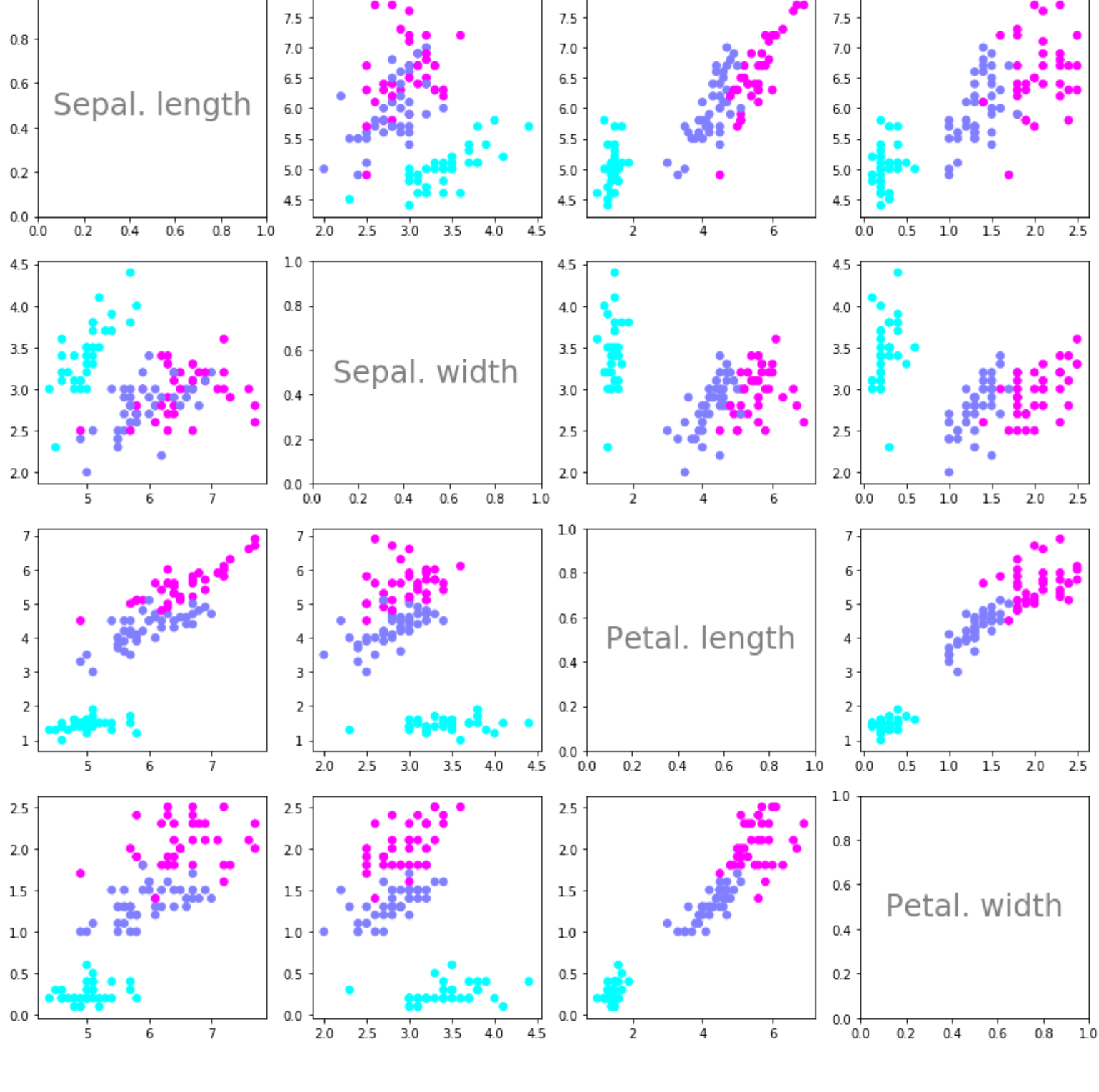
Returns
-----
X\_train : array, shape (N\_train, 4)
 Training features.
y\_train : array, shape (N\_train)
 Training labels.
X\_test : array, shape (N\_test, 4)
 Test features.
y\_test : array, shape (N\_test)
 Test labels.
"""
dataset = datasets.load\_iris()
X, y = dataset['data'], dataset['target']
X\_train, X\_test, y\_train, y\_test = model\_selection.train\_test\_split(X, y, random\_state=123, test\_size=(1 - split))
return X\_train, X\_test, y\_train, y\_test

In [3]: # prepare data
split = 0.75
X\_train, X\_test, y\_train, y\_test = load\_dataset(split)

Plot dataset

Since the data has 4 features, 16 scatterplots (4x4) are plotted showing the dependencies between each pair of features.

In [4]: f, axes = plt.subplots(4, 4, figsize=(15, 15))
for i in range(4):
 for j in range(4):
 if j == 0 and i == 0:
 axes[i,j].text(0.5, 0.5, 'Sepal. length', ha='center', va='center', size=24, alpha=.5)
 elif j == 1 and i == 1:
 axes[i,j].text(0.5, 0.5, 'Sepal. width', ha='center', va='center', size=24, alpha=.5)
 elif j == 2 and i == 2:
 axes[i,j].text(0.5, 0.5, 'Petal. length', ha='center', va='center', size=24, alpha=.5)
 elif j == 3 and i == 3:
 axes[i,j].text(0.5, 0.5, 'Petal. width', ha='center', va='center', size=24, alpha=.5)
 else:
 axes[i,j].scatter(X\_train[:,j],X\_train[:,i], c=y\_train, cmap=plt.cm.cool)



Task 1: Euclidean distance

Compute Euclidean distance between two data points.

In [5]: def euclidean\_distance(x1, x2):
"""Compute Euclidean distance between two data points.

Parameters
-----
x1 : array, shape (4)
 First data point.
x2 : array, shape (4)
 Second data point.

Returns
-----
distance : float
 Euclidean distance between x1 and x2.
"""
# TODO
d = sum(np.array(x1-x2)\*\*2)
distance = math.sqrt(d)
return distance

Task 2: get k nearest neighbors' labels

Get the labels of the k nearest neighbors of the datapoint x\_new.

In [6]: def get\_neighors\_labels(X\_train, y\_train, x\_new, k):
"""Get the labels of the k nearest neighbors of the datapoint x\_new.

Parameters
-----
X\_train : array, shape (N\_train, 4)
 Training features.
y\_train : array, shape (N\_train)
 Training labels.
x\_new : array, shape (4)
 Data point for which the neighbors have to be found.
k : int
 Number of neighbors to return.

Returns
-----
neighbors\_labels : array, shape (k)
 Array containing the labels of the k nearest neighbors.
"""
# TODO
distance\_new = []
#Calculate Euclidean distance of the new data point with all the training data
for x in X\_train:
 distance\_new.append(euclidean\_distance(x,x\_new))
#Create a sorted list of dist,label
distance\_label = list(zip(distance\_new,y\_train))
sorted\_distance = sorted(distance\_label, key=lambda distance\_label: distance\_label[0])
#Return the first k labels (second column of the created list)
neighbours\_labels = np.array(sorted\_distance)[:k]
labels = [int(x) for x in neighbours\_labels ] #Making sure the labels are int
return labels[:k]

Task 3: get the majority label

For the previously computed labels of the k nearest neighbors, compute the actual response. I.e. give back the class of the majority of nearest neighbors. In case of a tie, choose the "lowest" label (i.e. the order of tie resolutions is 0 > 1 > 2).

In [7]: def get\_response(neighbors\_labels, num\_classes=3):
"""Predict label given the set of neighbors.

Parameters
-----
neighbors\_labels : array, shape (k)
 Array containing the labels of the k nearest neighbors.
num\_classes : int
 Number of classes in the dataset.

Returns
-----
y : int
 Majority class among the neighbors.
"""
# TODO
class\_votes = np.zeros(num\_classes)
#Count the number of instances of each class label in the neighbors\_label
for i in range(num\_classes):
 class\_votes[i] = list(neighbors\_labels).count(i)
#Get the max in that and return the class it belongs to
majority\_class = np.where(class\_votes==max(class\_votes))
#There could be tie cases, so the majority class sometimes have a set of numbers.
#So choose the first one in that as that is what is needed (lowest label)
return int(majority\_class[0][0])

Task 4: compute accuracy

Compute the accuracy of the generated predictions.

In [8]: def compute\_accuracy(y\_pred, y\_test):
"""Compute accuracy of prediction.

Parameters
-----
y\_pred : array, shape (N\_test)
 Predicted labels.
y\_test : array, shape (N\_test)
 True labels.
"""
# TODO
correct\_pred = 0
#Subtract each element in the array, the 0 ones in the result are the correct predictions
#So count those and divide by the total number of tested instances
pred = y\_pred - y\_test
correct\_pred = list(pred).count(0)
accuracy = correct\_pred/len(y\_pred)
return accuracy

In [9]: # This function is given, nothing to do here.
def predict(X\_train, y\_train, X\_test, k):
"""Generate predictions for all points in the test set.

Parameters
-----
X\_train : array, shape (N\_train, 4)
 Training features.
y\_train : array, shape (N\_train)
 Training labels.
X\_test : array, shape (N\_test, 4)
 Test features.
k : int
 Number of neighbors to consider.

Returns
-----
y\_pred : array, shape (N\_test)
 Predictions for the test data.
"""
y\_pred = []
for x\_new in X\_test:
 neighbors = get\_neighors\_labels(X\_train, y\_train, x\_new, k)
 y\_pred.append(get\_response(neighbors))
return y\_pred

Testing

Should output an accuracy of 0.9473684210526315.

In [10]: # prepare data
split = 0.75
X\_train, X\_test, y\_train, y\_test = load\_dataset(split)
print('Training set: {0} samples'.format(X\_train.shape[0]))
print('Test set: {0} samples'.format(X\_test.shape[0]))

# generate predictions
k = 3
y\_pred = predict(X\_train, y\_train, X\_test, k)
accuracy = compute\_accuracy(y\_pred, y\_test)
print('Accuracy = {0}'.format(accuracy))

Training set: 112 samples
Test set: 38 samples
Accuracy = 0.9473684210526315

Problem 4 and 5

In [11]: #Load the dataset from the csv file and create an array of x1, x2, x3 and another array of y (cls)
data = pd.read\_csv("01\_homework\_dataset.csv")
cls = data['x1'].tolist()
x1 = data['x1'].tolist()
x2 = data['x2'].tolist()
x3 = data['x3'].tolist()
x = []
for i in range(len(x1)):
 l = [x1[i],x2[i],x3[i]]
 x.append(l)
x, cls

Out[11]: ([[5.5, 0.5, 4.5],
[7.4, 1.1, 3.6],
[5.9, 0.2, 3.4],
[9.9, 0.1, 0.8],
[6.9, -0.1, 0.6],
[6.8, -0.3, 5.1],
[4.1, 0.3, 5.1],
[1.3, -0.2, 1.8],
[4.5, 0.4, 2.0],
[0.5, 0.0, 2.3],
[5.9, -0.1, 4.4],
[9.3, -0.2, 3.2],
[1.0, 0.1, 2.8],
[0.4, 0.1, 4.3],
[2.7, -0.5, 4.2]],
[2, 0, 2, 0, 2, 2, 1, 1, 0, 1, 0, 0, 1, 1, 1])

In [12]: k = 3
X\_test = [[4.1,-0.1,1.2,2],[6.1,0.4,1.3]]
y\_pred = predict(x, cls, np.asarray(X\_test), k)
y\_pred

Out[12]: [0, 2]

Answer Problem 4:

(Lowest class label is chosen in case of tie)

xa is classified as 0 and xb is classified as 2

In [13]: #This function returns the distance with the class label
def get\_neighors\_distance\_regression(X\_train, y\_train, x\_new, k):
distance\_new = []
x\_new = np.asarray(x\_new)
#Calculate Euclidean distance between the new point and the training data
for x in X\_train:
 x = np.asarray(x)
 distance\_new.append(euclidean\_distance(x,x\_new))
#Add the class label with the distance, sort and return the first k values
distance\_label = list(zip(distance\_new,y\_train))
sorted\_distance = sorted(distance\_label, key=lambda distance\_label: distance\_label[0])
return np.asarray(sorted\_distance[:k])

In [14]: #This function computes the normalization term (1/2 in the knn regression equation
def get\_normalization(dist):
return np.reciprocal(sum(np.reciprocal(dist)))

In [15]: #This function computes sum of yi/d(x,xi)
def get\_y(dist\_label):
dist\_inv = np.reciprocal(dist\_label[:0])
lab = dist\_label[:1]
return sum(dist\_inv \* lab)

In [16]: #Calculates y
def get\_response\_regression(distance\_label):
term1 = get\_normalization(distance\_label[:0])
term2 = get\_sum(distance\_label)
return term1\*term2

Answers for Problem 5

In [17]: xa = [4.1,-0.1,1.2,2]
dist\_lab = get\_neighors\_distance\_regression(x, cls, xa, k)
print("3 Nearest Neighbors distances and labels \n", dist\_lab)
print("The predicted class is", get\_response\_regression(dist\_lab))

3 Nearest Neighbors distances and labels
[[0.67082039 0. ]
 [2.18403297 2. ]
 [2.47386339 1. ] ]
The predicted class is 0.5610164259744004

In [18]: xb = [6.1,0.4,1.3]
dist\_lab = get\_neighors\_distance\_regression(x, cls, xb, k)
print("3 Nearest Neighbors distances and labels \n", dist\_lab)
print("The predicted class is", get\_response\_regression(dist\_lab))

3 Nearest Neighbors distances and labels
[[1.17473401 2. ]
 [1.74642492 0. ]
 [2.11896201 2. ] ]
The predicted class is 1.39592451328945