

## HW 2

Please start working on this assignment as soon as possible. If you are a beginner in Python this might take a long time. The **objective** of this assignment is to help you familiarize w python packages related to machine learning, namely scikit-learn package.

**DO NOT ERASE MARKDOWN CELLS AND INSTRUCTIONS IN YOUR HW submission**

## Instructions

This assignment covers several aspects of KNN Classifier and performance evaluation we have covered in [m1 \(../practice/m1/README.md\)](#) module. eep the following in mind:

- Structure your [notebook \(https://git.txstate.edu/ML/2022Spring/blob/master/tutorials/notebook-checklist.md\)](https://git.txstate.edu/ML/2022Spring/blob/master/tutorials/notebook-checklist.md) cells as sugested
- **Q** - QUESTION posted in a markdown cell
  - it explains the task in details
  - it is marked with **Q1**, ... **Q10** ...
- **A** - Marks the location where you need to enter your answer below
  - it can be `python code` (more often) or markdown cell (less often)
  - it is marked with **A1**, ... **A10** ... and you enter your answers **below**
  - make sure the cell is running and produces no errors
- Before you submit the HW:
  - Make sure your notebook can always be rerun from top to bottom.
- Follow [README.md \(README.md\)](#) for homework submission instructions

## Tutorials

- [KNN with sklearn \(https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html)
- [Confusion Matrix \(https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion\\_matrix.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html)
- [Plot Confursion Matrix with Sklearn \(https://scikit-learn.org/stable/auto\\_examples/model\\_selection/plot\\_confusion\\_matrix.html\)](https://scikit-learn.org/stable/auto_examples/model_selection/plot_confusion_matrix.html)

## 1. CLASSIFICATION USING KNN ALGORITHM

**Data** Get the exploratory data and the folowwing files:

<http://archive.ics.uci.edu/ml/machine-learning-databases/statlog/heart/>

- Link contains the original data and the metadata both
- copy them in your HW folder

- If you are using command line: `>> wget http://archive.ics.uci.edu/ml/machine-learning-databases/statlog/heart/heart.dat`
  - If wget is not working
    - download it from [link \(https://eternallybored.org/misc/wget/\)](https://eternallybored.org/misc/wget/)
    - follow [steps \(https://stackoverflow.com/questions/29113456/wget-not-recognized-as-internal-or-external-command\)](https://stackoverflow.com/questions/29113456/wget-not-recognized-as-internal-or-external-command)

**Q1** use pandas to read heart.dat

- NOTE : use separator as space while reading this data
- Use column names from metadata in given order
- NOTE : YOU WON'T SEE 'PRESENCE' in metadata (in attribute information)

**A1** Replace the ? mark with your answer

```
In [25]: import pandas as pd

columns = ["age", "sex", "chestpain", "bp", "cholestrel", "sugar", "ecg", "
df = pd.read_csv("heart.dat", names=columns, sep=' ')
```

**Q2**

1. Have a look at head and tail of your data

- N.B: You can use `.tail` and `.head` methods
- N.B: Print both of them, if you just run `without printing` only output from last command will be printed

2. Let us view the size of dataset as well

- print data shape
3. Now let us see if there is some missing value
4. If there is any na values drop it

**A2** Replace ??? with code in the code cell below

```
In [18]: # Code goes below
df.head()

df.tail()

df.shape

df.isna().sum().sum()
```

```
Out[18]:
```

	age	sex	bp	cholestrel	sugar	heartrate	angina	oldpeak	vessels	chestpain-1.0	.
0	70.0	1.0	130.0	322.0	0.0	109.0	0.0	2.4	3.0	0	.
1	67.0	0.0	115.0	564.0	0.0	160.0	0.0	1.6	0.0	0	.

	age	sex	bp	cholestrel	sugar	heartrate	angina	oldpeak	vessels	chestpain-1.0	.
2	57.0	1.0	124.0	261.0	0.0	141.0	0.0	0.3	0.0	0	.
3	64.0	1.0	128.0	263.0	0.0	105.0	1.0	0.2	1.0	0	.
4	74.0	0.0	120.0	269.0	0.0	121.0	1.0	0.2	1.0	0	.
...	...	...	...	...	...	...	...	...	...	...	.
265	52.0	1.0	172.0	199.0	1.0	162.0	0.0	0.5	0.0	0	.
266	44.0	1.0	120.0	263.0	0.0	173.0	0.0	0.0	0.0	0	.
267	56.0	0.0	140.0	294.0	0.0	153.0	0.0	1.3	0.0	0	.
268	57.0	1.0	140.0	192.0	0.0	148.0	0.0	0.4	0.0	0	.
269	67.0	1.0	160.0	286.0	0.0	108.0	1.0	1.5	3.0	0	.

**Q3** Now we will look deeper into the dataset

- Use pairplot from sns to plot this data frame
- See the statistics of the data by describing dataframe

**A3** Replace ??? with code in the code cell below

```
In [17]: ▶ import seaborn as sns

sns.set(style="ticks", color_codes=True)
g = sns.pairplot(df)

import matplotlib.pyplot as plt
plt.show()
```

```
#describe dataframe
```



```
Out[17]:
```

	age	sex	bp	cholestrel	sugar	heartrate	angina	
count	270.000000	270.000000	270.000000	270.000000	270.000000	270.000000	270.000000	2
mean	54.433333	0.677778	131.344444	249.659259	0.148148	149.677778	0.329630	
std	9.109067	0.468195	17.861608	51.686237	0.355906	23.165717	0.470952	
min	29.000000	0.000000	94.000000	126.000000	0.000000	71.000000	0.000000	
25%	48.000000	0.000000	120.000000	213.000000	0.000000	133.000000	0.000000	
50%	55.000000	1.000000	130.000000	245.000000	0.000000	153.500000	0.000000	

	age	sex	bp	cholestorel	sugar	heartrate	angina
75%	61.000000	1.000000	140.000000	280.000000	0.000000	166.000000	1.000000
max	77.000000	1.000000	200.000000	564.000000	1.000000	202.000000	1.000000

**Q4** If you go through metadata (Attribute Information:) you will see that all data in our dataframe are not of same types.

- So we should deal them accordingly.
- We don't have to do anything to 'real' data. However we have to deal with ordered data and nominal data
- We only need to convert all nominal and ordered data to dummy variables

**A4** Replace ??? with code in the code cell below

In [14]:

```
dummy_list = ['chestpain', 'slope', 'ecg', 'thal']
df = pd.get_dummies(df, columns=dummy_list, prefix=['chestpain', 'slope', 'ecg', 'thal'])
df.head()
```

Out[14]:

	age	sex	bp	cholestorel	sugar	heartrate	angina	oldpeak	vessels	presence	...	chestpain	slope	ecg	thal
0	70.0	1.0	130.0	322.0	0.0	109.0	0.0	2.4	3.0	2	...	0	1	1	1
1	67.0	0.0	115.0	564.0	0.0	160.0	0.0	1.6	0.0	1	...	0	1	1	1
2	57.0	1.0	124.0	261.0	0.0	141.0	0.0	0.3	0.0	2	...	0	1	1	1
3	64.0	1.0	128.0	263.0	0.0	105.0	1.0	0.2	1.0	1	...	0	1	1	1
4	74.0	0.0	120.0	269.0	0.0	121.0	1.0	0.2	1.0	1	...	0	1	1	1

5 rows × 23 columns

## KNN Model from sklearn

**Q5** Get training data from the dataframe

1. Assign values of `presence` column to `y`, note you have to use `.values` method
2. Drop 'presence' column from data frame,
3. Assign `df` values to `x`

Split dataset into train and test data use `train_test_split`

1. Use `stratify = y` and `test_size = 0.2` and `random_state = 1`
2. Create a KNN model using sklearn library, Initialize `n_neighbors = 3`, (See the documenttaion for details)
3. Fit the model with the train data

**A5** Replace ??? with code in the code cell below

```
In [26]: ▶ import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier

# Assign values of ``presence`` column to y, note you have to use .v
y = df.ppresence.values
# Drop 'presence' column from data frame,
df.drop(columns=['presence'], inplace=True)
# Assign df values to x
x = df.values
# View shape of x and y
x.shape, y.shape

# Use stratify = y and test_size = 0.2 and random_state = 1

xtrain, xtest, ytrain, ytest = train_test_split(x,y, test_size = 0.2,
xtrain.shape, xtest.shape, ytrain.shape, ytest.shape
# Create a KNN model using sklearn library, k=3
#knn = KNeighborsClassifier(n_neighbors=3)

# Fit the model with the train data
#knn.fit(xtrain, ytrain)
```

Out[26]: ((216, 13), (54, 13), (216,), (54,))

#### Q6 Analysis

- Predict xtest and view first 25 predictions
- Compare prediction with real ytest 25 predictions
- Print the score with test data

The way we fit the dataset is not good

#### Normalization

- rescale only real value columns
- For each column normalize  $df[col]$  as  $(x - \text{mean}) / \text{standard\_deviation}$

A6 Replace ??? with code in the code cell below

```
In [6]: ▶ # Predict xtest and view first 25 predictions
print(knn.predict(xtest)[0:25])

# Compare prediction with real ytest 25 predictions
print(ytest[0:25])

# Print the score with test data
print(knn.score(xtest, ytest))

#rescale only real value columns
realcols = ['age', 'bp', 'cholestorel', 'heartrate', 'oldpeak', 'vessels']

# For each column normalize ``df[col] as (x - mean) / standard_deviat
for col in realcols:
    mean = df[col].mean()
```

```
std = df[col].std()
1.6511181111111111
[2 2 1 1 1 2 1 1 2 2 2 1 1 1 2 1 1 2 1 1 1 1]
[2 2 2 1 1 2 1 1 2 1 2 1 1 1 2 1 1 1 2 2 1 1]
0.7222222222222222
```

**Q7** Write the code to train new model using KNN classifier, k=3 (same as above)

**A7** Replace ??? with code in the code cell below

```
In [7]: # update x
x = df.values

# Train test Split
xtrain, xtest, ytrain, ytest = train_test_split(x,y, test_size = 0.2,

# Model Initialization
knn = KNeighborsClassifier(n_neighbors=3)

# Model fitting with training data
knn.fit(xtrain, ytrain)

# Now print score on test data
knn.score(xtest, ytest)
```

Out[7]: 0.8703703703703703

**Q8** Lets analyze the difference between two modeling strategies (data normalization) Compare score with and without data normalization process and explain

**A8**

Write your answer replacing this line

**Q9** Now we will write a function that will initialize, fit and return score on test data for given values of k and Plot result

1. Use values from 1 to 25(inclusive) and get score and plot as a bar graph

- Hint : For advance method you can use map (recall functional programming from last exercise) or you can use simple loops

2. Finally you can print the best value of k by getting the index

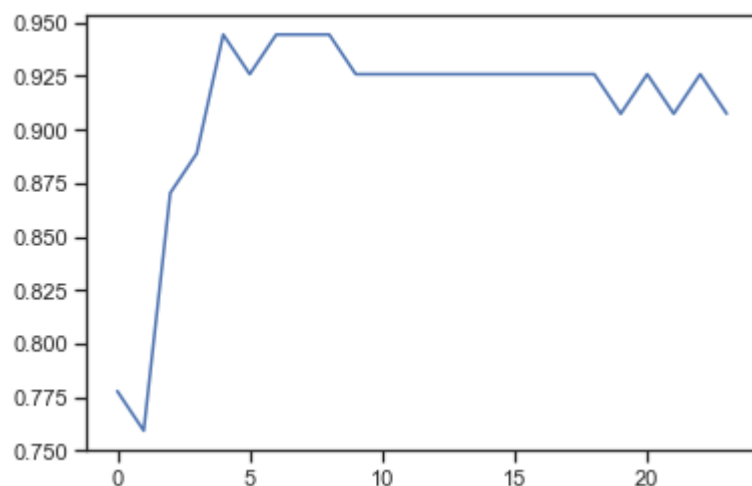
- N.B: Note index starts with 0 but values of k starts with 1 so actual value of k will be 1 more
- You can use `np.argmax()` function

3. Now define your best model as bestknn and print score

**A9** Write the code below (replace??)

```
In [8]: import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix
```

```
[0.7777777777777778, 0.7592592592592593, 0.8703703703703703, 0.888888  
8888888888, 0.9444444444444444, 0.9259259259259259, 0.9444444444444444  
4, 0.9444444444444444, 0.9444444444444444, 0.9259259259259259, 0.9259  
259259259259, 0.9259259259259259, 0.9259259259259259, 0.9259259259259  
259, 0.9259259259259259, 0.9259259259259259, 0.9259259259259259, 0.92  
59259259259259, 0.9259259259259259, 0.9074074074074074, 0.92592592592  
59259, 0.9074074074074074, 0.9259259259259259, 0.9074074074074074]  
BEST VALUE OF K 5  
0.9444444444444444  
[[29 1]  
 [ 2 22]]
```



- Please review following examples in documentation plot confusion plots
- [https://scikit-learn.org/stable/auto\\_examples/model\\_selection/plot\\_confusion\\_matrix.html](https://scikit-learn.org/stable/auto_examples/model_selection/plot_confusion_matrix.html)  
([https://scikit-learn.org/stable/auto\\_examples/model\\_selection/plot\\_confusion\\_matrix.html](https://scikit-learn.org/stable/auto_examples/model_selection/plot_confusion_matrix.html))

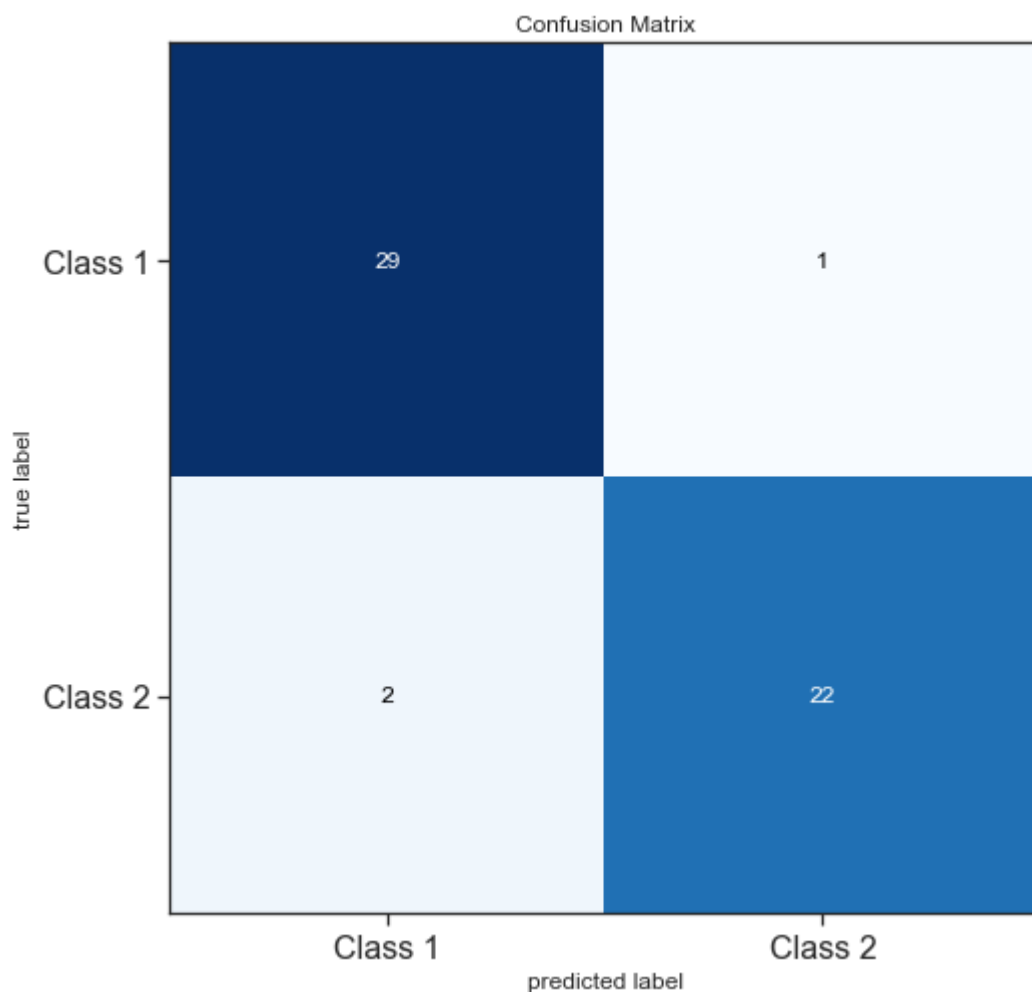


**A10** Replace ??? with code in the code cell below

```
In [9]:  from sklearn.metrics import plot_confusion_matrix
         from sklearn.metrics import confusion_matrix
         from mlxtend.plotting import plot_confusion_matrix
         import matplotlib.pyplot as plt

         cm = confusion_matrix(ytest, ypred)
         plt.figure()
         plot_confusion_matrix(cm, figsize=(12,8), hide_ticks=True, cmap=plt.cm
         plt.title("Confusion Matrix")
         plt.xticks(range(2), ["Class 1", "Class 2"], fontsize=16)
         plt.yticks(range(2), ["Class 1", "Class 2"], fontsize=16)
```

<Figure size 432x288 with 0 Axes>



**Q11:**

1. Calculate the test MSE
2. Get the score from the model using test data
3. Plot Precision-Recall Curve from the true & predicted test data

**A11** Replace ??? with code in the code cell below

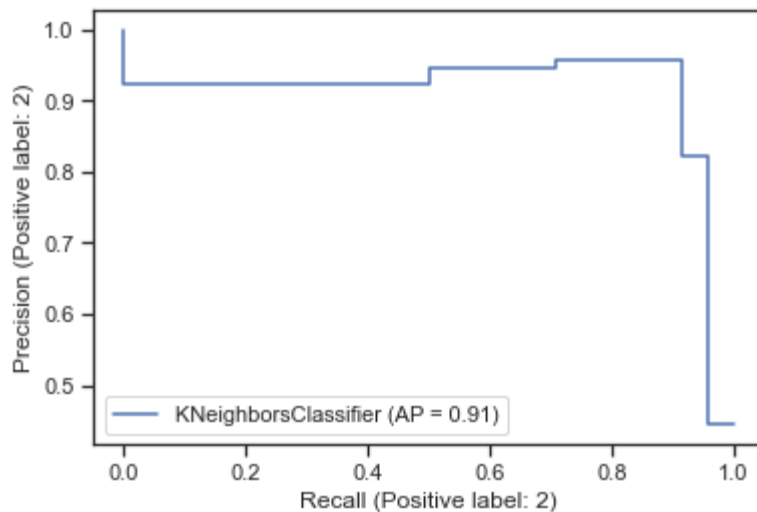
```
In [10]: ▶ from sklearn.metrics import mean_squared_error
from sklearn.metrics import PrecisionRecallDisplay
import matplotlib.pyplot as plt

mse = mean_squared_error(ytest, ypred) # Calculate the test MSE
print("Test mean squared error (MSE): {:.2f}".format(mse))

print(bestknn.score(xtest, ytest))

PrecisionRecallDisplay.from_estimator(bestknn, xtest, ytest)

Test mean squared error (MSE): 0.06
0.9444444444444444
```



## Further reading

- [KNN model creation \(https://towardsdatascience.com/building-a-k-nearest-neighbors-k-nn-model-with-scikit-learn-51209555453a\)](https://towardsdatascience.com/building-a-k-nearest-neighbors-k-nn-model-with-scikit-learn-51209555453a)
- [Example of KNN \(https://github.com/a-martyn/ISL-python/blob/master/Notebooks/ch4\\_classification\\_applied.ipynb\)](https://github.com/a-martyn/ISL-python/blob/master/Notebooks/ch4_classification_applied.ipynb)