# Machine Learning

## Course Work

MAHDSE212F-004     S.B.C. Sanjaya

Higher Diploma in Software Engineering 21.2F

Supervisor: Mr. Tharindu Adikari

National Institute of Business Management

2023

"The course work report is submitted in partial fulfilment of the requirement of the Machine Learning subject for Higher Diploma in Software Engineering of National Institute of Business Management."
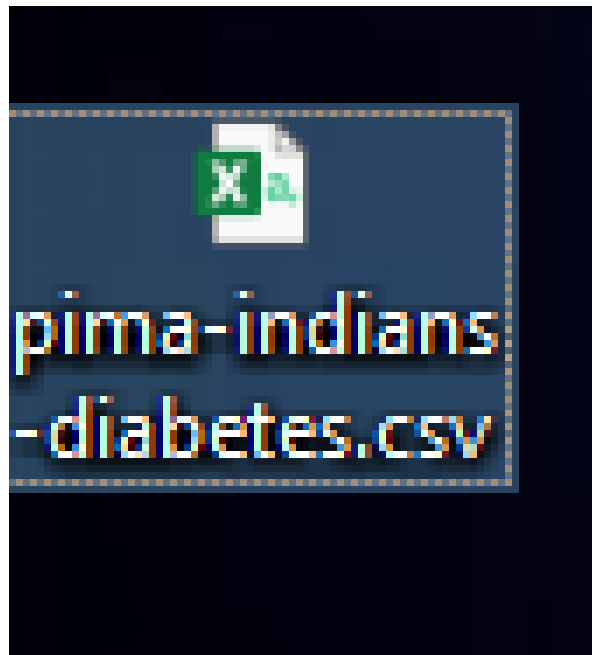
i

# Contents

# Introduction

I select a dataset with 1000> records data included one.

# Course Work Question Answers

- First I found a dataset with more than 1000 dataset records.

- This is the Prima-Indians_Diabetic.csv file with more than 1000 rows.

- First, I go to Google Colab and open a new notebook.
- Then I upload my downloaded Prima-Indian-Diabeties.csv file to the Google Colab.



Uploaded dataset.

- Then I open a new code and import pandas as following.

- Then I run this code as below.



- Then the compilation of above code completed successfully.

- This is the code I mention above.

```
#Importing
import pandas as pd
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn import metrics
```

- Then I read and test the file I uploaded. It can be coded by click using a new code + icon.



- Then I run this code as below.

- Then the compilation of above code completed successfully.

- This is the code I mention above.

```
col_names = ['pregnant', 'glucose', 'bp', 'skin', 'insulin', 'bmi', 'pedigre
e', 'age', 'label']
pima = pd.read_csv("pima-indians-
diabetes.csv", header=None, names=col_names)
```

- Then I enter a head and get the details of the table. It can be coded by click using a new code + icon.

- Then I run this code as below.



- This is the code I mention above.

```
pima.head()
```

- Then I split dataset in feature & finding target variable. It can be coded by click using a new code + icon.



- Then I run this code as below.



- This is the code I mention above.
  - `#split #dataset at feature & finding target variable`

```
feature_cols = ['pregnant', 'insulin', 'bmi', 'age', 'glucose', 'bp',
'pedigree']
X = pima[feature_cols] #features
y = pima.label #target variables
```

- Then I slit dataset into training set and test set. It can be coded by click using a new code +
  icon.

- Then I run this code as below.



- This is the code I mention above.

```
#slit dataset into training set and test set
X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.3
,random_state=1) #70% training
```

- Then I create decision tree classifier object, train decision tree classifier and predict the response for each dataset. It can be coded by click using a new code + icon.



- Then I run this code as below.



11

- This is the code I mention above.

```
#create dicision tree classifier object
clf = DecisionTreeClassifier()

#train decision tree classifier
clg = clf.fit(X_train, y_train)

#predict the response for each dataset
y_pred = clf.predict(X_test)
```

- Then I Model accuracy, how often is the classifier correct. It can be coded by click using a new code + icon.

- Then I run this code.



- This is the code I mention above.

```
#Model accuracy, how often is the classifier correct?
print("Accuracy: ", metrics.accuracy_score(y_test, y_pred))
```

- Finally, I got the accuracy value of this dataset as below.

```
Accuracy:   0.8590785907859079
```

- And I did some Machine Learning Algorithms as additional.

- Then I run codes as follows and did some Machine Learning Algorithms.





14

```python
from sklearn.cluster import KMeans
data = list(zip(x, y))
inertias = []
for i in range(1,11):
    kmeans = KMeans (n_clusters=i)
    kmeans.fit (data)
    inertias.append(kmeans.inertia_)
plt.plot (range(1, 11), inertias, marker='o')
plt.title('Elbow method')
plt.xlabel('Number of clusters')
plt.ylabel('Inertia')
plt.show()
```

```
<ipython-input-13-3541a499ede9>:6: ConvergenceWarning: Number of distinct clusters (9) found smaller than n_clusters (10). Possibly due to dup
  kmeans.fit (data)
```



```python
kmeans = KMeans(n_clusters = 2)
kmeans.fit(data)

plt.scatter(x,y, c = kmeans.labels_)
plt.show()
```

```
[16] knn.fit(data, classes)

     KNeighborsClassifier(n_neighbors=1)

[17] new_x = 8
     new_y = 21
     new_point = [(new_x, new_y)]

     prediction = knn.predict(new_point)

     plt.scatter(x + [new_x], y + [new_y], c=classes + [prediction[0]])
     plt.text(x=new_x-1.7, y=new_y-0.7, s=f"new point, class: {prediction[0]}")
     plt.show()
```



```
[18] KNeighborsClassifier(n_neighbors=5)

     knn.fit(data,classes)
     prediction = knn.predict(new_point)

     plt.scatter(x + [new_x], y + [new_y], c=classes + [prediction[0]])
     plt.text(x=new_x-1.7, y=new_y-0.7, s=f"new point, class: {prediction[0]}")
     plt.show()
```

```python
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import AgglomerativeClustering

x=[4,5,10,4,3,11,14,6,10,12]
y=[21,19,24,17,16,25,24,22,21,21]

data=list(zip(x,y))

heirarchical_cluster=AgglomerativeClustering(n_clusters=2,affinity='euclidea
n',linkage='ward')

labels=heirarchical_cluster.fit_predict(data)
plt.scatter(x,y, c=labels)

plt.show()
```

```python
import matplotlib.pyplot as plt

x = [4,5,10,4,3,11,14,6,10,12]
y = [21,19,24,17,16,25,24,22,24,21]

plt.scatter(x,y)
plt.show()
```

```python
from sklearn.cluster import KMeans
data = list(zip(x, y))
inertias = []
for i in range(1,11):
    kmeans = KMeans (n_clusters=i)
    kmeans.fit (data)
    inertias.append(kmeans.inertia_)
plt.plot (range(1, 11), inertias, marker='o')
plt.title('Elbow method')
plt.xlabel('Number of clusters')
plt.ylabel('Inertia')
plt.show()
```

```python
kmeans = KMeans(n_clusters = 2)
kmeans.fit(data)

plt.scatter(x,y, c = kmeans.labels_)
plt.show()
```

18

```python
import matplotlib.pyplot as plt

x = [4,5,10,4,3,11,14,8,10,12]
y = [21,19,24,17,16,25,24,22,21,21]
classes = [0,0,1,0,0,1,1,0,1,1]

plt.scatter(x,y, c=classes)
plt.show()
```

```python
from sklearn.neighbors import KNeighborsClassifier

data = list(zip(x,y))
knn = KNeighborsClassifier(n_neighbors=1)

knn.fit(data, classes)
```

```python
new_x = 8
new_y = 21
new_point = [(new_x, new_y)]

prediction = knn.predict(new_point)

plt.scatter(x + [new_x], y + [new_y], c=classes + [prediction[0]])
plt.text(x=new_x-1.7, y=new_y-0.7, s=f"new point, class: {prediction[0]}")
plt.show()
```

```python
KNeighborsClassifier(n_neighbors=5)

knn.fit(data,classes)
prediction = knn.predict(new_point)

plt.scatter(x + [new_x], y + [new_y], c=classes + [prediction[0]])
plt.text(x=new_x-1.7, y=new_y-0.7, s=f"new point, class: {prediction[0]}")
plt.show()
```

19

## Source Code

https://colab.research.google.com/drive/1OnBrYH5M064vyyyKxcBvUwLfo_keV25Z?usp=sharing

## Dataset File

https://drive.google.com/file/d/1GyTkel27AyonSxZHriUZyu05t9cbiJ8k/view?usp=share_link

# Conclusion

The conclusion is, I train some machine learning models and get output by finding the Accuracy value of the Prima-Indians-Diabestes.csv dataset file.

# References

Kaggle (2022). *Kaggle: Your Home for Data Science*. [online] Kaggle.com. Available at: https://www.kaggle.com/.

Google (2019). *Google Colaboratory*. [online] Google.com. Available at: https://colab.research.google.com/.