

Name: Allan Rodrigues
Class: TE IT A
Roll no: 59
Pid:191104

St. Francis Institute of Technology, Mumbai-400 103
Department of Information Technology

A.Y. 2021-2022
Class: TE-ITA/B, Semester: VI

Subject: **Data Science Lab**

Experiment – 3: To implement Data Modelling.

1. **Aim:** To implement Data Modelling
2. **Objectives:** After study of this experiment, the student will be able to
 - Understand how data to be re-scaled
 - Understand how data partitioning in done using sklearn
3. **Outcomes:** After study of this experiment, the student will be able to
 - Understand data rescaling and data partitioning using sklearn
4. **Prerequisite:** Fundamentals of Python Programming and Database Management System.
5. **Requirements:** Python Installation, Personal Computer, Windows operating system, Internet Connection, Microsoft Word.
6. **Pre-Experiment Exercise:**

Brief Theory:

- Concept of sklearn package for machine learning.

Scikit-learn provides a range of supervised and unsupervised learning algorithms via a consistent interface in Python. Extensions or modules for SciPy are conventionally named SciKits. As such, the module provides learning algorithms and is named scikit-learn. The library is focused on modeling data. It is not focused on loading, manipulating and summarizing data. For these features, refer to NumPy and Pandas.

Laboratory Exercise

- A. **Procedure:** (the sheet for commands in attached with the file)
- B. Paste Screenshots of above commands.

Exp 3.ipynb

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

RAM Disk

Editing

```
[1] import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd
import numpy as np

sns.get_dataset_names()
```

```
['anagrams',
 'anscombe',
 'attention',
 'brain_networks',
 'car_crashes',
 'diamonds',
 'dolls',
 'exercise',
 'flights',
 'fmri',
 'gammas',
 'geyser',
 'iris',
 'mpg',
 'penguins',
 'planets',
 'taxis',
 'tips',
```

Exp 3.ipynb

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

RAM Disk

Editing

```
[2] df= sns.load_dataset("titanic")
df
```

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male	deck	embark_town	alive	alone
0	0	3	male	22.0	1	0	7.2500	S	Third	man	True	NaN	Southampton	no	False
1	1	1	female	38.0	1	0	71.2833	C	First	woman	False	C	Cherbourg	yes	False
2	1	3	female	26.0	0	0	7.9250	S	Third	woman	False	NaN	Southampton	yes	True
3	1	1	female	35.0	1	0	53.1000	S	First	woman	False	C	Southampton	yes	False
4	0	3	male	35.0	0	0	8.0500	S	Third	man	True	NaN	Southampton	no	True
...
886	0	2	male	27.0	0	0	13.0000	S	Second	man	True	NaN	Southampton	no	True
887	1	1	female	19.0	0	0	30.0000	S	First	woman	False	B	Southampton	yes	True
888	0	3	female	NaN	1	2	23.4500	S	Third	woman	False	NaN	Southampton	no	False
889	1	1	male	26.0	0	0	30.0000	C	First	man	True	C	Cherbourg	yes	True
890	0	3	male	32.0	0	0	7.7500	Q	Third	man	True	NaN	Queenstown	no	True

Exp 3.ipynb

File Edit View Insert Runtime Tools Help All changes saved

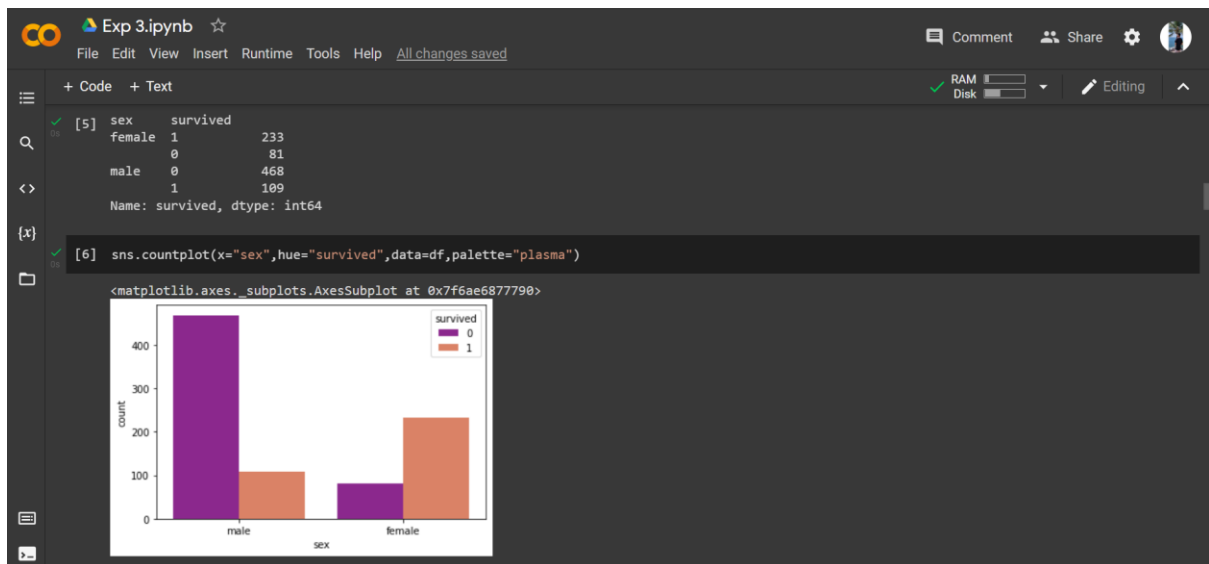
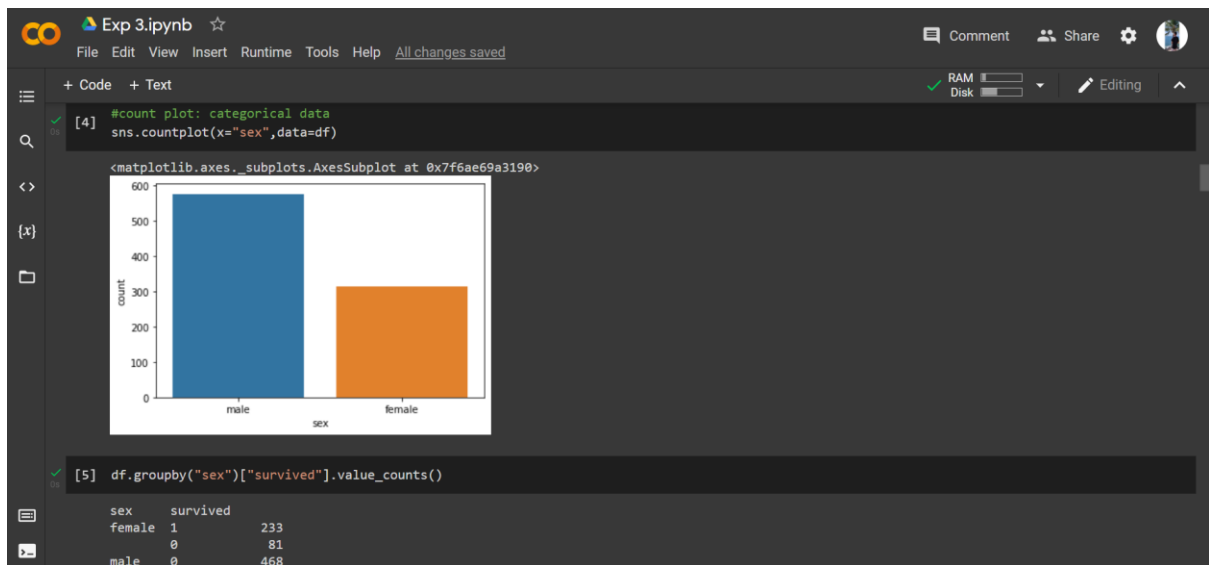
+ Code + Text

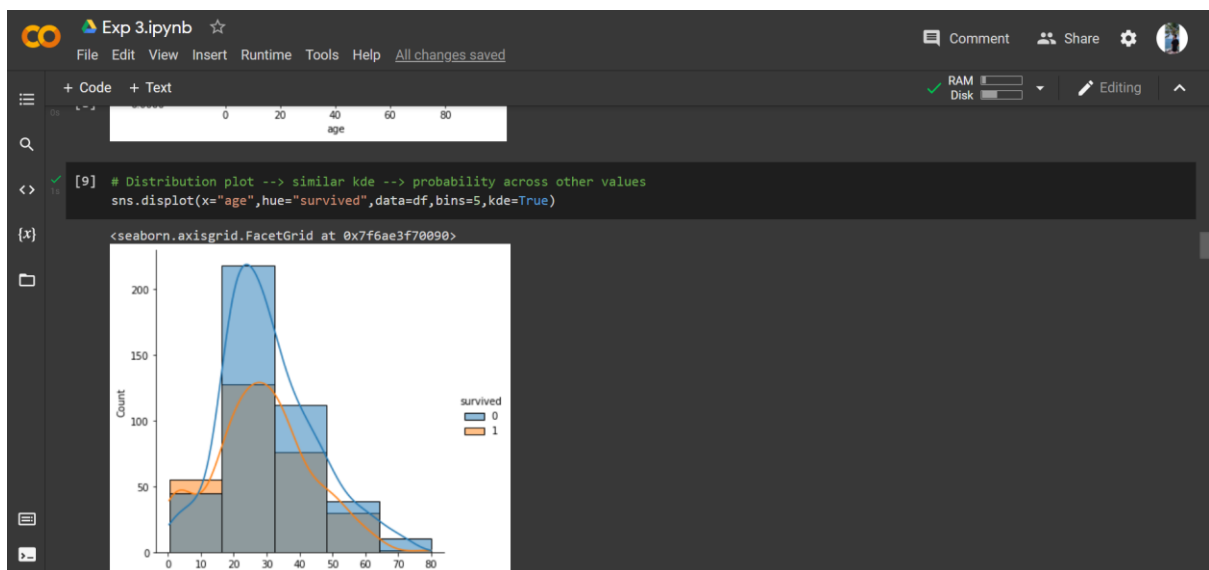
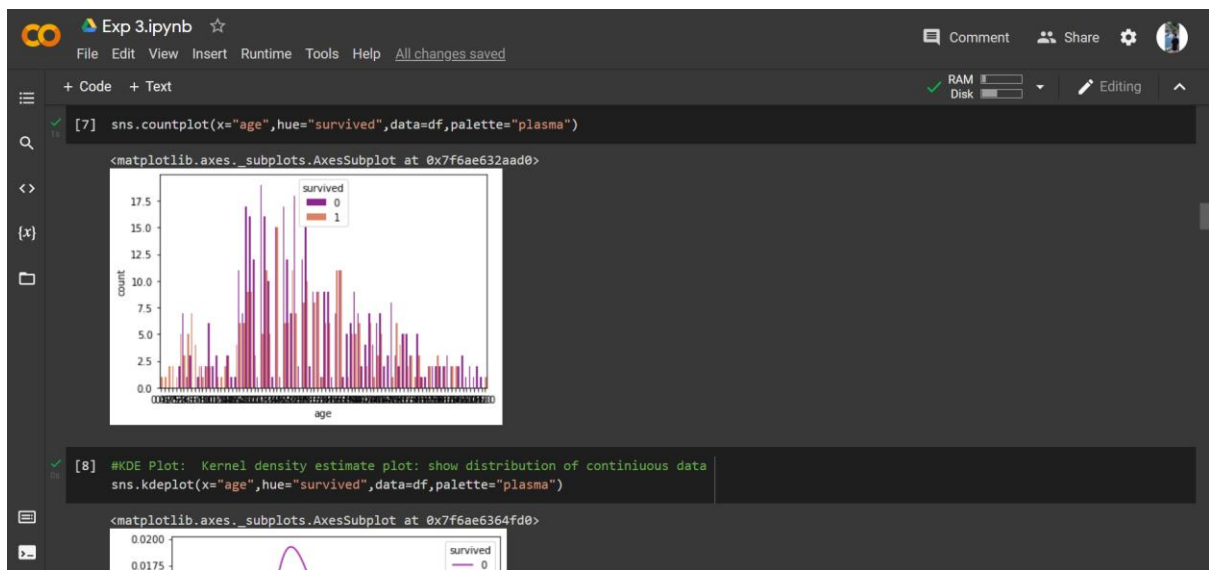
RAM Disk

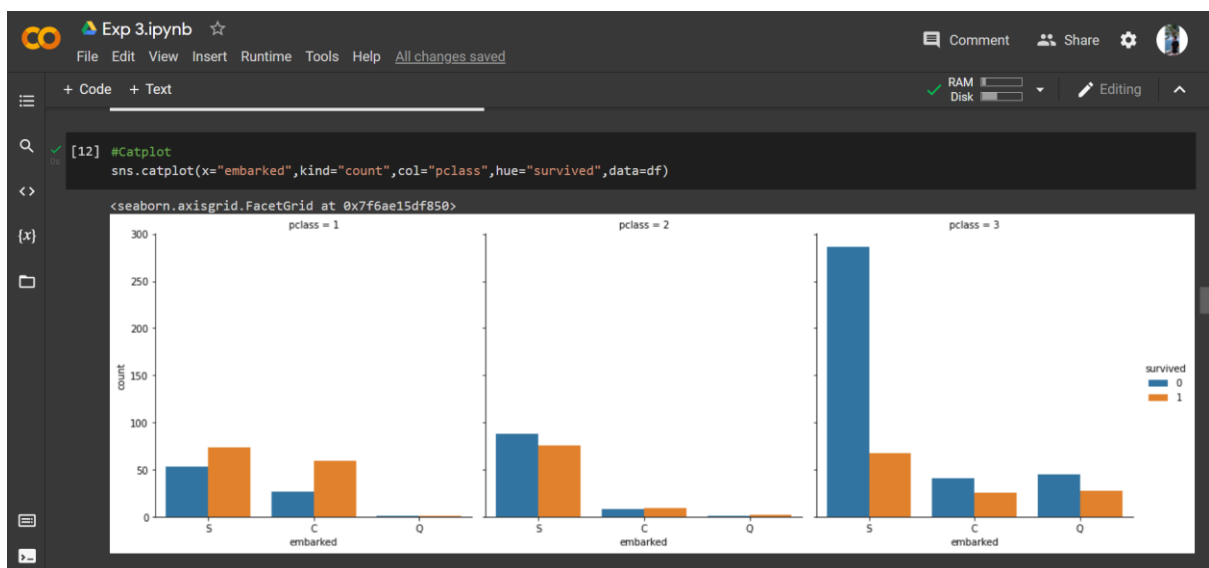
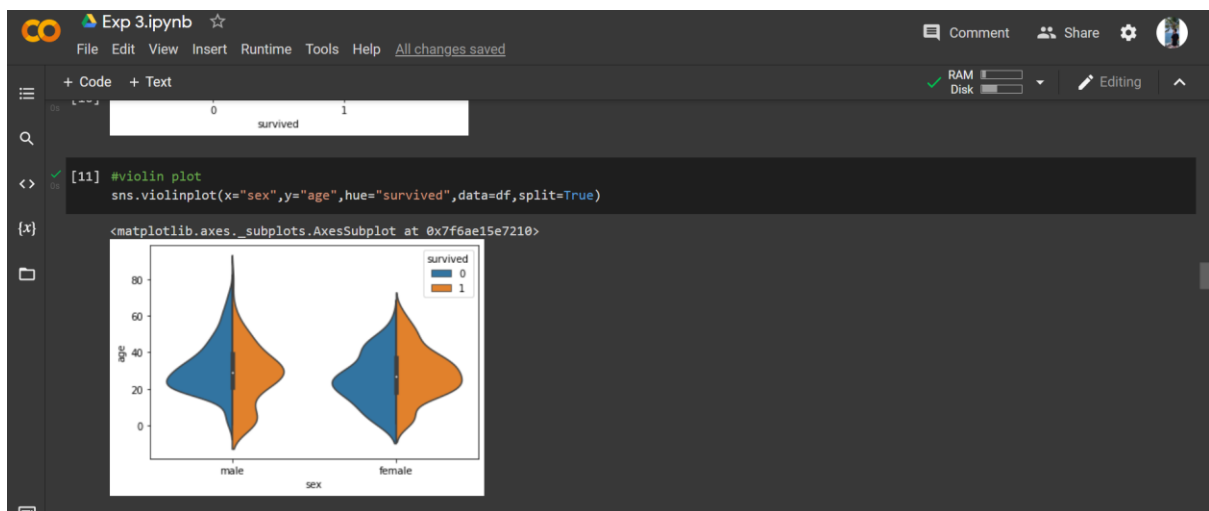
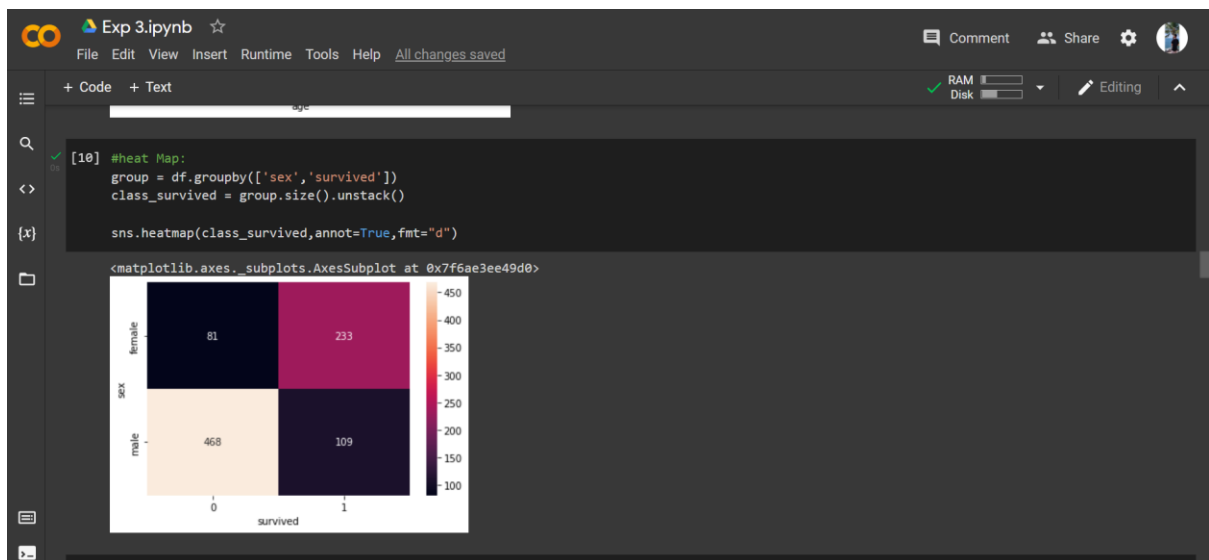
Editing

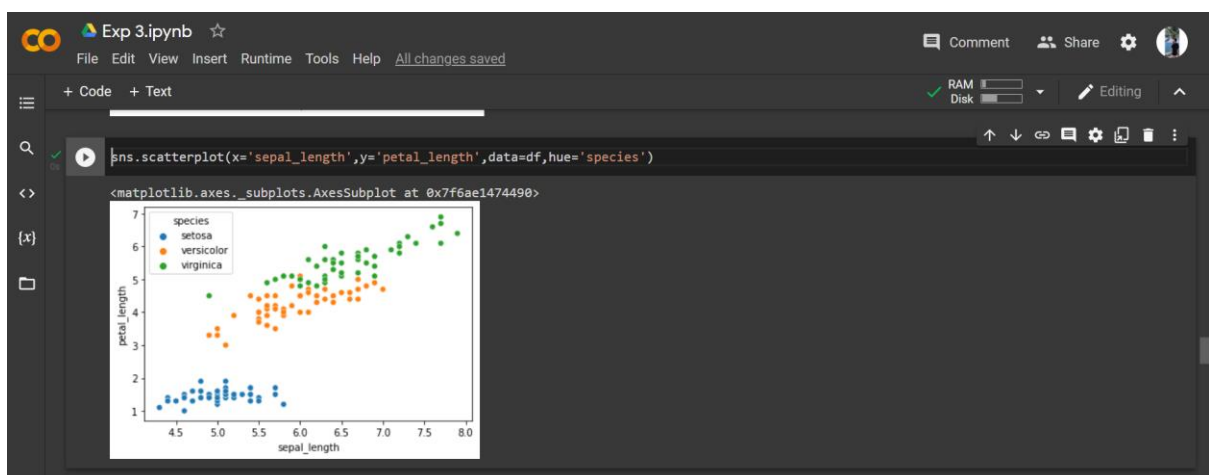
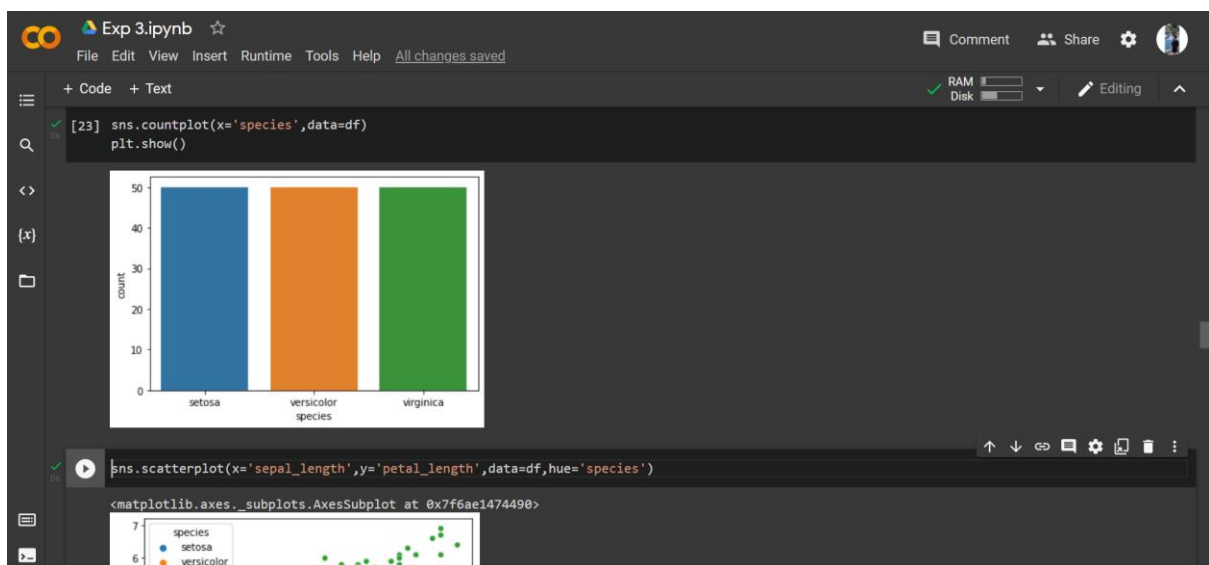
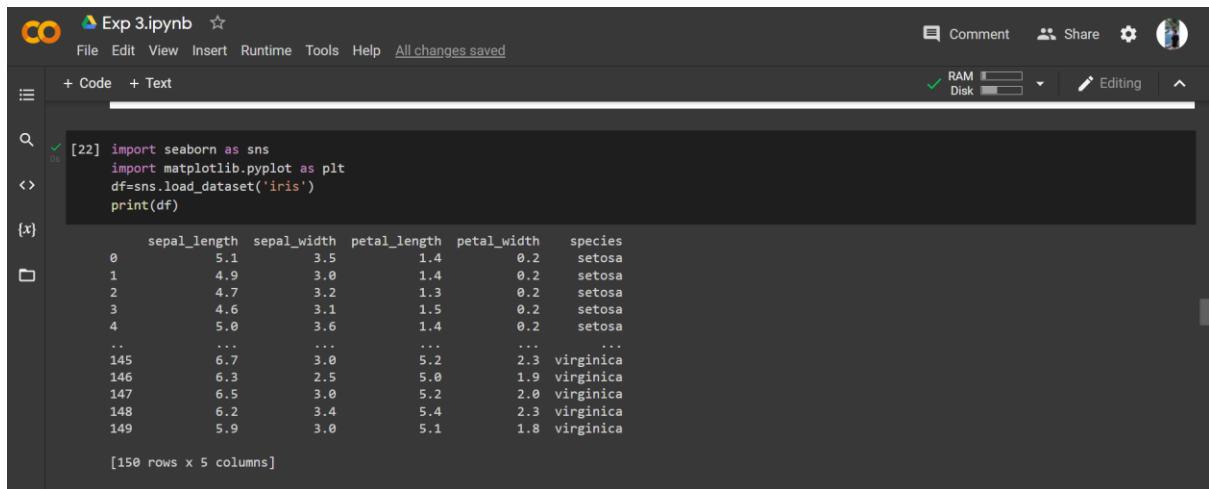
```
[3] df.info()
```

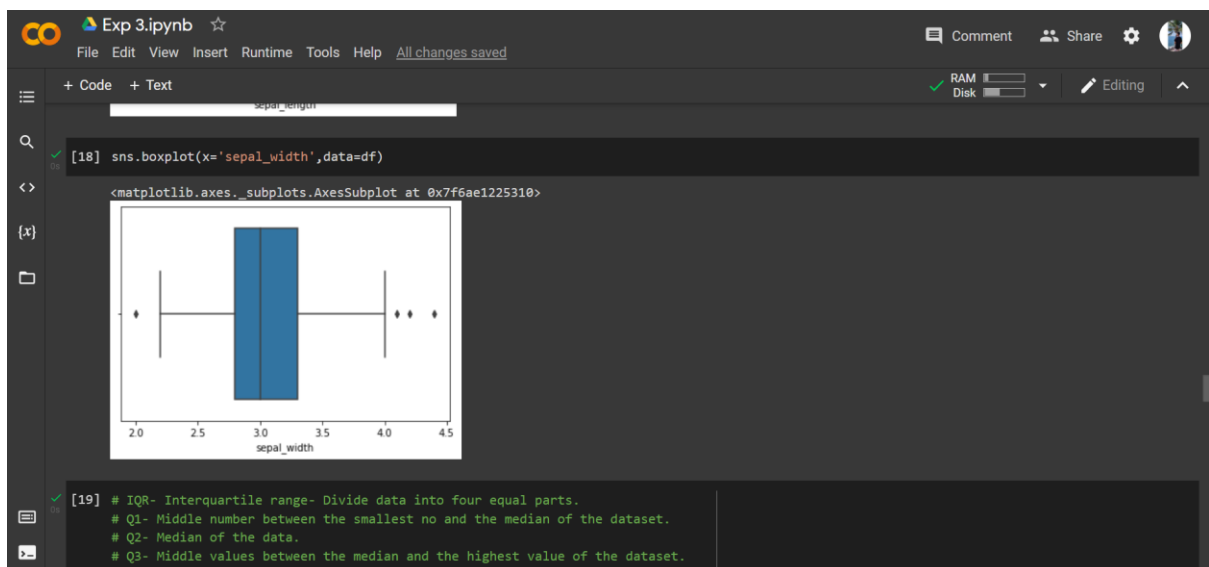
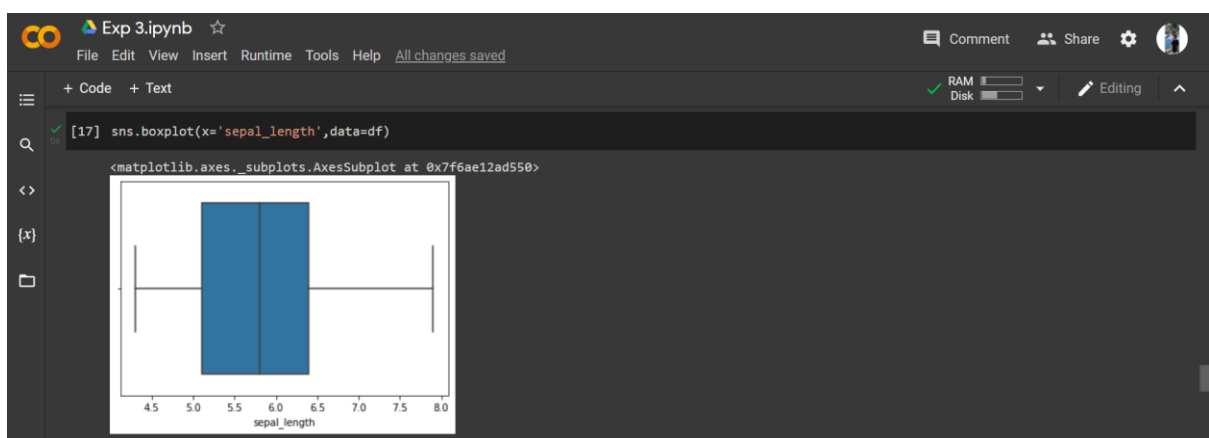
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 15 columns):
#   Column      Non-Null Count  Dtype
---  -
0   survived    891 non-null    int64
1   pclass      891 non-null    int64
2   sex         891 non-null    object
3   age         714 non-null    float64
4   sibsp       891 non-null    int64
5   parch       891 non-null    int64
6   fare        891 non-null    float64
7   embarked    889 non-null    object
8   class       891 non-null    category
9   who         891 non-null    object
10  adult_male  891 non-null    bool
11  deck        203 non-null    category
12  embark_town 889 non-null    object
13  alive       891 non-null    object
14  alone       891 non-null    bool
dtypes: bool(2), category(2), float64(2), int64(4), object(5)
memory usage: 80.7+ KB
```

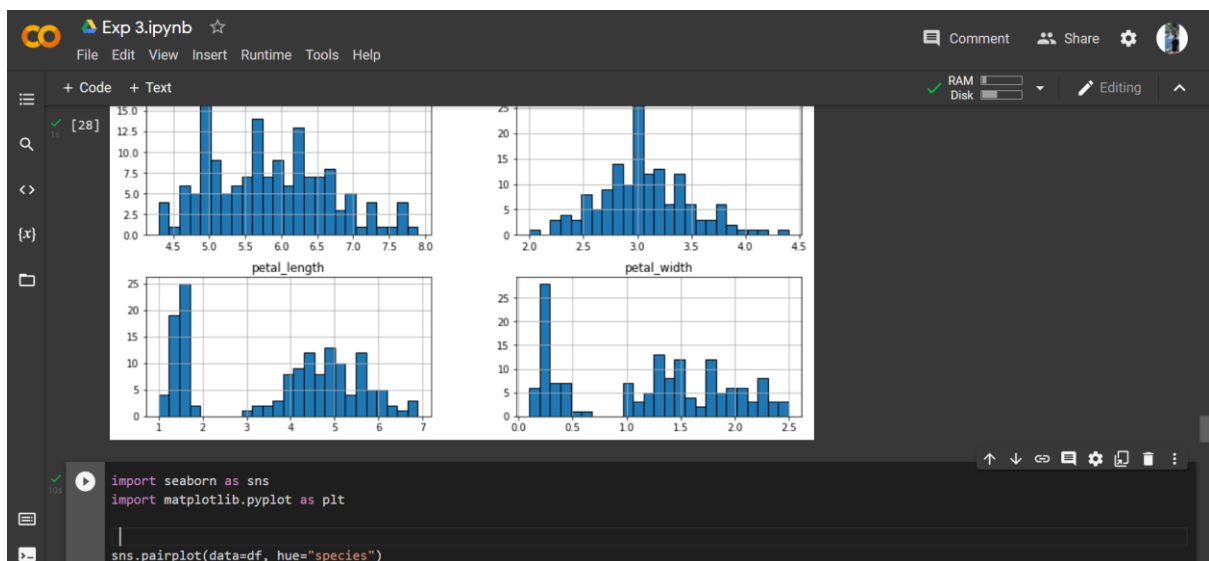
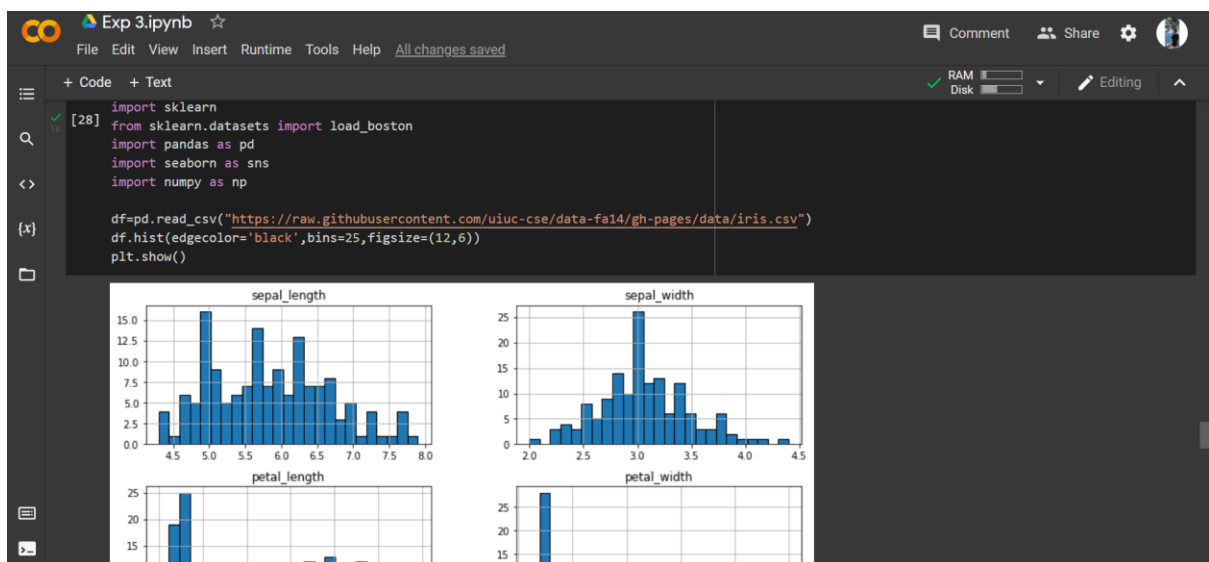
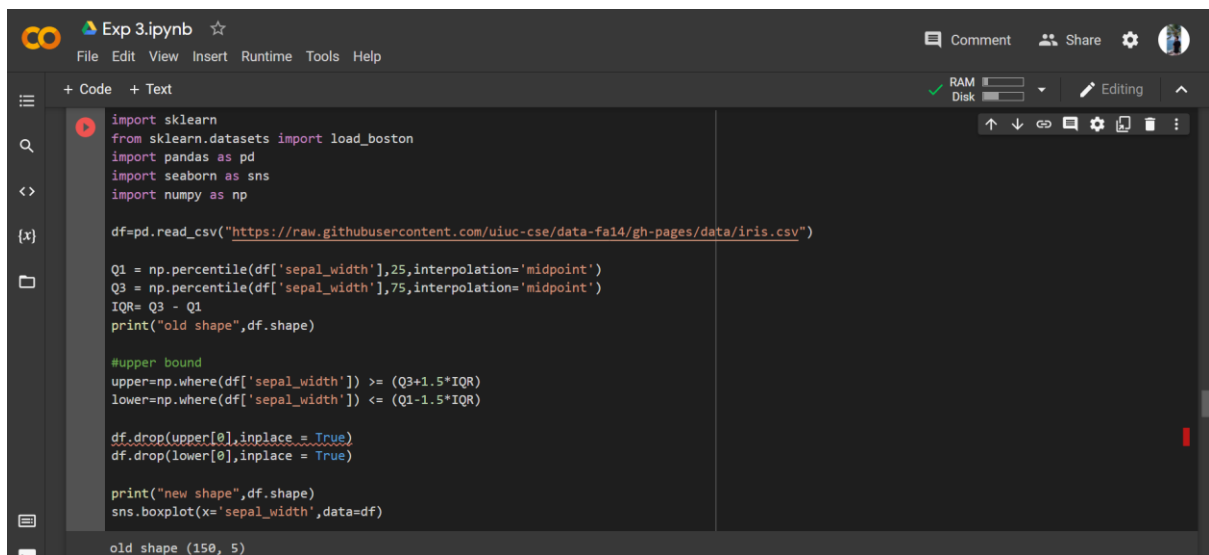


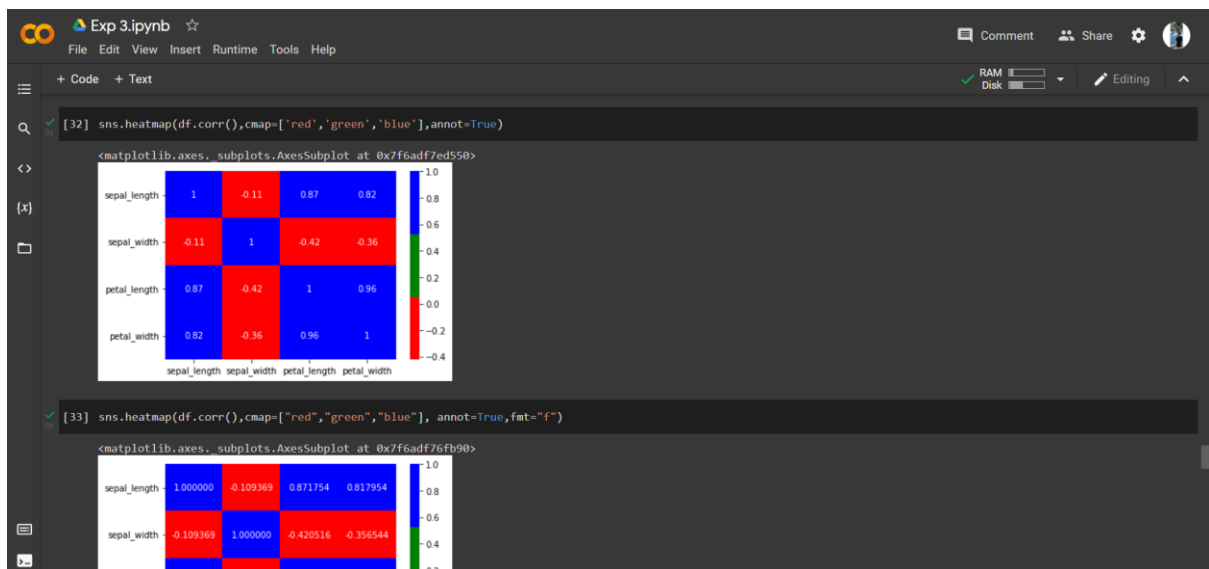
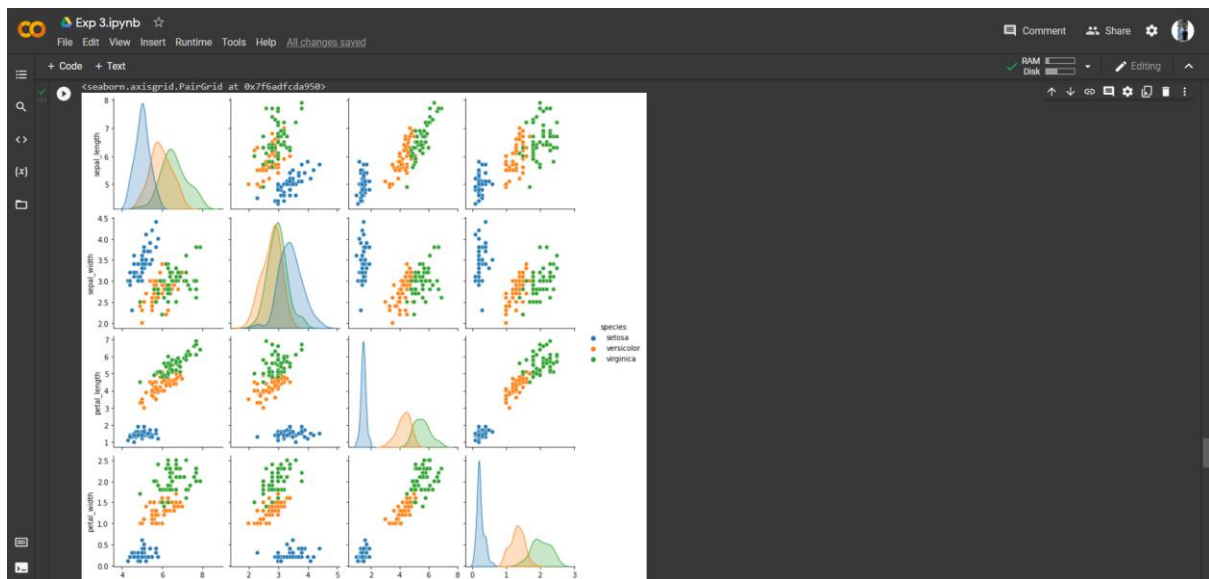


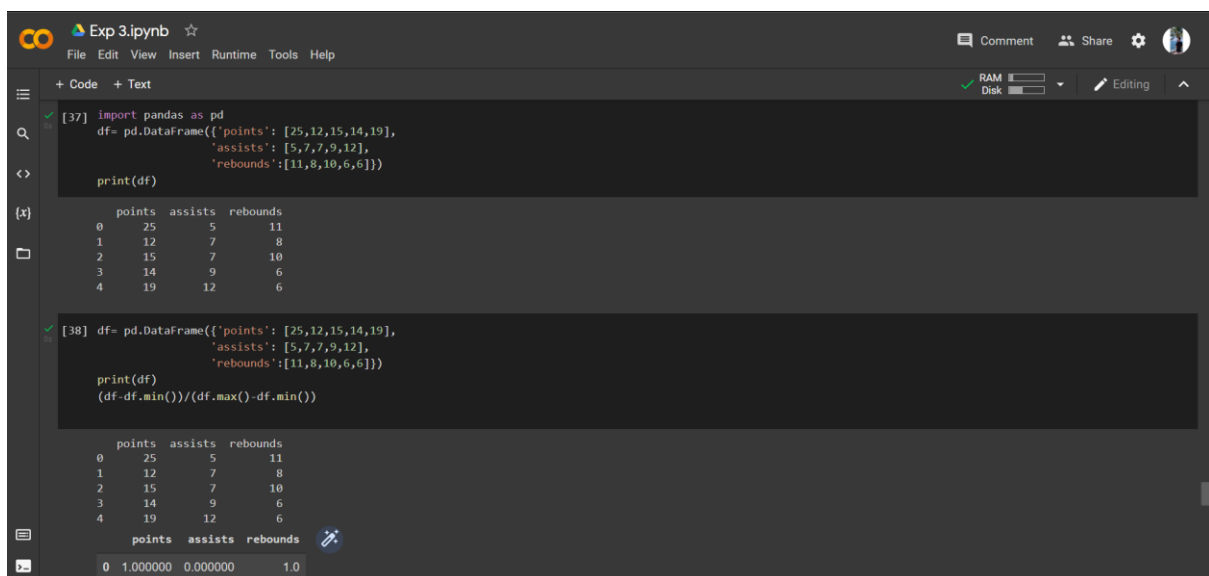
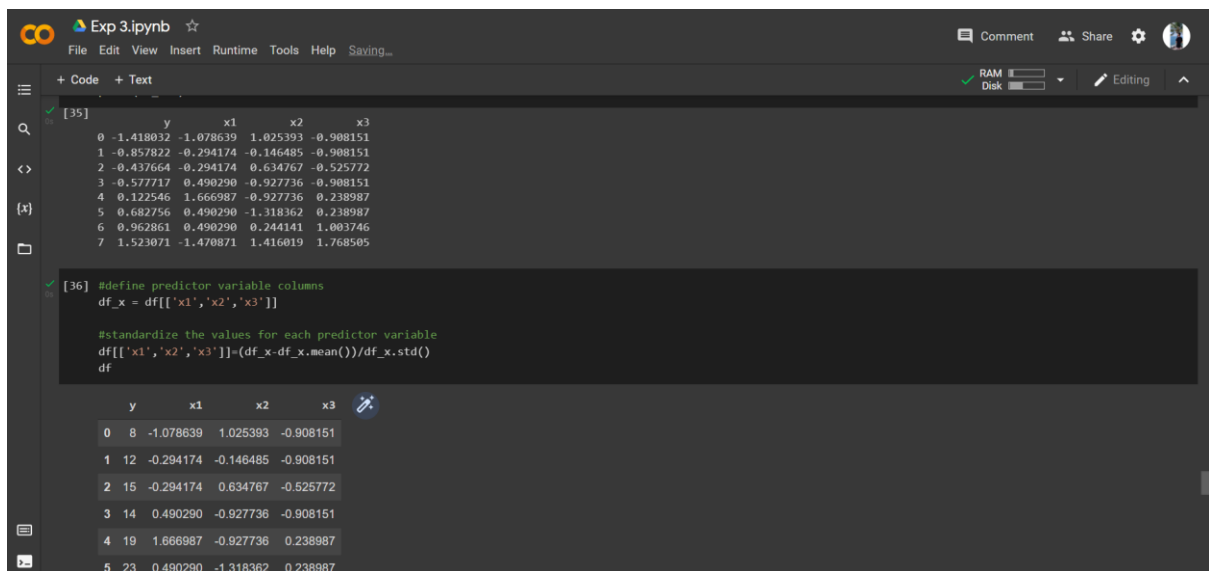
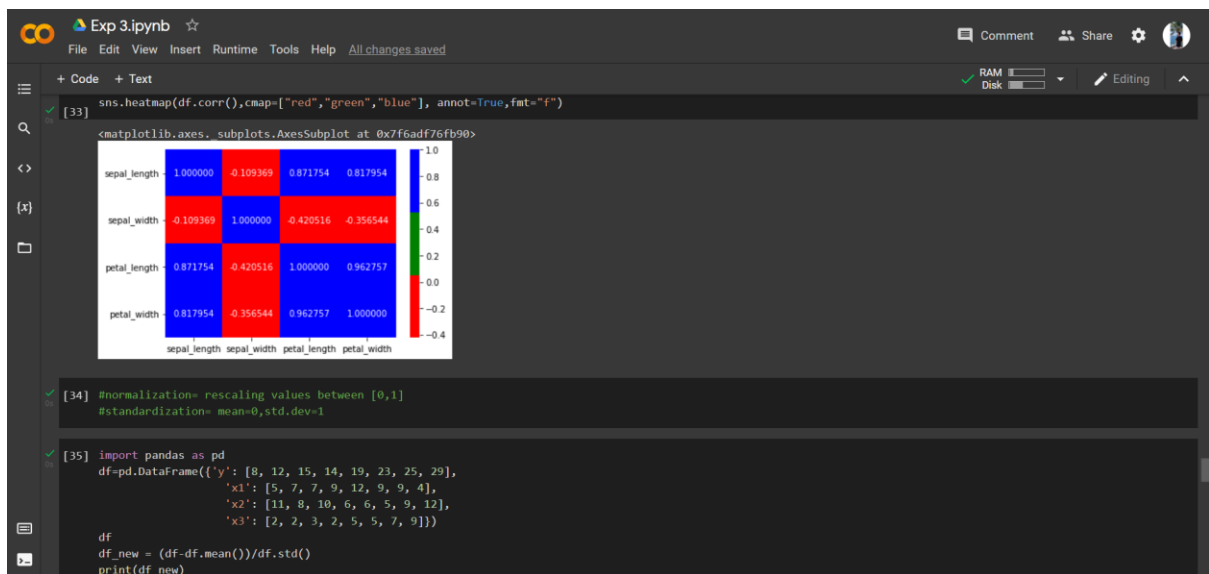












Exp 3.ipynb

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

RAM Disk

Editing

```
[38] df= pd.DataFrame({'points': [25,12,15,14,19],
                        'assists': [5,7,7,9,12],
                        'rebounds':[11,8,10,6,6]})

print(df)
(df-df.min())/(df.max()-df.min())
```

	points	assists	rebounds
0	25	5	11
1	12	7	8
2	15	7	10
3	14	9	6
4	19	12	6

	points	assists	rebounds
0	1.000000	0.000000	1.0
1	0.000000	0.285714	0.4
2	0.230769	0.285714	0.8
3	0.153846	0.571429	0.0
4	0.538462	1.000000	0.0

Exp 3.ipynb

File Edit View Insert Runtime Tools Help

+ Code + Text

RAM Disk

Editing

```
[39] import pandas as pd
df=pd.read_csv("https://raw.githubusercontent.com/uluc-cse/data-fal4/ph-pages/data/iris.csv")
pd.crosstab(index=df['species'],columns=df['sepal_length'],margins=True)
```

sepal_length	4.3	4.4	4.5	4.6	4.7	4.8	4.9	5.0	5.1	5.2	5.3	5.4	5.5	5.6	5.7	5.8	5.9	6.0	6.1	6.2	6.3	6.4	6.5	6.6	6.7	6.8	6.9	7.0	7.1	7.2	7.3	7.
species																																
setosa	1	3	1	4	2	5	4	8	8	3	1	5	2	0	2	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
versicolor	0	0	0	0	0	0	1	2	1	1	0	1	5	5	5	3	2	4	4	2	3	2	1	2	3	1	1	1	0	0	0	
virginica	0	0	0	0	0	0	1	0	0	0	0	0	0	1	1	3	1	2	2	2	6	5	4	0	5	2	3	0	1	3	1	
All	1	3	1	4	2	5	6	10	9	4	1	6	7	6	8	7	3	6	6	4	9	7	5	2	8	3	4	1	1	3	1	

```
[40] import pandas as pd
from sklearn.datasets import load_iris
iris_data=load_iris()
df=pd.DataFrame(iris_data.data,columns=iris_data.feature_names)
print(df)
```

Exp 3.ipynb

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

RAM Disk

Editing

```
[40] import pandas as pd
from sklearn.datasets import load_iris
iris_data=load_iris()
df=pd.DataFrame(iris_data.data,columns=iris_data.feature_names)
print(df)
```

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2
...
145	6.7	3.0	5.2	2.3
146	6.3	2.5	5.0	1.9
147	6.5	3.0	5.2	2.0
148	6.2	3.4	5.4	2.3
149	5.9	3.0	5.1	1.8

[150 rows x 4 columns]

```
[41] training_data=df.sample(frac=0.8,random_state=25)
testing_data=df.drop(training_data.index)
print(f"no of training examples: {training_data.shape[0]}")
print(f"no of testing examples: {testing_data.shape[0]}")

no of training examples: 120
no of testing examples: 30

[42] from sklearn.model_selection import train_test_split
training_data,testing_data = train_test_split(df,test_size=0.2,random_state=25)
print(f"no of training examples: {training_data.shape[0]}")
print(f"no of testing examples: {testing_data.shape[0]}")

no of training examples: 120
no of testing examples: 30
```

```
[43] from sklearn import datasets
from sklearn.preprocessing import StandardScaler, MinMaxScaler
from sklearn.model_selection import train_test_split

# loading the Iris data
iris = datasets.load_iris()
print(iris)
X = iris.data # array for the features
y = iris.target # array for the target
feature_names = iris.feature_names # feature names
target_names = iris.target_names # target names
# splitting the data into training and testing data sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.333,random_state=2020)
X_train

[[5.4, 3.9, 1.7, 0.4],
 [6.6, 3. , 4.4, 1.4],
 [5.9, 3. , 5.1, 1.8],
 [5.8, 4. , 1.2, 0.2],
 [6.7, 3. , 5.2, 2.3],
 [5.6, 2.7, 4.2, 1.3],
 [6.3, 2.5, 4.9, 1.5],
 [7.3, 2.9, 6.3, 1.8],
 [5.5, 2.4, 3.8, 1.1],
 [5.5, 4.2, 1.4, 0.2],
 [6.8, 3. , 5.5, 2.1],
 [5.6, 3. , 4.5, 1.5],
 [4.4, 3.2, 1.3, 0.2],
 [5.8, 3.2, 4.6, 1.4],
 [6.4, 3.2, 4.5, 1.4]]
```

8. Post-Experiments Exercise

A. Extended Theory: (Soft Copy)

- Use iris dataset and perform rescaling using sklearn package using normalization

```
[48] - Gates, G.W. (1972) "The Reduced Nearest Neighbor Rule". IEEE Transactions
on Information Theory, May 1972, 431-433.
- See also: 1988 MLC Proceedings, 54-64. Cheeseman et al's AUTOCLASS II
conceptual clustering system finds 3 classes in the data.
- Many, many more ...

sepal length (cm) sepal width (cm) petal length (cm) petal width (cm) target
0 5.1 3.5 1.4 0.2 0
1 4.9 3.0 1.4 0.2 0
2 4.7 3.2 1.3 0.2 0
3 4.6 3.1 1.5 0.2 0
4 5.0 3.6 1.4 0.2 0
...
145 6.7 3.0 5.2 2.3 2
146 6.3 2.5 5.0 1.9 2
147 6.5 3.0 5.2 2.0 2
148 6.2 3.4 5.4 2.3 2
149 5.9 3.0 5.1 1.8 2
```

150 rows x 5 columns

The screenshot shows a JupyterLab environment with a dark theme. The top bar includes the 'Exp 3.ipynb' title and a menu with 'File', 'Edit', 'View', 'Insert', 'Runtime', 'Tools', 'Help', and 'All changes saved'. The left sidebar has icons for a file explorer, search, and other tools. The main area displays a Python script for normalizing data using sklearn's MinMaxScaler. Below the script, a table shows the normalized data for the first five rows of the Iris dataset. The table has columns for sepal length, sepal width, petal length, petal width, and target. The bottom of the image shows a snippet of the script's output, which includes a description of the dataset and the first few rows of the data array.

```
# NORMALIZATION
from sklearn.preprocessing import MinMaxScaler
df_n = df.copy()
min_max_scaler = MinMaxScaler()
df_n.iloc[:, [0, 1, 2, 3]] = min_max_scaler.fit_transform(df_n.iloc[:, [0, 1, 2, 3]])
df_n.head()
```

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
0	0.222222	0.625000	0.067797	0.041667	0
1	0.186867	0.416667	0.067797	0.041667	0
2	0.111111	0.500000	0.050847	0.041667	0
3	0.083333	0.458333	0.084746	0.041667	0
4	0.194444	0.886867	0.067797	0.041667	0

```
[46] from sklearn.model_selection import train_test_split
from sklearn.datasets import load_iris
iris = load_iris()
x, y = iris.data, iris.target
iris

{'DESCR': '.. _iris_dataset:\n\nIris plants dataset\n-----\n\nData Set Characteristics: **\n\n :Number of Instances: 150 (50 in each of three cl
'data': array([[5.1, 3.5, 1.4, 0.2],
               [4.9, 3. , 1.4, 0.2],
               [4.7, 3.2, 1.3, 0.2],
```

- Partition the iris dataset such that 80% data to be taken for training purpose

```
[46]: from sklearn.model_selection import train_test_split
      from sklearn.datasets import load_iris
      iris = load_iris()
      x, y = iris.data, iris.target
      iris

{'DESCR': '.._iris_dataset:\n\nIris plants dataset\n-----\n\nData Set Characteristics: **\n\n'
'data': array([[5.1, 3.5, 1.4, 0.2],
               [4.9, 3. , 1.4, 0.2],
               [4.7, 3.2, 1.3, 0.2],
```

The screenshot shows a Jupyter Notebook interface. The top bar indicates the notebook is named "Exp 3.ipynb" and is running. The menu bar includes File, Edit, View, Insert, Runtime, Tools, Help, and a link to "All changes saved". The left sidebar shows the notebook's structure with a "Code" tab selected. The main area displays a code cell with the following content:

```
[47]: # Splitting dataset into 80/20 ratio
      x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=123)
      x_train
```

The output of the code cell is a list of 20 data points, each represented as a list of five values:

```
array([[7.4, 2.8, 6.1, 1.9],
       [6. , 2.2, 5. , 1.5],
       [4.7, 3.2, 1.6, 0.2],
       [5.1, 3.5, 1.4, 0.2],
       [6. , 2.2, 4. , 1. ],
       [5. , 2.3, 3.3, 1. ],
       [7.9, 3.8, 6.4, 2. ],
       [5.4, 3.9, 1.7, 0.4],
       [5.4, 3.9, 1.3, 0.4],
       [5.8, 2.7, 3.9, 1.2],
       [5. , 2. , 3.5, 1. ],
       [5. , 3.2, 1.2, 0.2],
       [6.8, 3.2, 5.9, 2.3],
       [6.7, 3. , 5.2, 2.3],
       [5.8, 2.7, 5.1, 1.9],
       [5.8, 2.8, 5.1, 2.4],
       [6.3, 3.4, 5.6, 2.4],
       [5.5, 2.3, 4. , 1.3],
       [5.1, 3.8, 1.5, 0.3],
       [4.4, 3. , 1.3, 0.2],
       [6.5, 3.2, 5.1, 2. ],
       [5.1, 3.3, 1.7, 0.5],
       [4.9, 3.1, 1.5, 0.1],
       [6.7, 3.1, 4.7, 1.5],
```

B. Questions:

- Mention three differences between normalization and standardization.
- Describe train_test_split function

Q.8

B

1) Mention three difference between normalization & standardization

Normalization	Standardization
1) Minimum & maximum value of features are used for scaling	1) Mean & standard deviation is used for scaling
2) It is used when features are of different scales	2) It is used when we want to ensure zero mean & unit standard deviation
3) Scales values between $[0, 1]$ or $[-1, 1]$	3) It is not bounded to a certain range

2) Describe train-test-split function.

~~test~~ train-test-split is a function in sklearn model selection for splitting data arrays into two subsets for training data & for testing data. With this function, you don't need to manually divide the dataset manually. By default, sklearn train-test-split will make random partitions for the two subsets.

C. Conclusion:

Write the significance of the topic studied in the experiment.

C) Conclusion.

In this experiment we used the sklearn library for various plots like countplot, distribution plot and heatmap. We also learnt about normalization & standardization. We learnt to split the dataset into training and testing using train-test-split function from sklearn.

D. References:

1. <https://www.geeksforgeeks.org/exploratory-data-analysis-on-iris-dataset/>
2. <https://www.statology.org/normalize-columns-pandas-dataframe/Normalization>
3. <https://www.datascienceguide.org/python-code-snippets.html>
4. <https://machinelearningmastery.com/standardscaler-and-minmaxscaler-transforms-in-python/>
