

230701521

JABARAJ E

“A”

1.MATPLOT SEABORN

```
In [16]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: data=pd.read_csv('Iris.csv')
data
```

```
Out[2]:
```

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

150 rows × 5 columns

```
In [3]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   sepal.length    150 non-null    float64
1   sepal.width     150 non-null    float64
2   petal.length    150 non-null    float64
3   petal.width     150 non-null    float64
4   variety         150 non-null    object
dtypes: float64(4), object(1)
memory usage: 6.0+ KB
```

```
In [4]: data.describe()
```

```
Out[4]:
```

	sepal.length	sepal.width	petal.length	petal.width
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.057333	3.758000	1.199333
std	0.828066	0.435866	1.765298	0.762238
min	4.300000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	0.300000
50%	5.800000	3.000000	4.350000	1.300000
75%	6.400000	3.300000	5.100000	1.800000
max	7.900000	4.400000	6.900000	2.500000

```
In [5]: data.value_counts('variety')
```

```
Out[5]:
```

variety	count
Setosa	50
Versicolor	50
Virginica	50

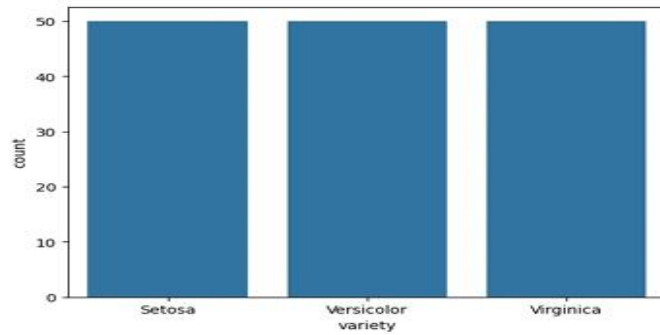
Name: count, dtype: int64

```
In [6]: sns.countplot(x='variety',data=data,)
plt.show()
```

230701521

JABARAJ E

"A"



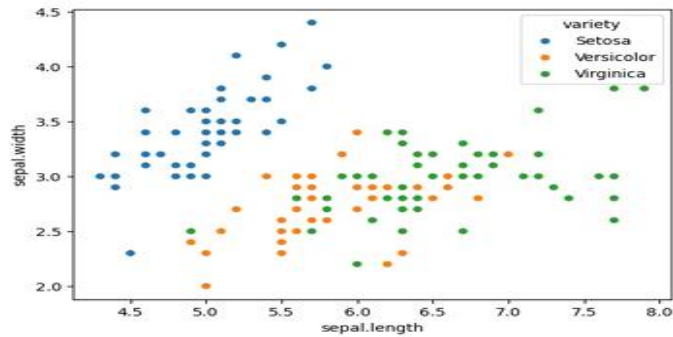
```
In [7]: dummies=pd.get_dummies(data.variety)
FinalDataset=pd.concat([pd.get_dummies(data.variety),data.iloc[:,[0,1,2,3]]],axis=1)
FinalDataset.head()
```

```
Out[7]:
```

	Setosa	Versicolor	Virginica	sepal.length	sepal.width	petal.length	petal.width
0	True	False	False	5.1	3.5	1.4	0.2
1	True	False	False	4.9	3.0	1.4	0.2
2	True	False	False	4.7	3.2	1.3	0.2
3	True	False	False	4.6	3.1	1.5	0.2
4	True	False	False	5.0	3.6	1.4	0.2

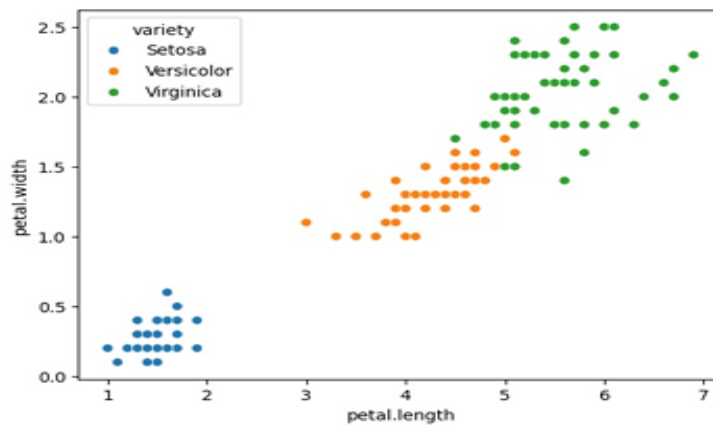
```
In [8]: sns.scatterplot(x='sepal.length',y='sepal.width',hue='variety',data=data,)
```

```
Out[8]: <Axes: xlabel='sepal.length', ylabel='sepal.width'>
```



```
In [9]: sns.scatterplot(x='petal.length',y='petal.width',hue='variety',data=data,)
```

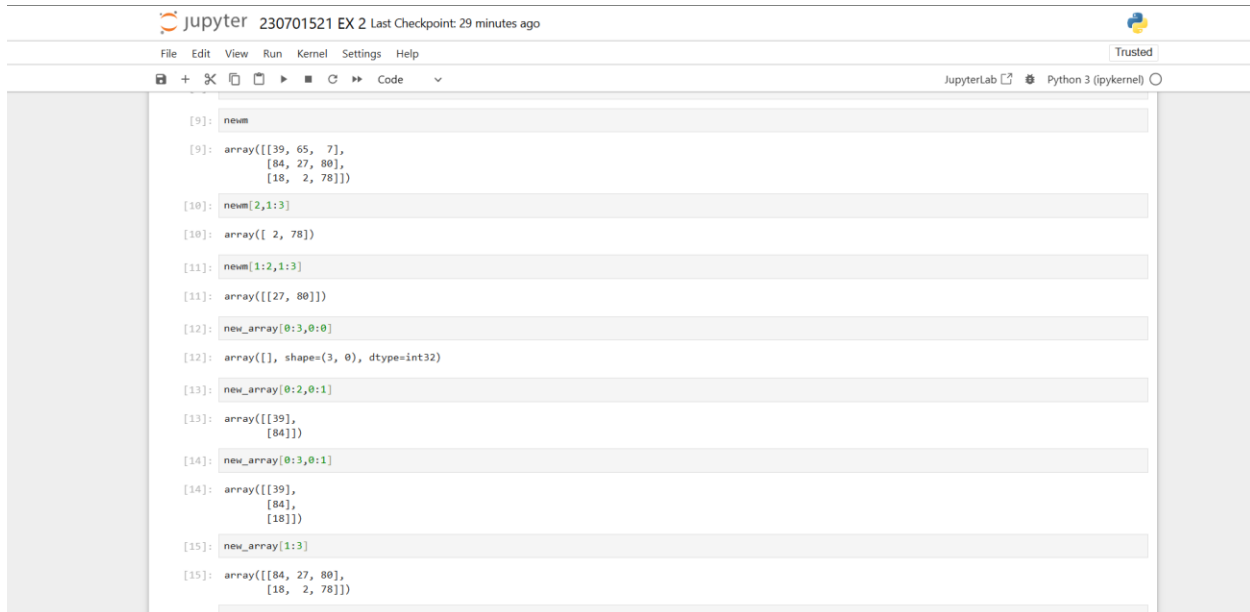
```
Out[9]: <Axes: xlabel='petal.length', ylabel='petal.width'>
```



```
In [10]: sns.pairplot(data,hue='variety',height=3);
```

230701521
JABARAJ E
“A”

2. NUMPY BUILT IN FUNCTION



The screenshot shows a JupyterLab interface with the title "230701521 EX 2 Last Checkpoint: 29 minutes ago". The interface includes a menu bar (File, Edit, View, Run, Kernel, Settings, Help) and a toolbar with icons for file operations, running, and code execution. The code area contains the following Python code:

```
[9]: newm
[9]: array([[39, 65, 7],
          [84, 27, 80],
          [18, 2, 78]])

[10]: newm[2,1:3]
[10]: array([ 2, 78])

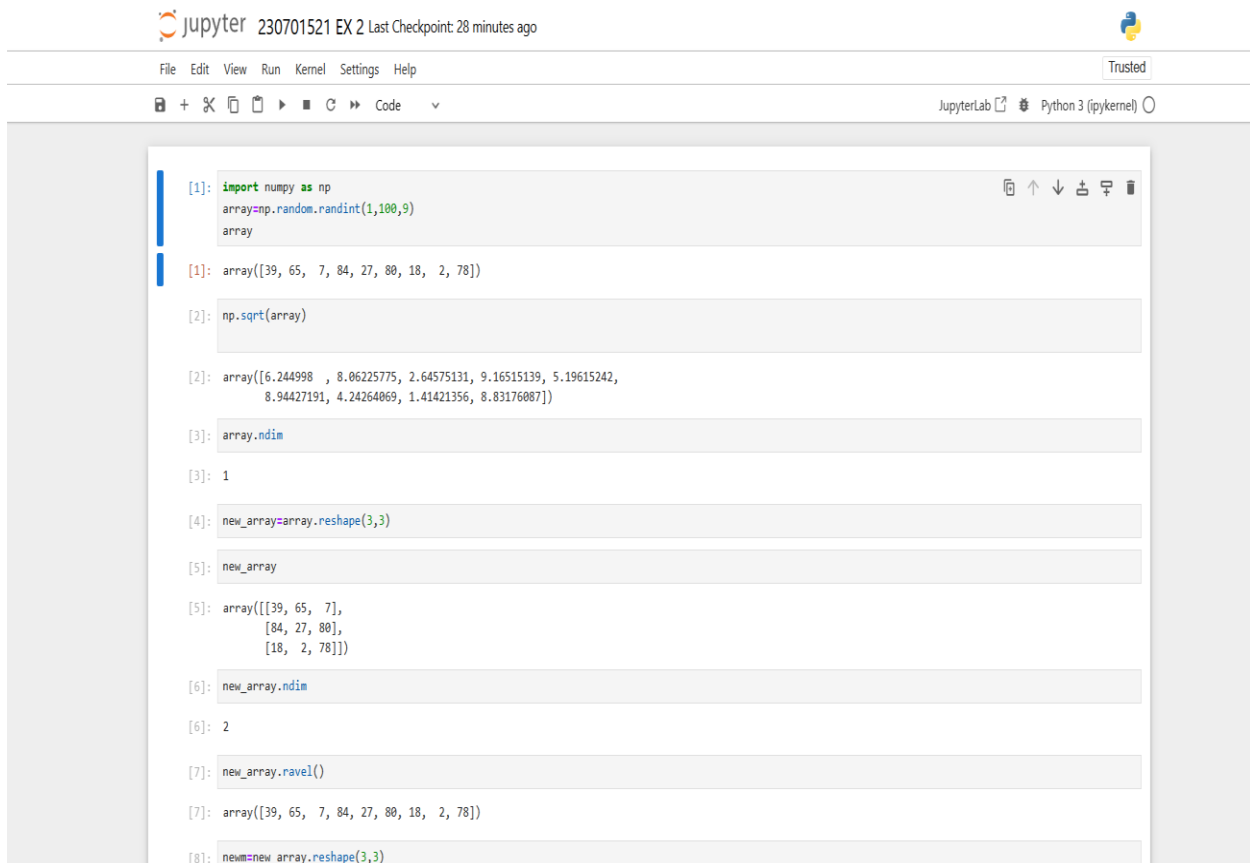
[11]: newm[1:2,1:3]
[11]: array([[27, 80]])

[12]: new_array[0:3,0:0]
[12]: array([], shape=(3, 0), dtype=int32)

[13]: new_array[0:2,0:1]
[13]: array([[39],
          [84]])

[14]: new_array[0:3,0:1]
[14]: array([[39],
          [84],
          [18]])

[15]: new_array[1:3]
[15]: array([[84, 27, 80],
          [18, 2, 78]])
```



The screenshot shows a JupyterLab interface with the title "230701521 EX 2 Last Checkpoint: 28 minutes ago". The interface includes a menu bar (File, Edit, View, Run, Kernel, Settings, Help) and a toolbar with icons for file operations, running, and code execution. The code area contains the following Python code:

```
[1]: import numpy as np
    array=np.random.randint(1,100,9)
    array

[1]: array([39, 65, 7, 84, 27, 80, 18, 2, 78])

[2]: np.sqrt(array)

[2]: array([6.244998 , 8.06225775, 2.64575131, 9.16515139, 5.19615242,
          8.94427191, 4.24264069, 1.41421356, 8.83176087])

[3]: array.ndim

[3]: 1

[4]: new_array=array.reshape(3,3)

[5]: new_array

[5]: array([[39, 65, 7],
          [84, 27, 80],
          [18, 2, 78]])

[6]: new_array.ndim

[6]: 2

[7]: new_array.ravel()

[7]: array([39, 65, 7, 84, 27, 80, 18, 2, 78])

[8]: newm=new_array.reshape(3,3)
```

“A”

Jupyter230701521 EX3Last Checkpoint: 13 minutes ago

FileEditViewRunKernelSettingsHelpTrusted

JupyterLabPython 3 (ipykernel)

[10]:import numpy as np
import pandas as pd
lists=[1,'Smith',50000],[2,'Jones',60000]
import warnings
warnings.filterwarnings('ignore')

[2]:df=pd.DataFrame(lists)
df

[2]:

	0	1	2
0	1	Smith	50000
1	2	Jones	60000

[3]:df.columns=['EmpId','Name','Salary']
df

[3]:

	EmpId	Name	Salary
0	1	Smith	50000
1	2	Jones	60000

[4]:df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2 entries, 0 to 1
Data columns (total 3 columns):
Column Non-Null Count Dtype

0 EmpId 2 non-null int64
1 Name 2 non-null object
2 Salary 2 non-null int64

Jupyter230701521 EX3Last Checkpoint: 14 minutes ago

FileEditViewRunKernelSettingsHelpTrusted

JupyterLabPython 3 (ipykernel)

[5]:df=pd.read_csv('50_Startups.csv')
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50 entries, 0 to 49
Data columns (total 5 columns):
Column Non-Null Count Dtype

0 R&D Spend 50 non-null float64
1 Administration 50 non-null float64
2 Marketing Spend 50 non-null float64
3 State 50 non-null object
4 Profit 50 non-null float64
dtypes: float64(4), object(1)
memory usage: 2.1+ KB

[6]:df.head(5)

[6]:

	R&D Spend	Administration	Marketing Spend	State	Profit
0	165349.20	136897.80	471784.10	New York	192261.83
1	162597.70	151377.59	443898.53	California	191792.06
2	153441.51	101145.55	407934.54	Florida	191050.39
3	144372.41	118671.85	383199.62	New York	182901.99
4	142107.34	91391.77	366168.42	Florida	166187.94

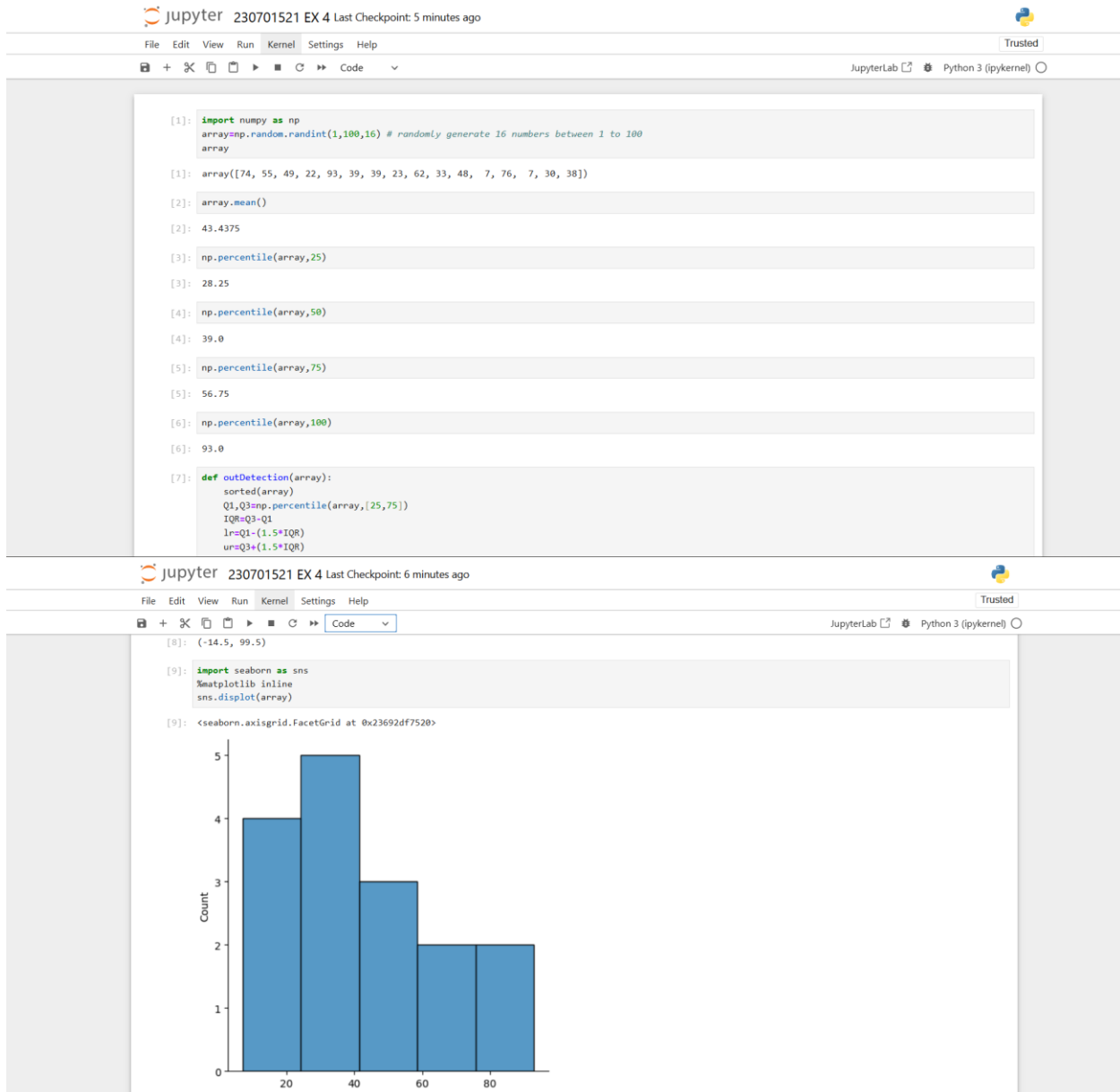
[7]:df.tail(5)

[7]:

	R&D Spend	Administration	Marketing Spend	State	Profit
45	1000.23	124153.04	1903.93	New York	64926.08
46	1215.45	115016.24	207114.45	Florida	40400.75

230701521
JABARAJ E
“A”

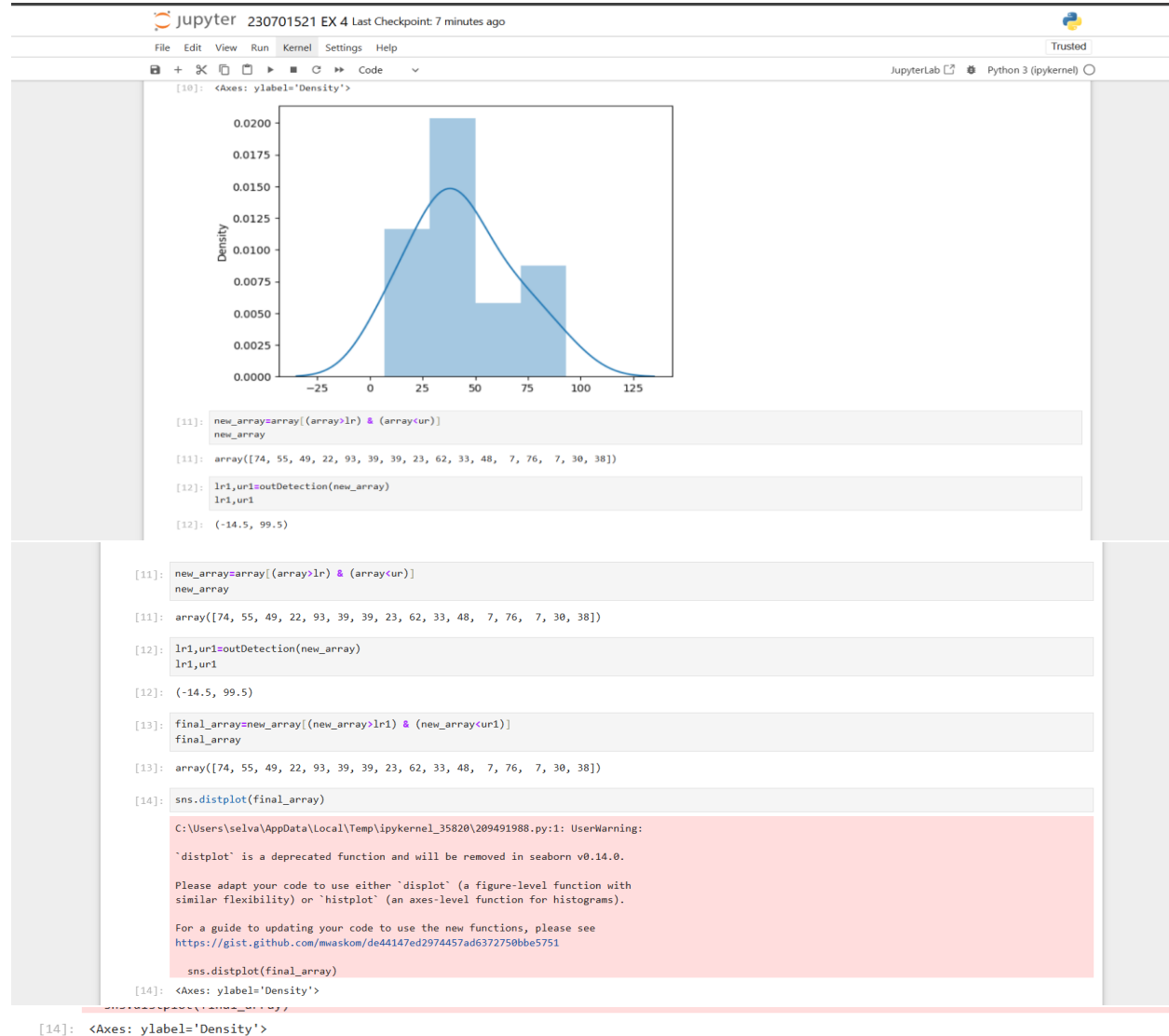
4.OUTLIER DETECTION



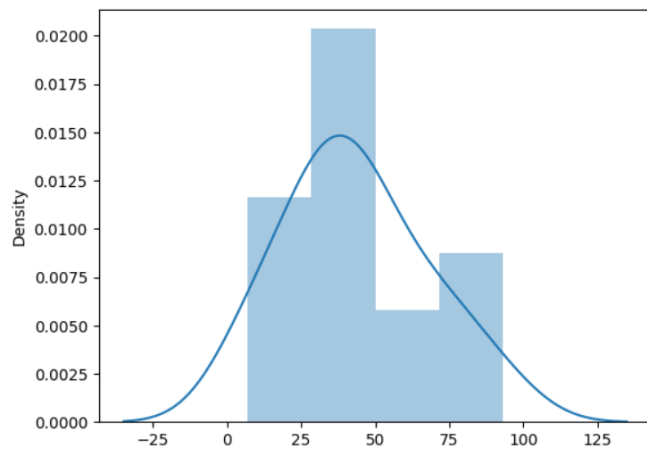
230701521

JABARAJ E

"A"



[14]: <Axes: ylabel='Density'>



230701521
JABARAJ E
“A”

5.MISSING AND INAPPROPRIATE DATA

Jupyter 230701521 EX 5 Last Checkpoint: 25 minutes ago

File Edit View Run Kernel Settings Help

JupyterLab Python 3 (ipykernel)

```
[16]: import numpy as np
import pandas as pd
import warnings
warnings.filterwarnings('ignore')
df=pd.read_csv("Hotel_Dataset.csv")
df
```

	CustomerID	Age_Group	Rating(1-5)	Hotel	FoodPreference	Bill	NoOfPax	EstimatedSalary	Age_Group.1
0	1	20-25	4	Ibis	veg	1300	2	40000	20-25
1	2	30-35	5	LemonTree	Non-Veg	2000	3	59000	30-35
2	3	25-30	6	RedFox	Veg	1322	2	30000	25-30
3	4	20-25	-1	LemonTree	Veg	1234	2	120000	20-25
4	5	35+	3	Ibis	Vegetarian	989	2	45000	35+
5	6	35+	3	Ibys	Non-Veg	1909	2	122220	35+
6	7	35+	4	RedFox	Vegetarian	1000	-1	21122	35+
7	8	20-25	7	LemonTree	Veg	2999	-10	345673	20-25
8	9	25-30	2	Ibis	Non-Veg	3456	3	-99999	25-30
9	9	25-30	2	Ibis	Non-Veg	3456	3	-99999	25-30
10	10	30-35	5	RedFox	non-Veg	-6755	4	87777	30-35

```
[2]: df.duplicated()
[2]: 0    False
1    False
2    False
```

230701521

JABARAJ E

“A”

Jupyter

230701521 EX 5 Last Checkpoint: 26 minutes ago

Trusted

File Edit View Run Kernel Settings Help

JupyterLab Python 3 (ipykernel)

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

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11 entries, 0 to 10
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype  
---  --
0   CustomerID            11 non-null     int64  
1   Age_Group             11 non-null     object  
2   Rating(1-5)           11 non-null     int64  
3   Hotel                 11 non-null     object  
4   FoodPreference         11 non-null     object  
5   Bill                  11 non-null     int64  
6   NoOfPax               11 non-null     int64  
7   EstimatedSalary        11 non-null     int64  
8   Age_Group.1           11 non-null     object  
dtypes: int64(5), object(4)
memory usage: 920.0+ bytes
```

```
[4]: df.drop_duplicates(inplace=True)
df
```

```
[4]:
```

	CustomerID	Age_Group	Rating(1-5)	Hotel	FoodPreference	Bill	NoOfPax	EstimatedSalary	Age_Group.1
0	1	20-25	4	Ibis	veg	1300	2	40000	20-25
1	2	30-35	5	LemonTree	Non-Veg	2000	3	59000	30-35
2	3	25-30	6	RedFox	Veg	1322	2	30000	25-30
3	4	20-25	-1	LemonTree	Veg	1234	2	120000	20-25
4	5	35+	3	Ibis	Vegetarian	989	2	45000	35+
5	6	35+	3	Ibys	Non-Veg	1909	2	122220	35+
6	7	35+	4	RedFox	Vegetarian	1000	-1	21122	35+

Jupyter

230701521 EX 5 Last Checkpoint: 30 minutes ago

Trusted

File Edit View Run Kernel Settings Help

JupyterLab Python 3 (ipykernel)

```
[14]: df.FoodPreference.replace(['Vegetarian','veg'],'Veg',inplace=True)
df.FoodPreference.replace(['non-Veg'],'Non-Veg',inplace=True)
```

```
[15]: df.EstimatedSalary.fillna(round(df.EstimatedSalary.mean()),inplace=True)
df.NoOfPax.fillna(round(df.NoOfPax.median()),inplace=True)
df['Rating(1-5)'].fillna(round(df['Rating(1-5)'].median()), inplace=True)
df.Bill.fillna(round(df.Bill.mean()),inplace=True)
df
```

C:\Users\selsva\AppData\Local\Temp\ipykernel_27116\3711388855.py:3: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
df['Rating(1-5)'].fillna(round(df['Rating(1-5)'].median()), inplace=True)
```

```
[15]:
```

	CustomerID	Age_Group	Rating(1-5)	Hotel	FoodPreference	Bill	NoOfPax	EstimatedSalary
0	1.0	20-25	4	Ibis	Veg	1300.0	2.0	40000.0
1	2.0	30-35	5	LemonTree	Non-Veg	2000.0	3.0	59000.0
2	3.0	25-30	6	RedFox	Veg	1322.0	2.0	30000.0
3	4.0	20-25	-1	LemonTree	Veg	1234.0	2.0	120000.0
4	5.0	35+	3	Ibis	Veg	989.0	2.0	45000.0
5	6.0	35+	3	Ibis	Non-Veg	1909.0	2.0	122220.0
6	7.0	35+	4	RedFox	Veg	1000.0	2.0	21122.0
7	8.0	20-25	7	LemonTree	Veg	2999.0	2.0	345673.0

230701521

JABARAJ E

“A”

Jupyter

230701521 EX 5

Last Checkpoint: 30 minutes ago

File

Edit

View

Run

Kernel

Settings

Help

Trusted

JupyterLab

Python 3 (ipykernel)

[14]:

df.FoodPreference.replace(['Vegetarian','veg'],'Veg',inplace=True)
df.FoodPreference.replace(['non-Veg'],'Non-Veg',inplace=True)

[15]:

df.EstimatedSalary.fillna(round(df.EstimatedSalary.mean()),inplace=True)
df.NoOfPax.fillna(round(df.NoOfPax.median()),inplace=True)
df['Rating(1-5)'].fillna(round(df['Rating(1-5)'].median()), inplace=True)
df.Bill.fillna(round(df.Bill.mean()),inplace=True)
df

C:\Users\selsva\AppData\Local\Temp\ipykernel_27116\371138885.py:3: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.
The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

df['Rating(1-5)'].fillna(round(df['Rating(1-5)'].median()), inplace=True)

[15]:

	CustomerID	Age_Group	Rating(1-5)	Hotel	FoodPreference	Bill	NoOfPax	EstimatedSalary
0	1.0	20-25	4	Ibis	Veg	1300.0	2.0	40000.0
1	2.0	30-35	5	LemonTree	Non-Veg	2000.0	3.0	59000.0
2	3.0	25-30	6	RedFox	Veg	1322.0	2.0	30000.0
3	4.0	20-25	-1	LemonTree	Veg	1234.0	2.0	120000.0
4	5.0	35+	3	Ibis	Veg	989.0	2.0	45000.0
5	6.0	35+	3	Ibis	Non-Veg	1909.0	2.0	122220.0
6	7.0	35+	4	RedFox	Veg	1000.0	2.0	21122.0
7	8.0	20-25	7	LemonTree	Veg	2999.0	2.0	345673.0

Jupyter

230701521 EX 5

Last Checkpoint: 29 minutes ago

File

Edit

View

Run

Kernel

Settings

Help

Trusted

JupyterLab

Python 3 (ipykernel)

A value is trying to be set on a copy of a slice from a DataFrame
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df['NoOfPax'].loc[(df['NoOfPax']<1) | (df['NoOfPax']>20)]=np.nan

[10]:

	CustomerID	Age_Group	Rating(1-5)	Hotel	FoodPreference	Bill	NoOfPax	EstimatedSalary
0	1.0	20-25	4	Ibis	veg	1300.0	2.0	40000.0
1	2.0	30-35	5	LemonTree	Non-Veg	2000.0	3.0	59000.0
2	3.0	25-30	6	RedFox	Veg	1322.0	2.0	30000.0
3	4.0	20-25	-1	LemonTree	Veg	1234.0	2.0	120000.0
4	5.0	35+	3	Ibis	Vegetarian	989.0	2.0	45000.0
5	6.0	35+	3	Ibys	Non-Veg	1909.0	2.0	122220.0
6	7.0	35+	4	RedFox	Vegetarian	1000.0	NaN	21122.0
7	8.0	20-25	7	LemonTree	Veg	2999.0	NaN	345673.0
8	9.0	25-30	2	Ibis	Non-Veg	3456.0	3.0	NaN
9	10.0	30-35	5	RedFox	non-Veg	NaN	4.0	87777.0

[11]:

df.Age_Group.unique()

[11]:

array(['20-25', '30-35', '25-30', '35+'], dtype=object)

[12]:

df.Hotel.unique()

[12]:

array(['Ibis', 'LemonTree', 'RedFox', 'Ibys'], dtype=object)

[13]:

df.Hotel.replace(['Ibys'],'Ibis',inplace=True)
df.FoodPreference.unique

230701521

JABARAJ E

“A”

jupyter

230701521 EX 5 Last Checkpoint: 28 minutes ago

File

Edit

View

Run

Kernel

Settings

Help

Trusted

+

✕

Code

JupyterLabPython 3 (ipykernel)

C:\Users\selsva\AppData\Local\Temp\ipykernel_27116\2080958306.py:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df.EstimatedSalary.loc[df.EstimatedSalary<0]=np.nan

[9]:

	CustomerID	Age_Group	Rating(1-5)	Hotel	FoodPreference	Bill	NoOfPax	EstimatedSalary
0	1.0	20-25	4	Ibhis	veg	1300.0	2	40000.0
1	2.0	30-35	5	LemonTree	Non-Veg	2000.0	3	59000.0
2	3.0	25-30	6	RedFox	Veg	1322.0	2	30000.0
3	4.0	20-25	-1	LemonTree	Veg	1234.0	2	120000.0
4	5.0	35+	3	Ibhis	Vegetarian	989.0	2	45000.0
5	6.0	35+	3	Ibhis	Non-Veg	1909.0	2	122220.0
6	7.0	35+	4	RedFox	Vegetarian	1000.0	-1	21122.0
7	8.0	20-25	7	LemonTree	Veg	2999.0	-10	345673.0
8	9.0	25-30	2	Ibhis	Non-Veg	3456.0	3	NaN
9	10.0	30-35	5	RedFox	non-Veg	NaN	4	87777.0

[10]:df['NoOfPax'].loc[(df['NoOfPax']<1) | (df['NoOfPax']>20)]=np.nan
df

C:\Users\selsva\AppData\Local\Temp\ipykernel_27116\2129877948.py:1: FutureWarning: ChainedAssignmentError: behaviour will change in pandas 3.0!
You are setting values through chained assignment. Currently this works in certain cases, but when using Copy-on-Write (which will become the default behaviour in pandas 3.0) this will never work to update the original DataFrame or Series, because the intermediate object on which we are setting values will behave as a copy.
A typical example is when you are setting values in a column of a DataFrame, like:

230701521

JABARAJ E

“A”

Jupyter 230701521 EX 5 Last Checkpoint: 27 minutes ago



File Edit View Run Kernel Settings Help

Trusted

Python 3 (ipykernel)

```
[5]: df.drop_duplicates(inplace=True)
df
```

```
[5]:
```

	CustomerID	Age_Group	Rating(1-5)	Hotel	FoodPreference	Bill	NoOfPax	EstimatedSalary	Age_Group.1
0	1	20-25	4	Ibis	veg	1300	2	40000	20-25
1	2	30-35	5	LemonTree	Non-Veg	2000	3	59000	30-35
2	3	25-30	6	RedFox	Veg	1322	2	30000	25-30
3	4	20-25	-1	LemonTree	Veg	1234	2	120000	20-25
4	5	35+	3	Ibis	Vegetarian	989	2	45000	35+
5	6	35+	3	Ibys	Non-Veg	1909	2	122220	35+
6	7	35+	4	RedFox	Vegetarian	1000	-1	21122	35+
7	8	20-25	7	LemonTree	Veg	2999	-10	345673	20-25
8	9	25-30	2	Ibis	Non-Veg	3456	3	-99999	25-30
10	10	30-35	5	RedFox	non-Veg	-6755	4	87777	30-35

```
[6]: index=np.array(list(range(0,len(df))))
df.set_index(index,inplace=True)
index
```

```
[6]: array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
```

Jupyter 230701521 EX 5 Last Checkpoint: 28 minutes ago



File Edit View Run Kernel Settings Help

Trusted

Python 3 (ipykernel)

```
[7]:
```

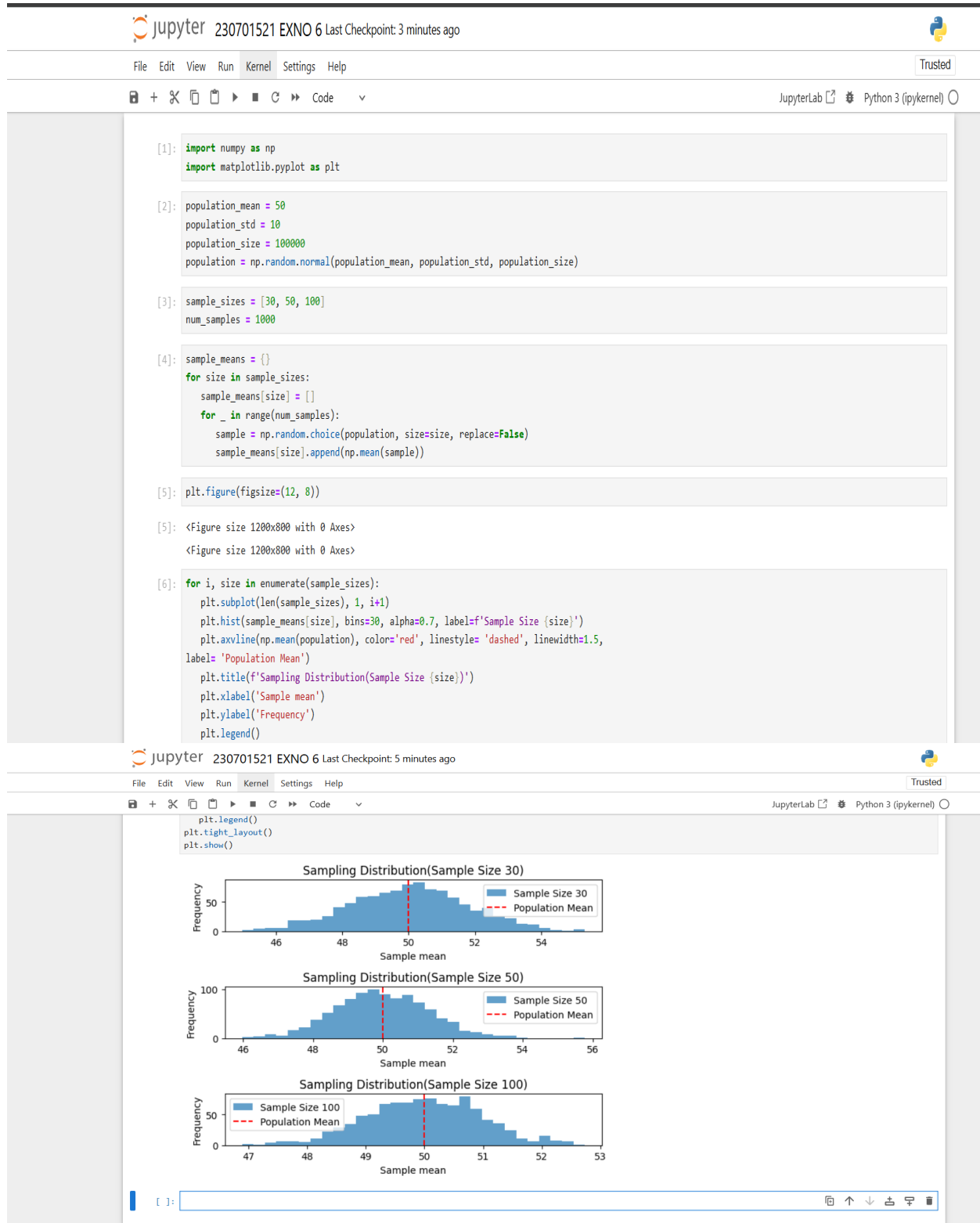
	CustomerID	Age_Group	Rating(1-5)	Hotel	FoodPreference	Bill	NoOfPax	EstimatedSalary	Age_Group.1
0	1	20-25	4	Ibis	veg	1300	2	40000	20-25
1	2	30-35	5	LemonTree	Non-Veg	2000	3	59000	30-35
2	3	25-30	6	RedFox	Veg	1322	2	30000	25-30
3	4	20-25	-1	LemonTree	Veg	1234	2	120000	20-25
4	5	35+	3	Ibis	Vegetarian	989	2	45000	35+
5	6	35+	3	Ibys	Non-Veg	1909	2	122220	35+
6	7	35+	4	RedFox	Vegetarian	1000	-1	21122	35+
7	8	20-25	7	LemonTree	Veg	2999	-10	345673	20-25
8	9	25-30	2	Ibis	Non-Veg	3456	3	-99999	25-30
9	10	30-35	5	RedFox	non-Veg	-6755	4	87777	30-35

```
[8]: df.drop(['Age_Group.1'],axis=1,inplace=True)
df
```

```
[8]:
```

	CustomerID	Age_Group	Rating(1-5)	Hotel	FoodPreference	Bill	NoOfPax	EstimatedSalary
0	1	20-25	4	Ibis	veg	1300	2	40000
1	2	30-35	5	LemonTree	Non-Veg	2000	3	59000
2	3	25-30	6	RedFox	Veg	1322	2	30000
3	4	20-25	-1	LemonTree	Veg	1234	2	120000
4	5	35+	3	Ibis	Vegetarian	989	2	45000
5	6	35+	3	Ibys	Non-Veg	1909	2	122220
6	7	35+	4	RedFox	Vegetarian	1000	-1	21122
7	8	20-25	7	LemonTree	Veg	2999	-10	345673
8	9	25-30	2	Ibis	Non-Veg	3456	3	-99999
9	10	30-35	5	RedFox	non-Veg	-6755	4	87777

6. RANDOM SAMPLING AND DISTRIBUTION

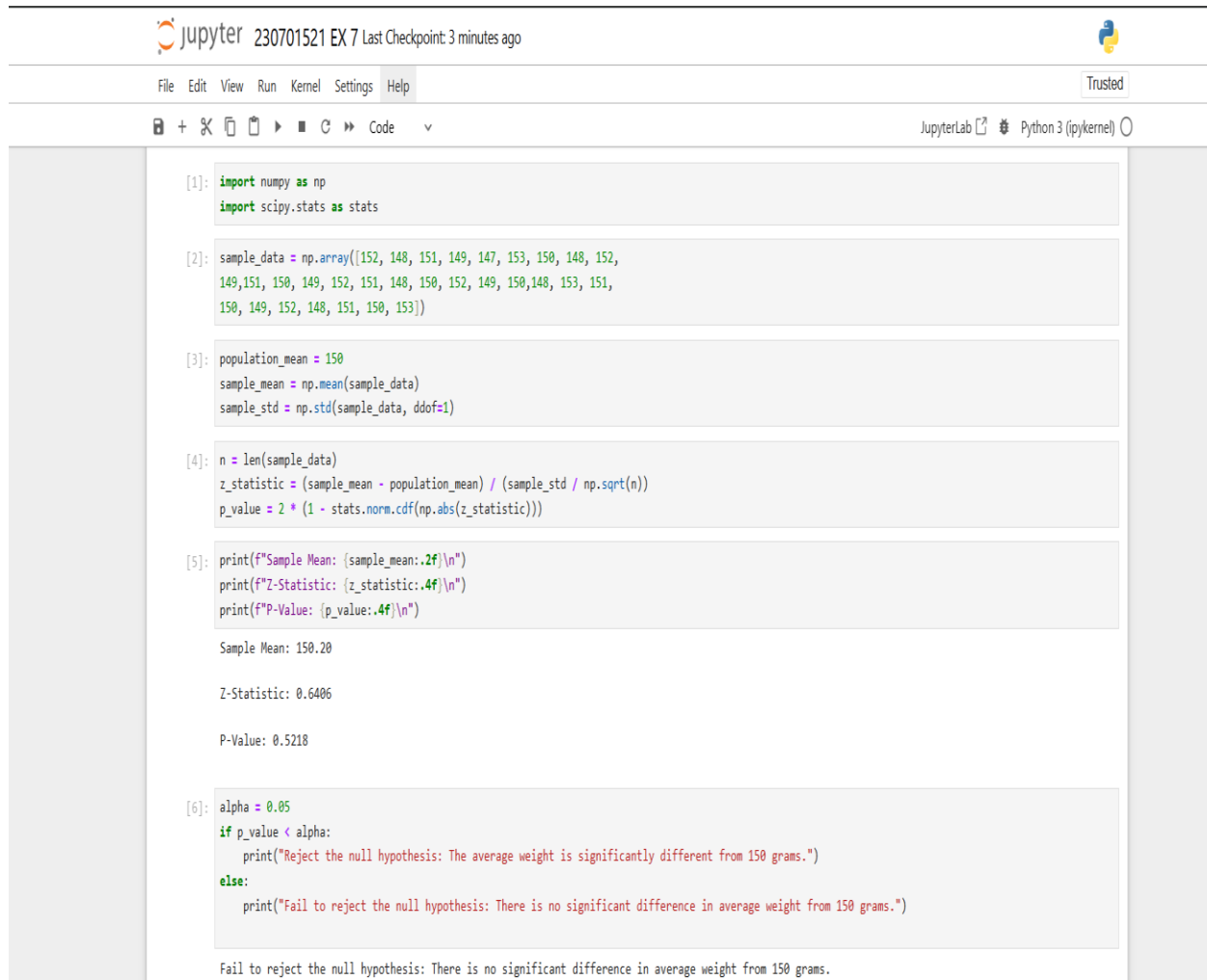


230701521

JABARAJ E

“A”

7. Z TEST



The image shows a JupyterLab interface with a single code cell containing Python code for a Z-test. The code imports numpy and scipy.stats, creates a sample data array, calculates the sample mean and standard deviation, and then computes the Z-statistic and p-value. The output of the code is displayed below the code cell, showing the sample mean, Z-statistic, and p-value. The code also includes a conditional statement to reject or fail to reject the null hypothesis based on the p-value.

```
[1]: import numpy as np
import scipy.stats as stats

[2]: sample_data = np.array([152, 148, 151, 149, 147, 153, 150, 148, 152,
149, 151, 150, 149, 152, 151, 148, 150, 152, 149, 150, 148, 153, 151,
150, 149, 152, 148, 151, 150, 153])

[3]: population_mean = 150
sample_mean = np.mean(sample_data)
sample_std = np.std(sample_data, ddof=1)

[4]: n = len(sample_data)
z_statistic = (sample_mean - population_mean) / (sample_std / np.sqrt(n))
p_value = 2 * (1 - stats.norm.cdf(np.abs(z_statistic)))

[5]: print(f"Sample Mean: {sample_mean:.2f}\n")
print(f"Z-Statistic: {z_statistic:.4f}\n")
print(f"P-Value: {p_value:.4f}\n")

Sample Mean: 150.20

Z-Statistic: 0.6406

P-Value: 0.5218

[6]: alpha = 0.05
if p_value < alpha:
    print("Reject the null hypothesis: The average weight is significantly different from 150 grams.")
else:
    print("Fail to reject the null hypothesis: There is no significant difference in average weight from 150 grams.")

Fail to reject the null hypothesis: There is no significant difference in average weight from 150 grams.
```

230701521
JABARAJ E
"A"

8. T- TEST



The image shows a JupyterLab interface with a code editor and output area. The code implements a T-test to check if a sample mean is significantly different from a population mean of 100. The output shows the sample mean, T-statistic, P-value, and a conclusion to fail to reject the null hypothesis.

```
[1]: import numpy as np
import scipy.stats as stats
np.random.seed(42)
sample_size = 25
sample_data = np.random.normal(loc=102, scale=15, size=sample_size)

[2]: population_mean = 100
sample_mean = np.mean(sample_data)
sample_std = np.std(sample_data, ddof=1)

[3]: n = len(sample_data)
t_statistic, p_value = stats.ttest_1samp(sample_data, population_mean)

[5]: print(f"Sample Mean: {sample_mean:.2f}\n")
print(f"T-Statistic: {t_statistic:.4f}\n")
print(f"P-Value: {p_value:.4f}\n")
alpha = 0.05
if p_value < alpha:
    print("Reject the null hypothesis: The average IQ score is significantly different from 100.")
else:
    print("Fail to reject the null hypothesis: There is no significant difference in average IQ score from 100.")
```

Sample Mean: 99.55

T-Statistic: -0.1577

P-Value: 0.8760

Fail to reject the null hypothesis: There is no significant difference in average IQ score from 100.

230701521

JABARAJ E

“A”

9.ANOVA TEST

```
import numpy as np
import scipy.stats as stats
from statsmodels.stats.multicomp import pairwise_tukeyhsd

np.random.seed(42)
n_plants = 25
```

[124]

```
growth_A = np.random.normal(loc=10, scale=2, size=n_plants)
growth_B = np.random.normal(loc=12, scale=3, size=n_plants)
growth_C = np.random.normal(loc=15, scale=2.5, size=n_plants)
```

[125]

```
all_data = np.concatenate([growth_A, growth_B, growth_C])
```

[126]

```
treatment_labels = ['A'] * n_plants + ['B'] * n_plants + ['C'] * n_plants
f_statistic, p_value = stats.f_oneway(growth_A, growth_B, growth_C)
```

[127]

```
mean_A = np.mean(growth_A)
mean_B = np.mean(growth_B)
mean_C = np.mean(growth_C)
print(f"Treatment A Mean Growth: {mean_A:.4f}")
print(f"Treatment B Mean Growth: {mean_B:.4f}")
print(f"Treatment C Mean Growth: {mean_C:.4f}")
print(f"F-Statistic: {f_statistic:.4f}")
```

[128]

```
print(f"F-Statistic: {f_statistic:.4f}")
print(f"P-Value: {p_value:.4f}")

alpha = 0.05
if p_value < alpha:
    print("Reject the null hypothesis: There is a significant difference in mean growth rates among the three treatments.")
else:
    print("Fail to reject the null hypothesis: There is no significant difference in mean growth rates among the three treatments.")

if p_value < alpha:
    tukey_results = pairwise_tukeyhsd(all_data, treatment_labels, alpha=0.05)

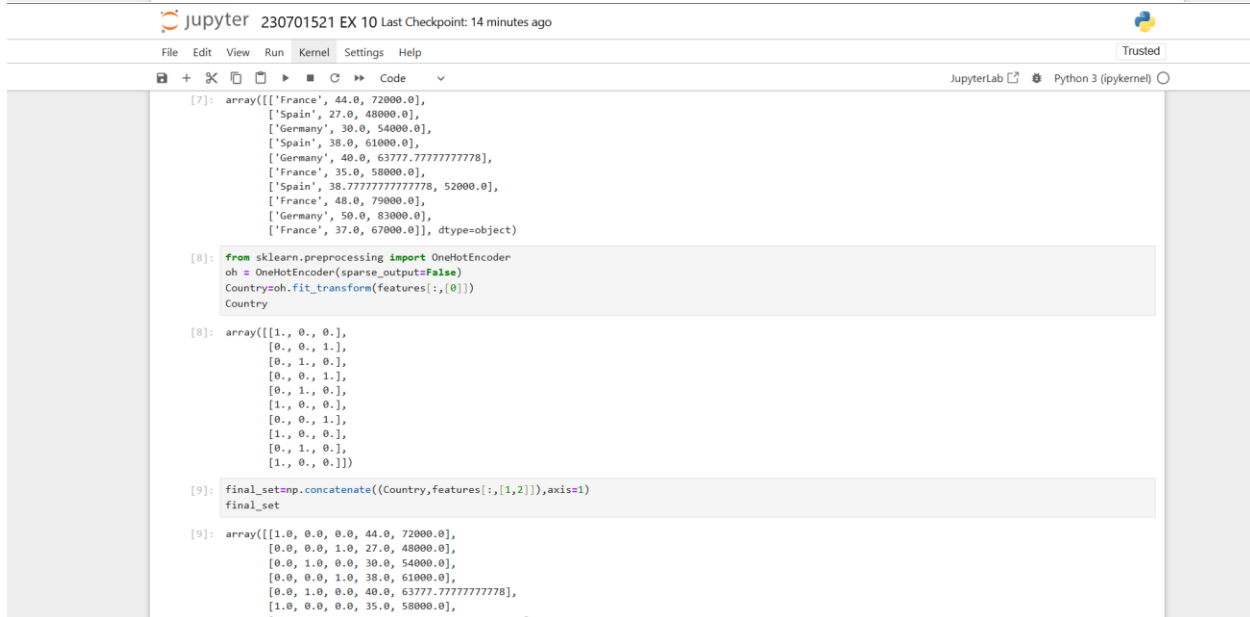
    print("\nTukey's HSD Post-hoc Test:")
    print(tukey_results)
```

[128]

```
... Treatment A Mean Growth: 9.6730
Treatment B Mean Growth: 11.1377
Treatment C Mean Growth: 15.2652
F-Statistic: 36.1214
P-Value: 0.0000
Reject the null hypothesis: There is a significant difference in mean growth rates among the three treatments.

Tukey's HSD Post-hoc Test:
Multiple Comparison of Means - Tukey HSD, FWER=0.05
=====
group1 group2 meandiff p-adj lower upper reject
-----
A B 1.4647 0.0877 -0.1683 3.0977 False
A C 5.5923 0.0 3.9593 7.2252 True
B C 4.1276 0.0 2.4946 5.7605 True
=====
```

“A”



230701521
JABARAJ E
“A”

```
jupyter 230701521 EX 10 Last Checkpoint: 13 minutes ago

File Edit View Run Kernel Settings Help Trusted

JupyterLab Python 3 (ipykernel)

[3]: array([[ 'France', 44.0, 72000.0],
          [ 'Spain', 27.0, 48000.0],
          [ 'Germany', 30.0, 54000.0],
          [ 'Spain', 38.0, 61000.0],
          [ 'Germany', 40.0, nan],
          [ 'France', 35.0, 58000.0],
          [ 'Spain', nan, 52000.0],
          [ 'France', 48.0, 79000.0],
          [ 'Germany', 50.0, 83000.0],
          [ 'France', 37.0, 67000.0]], dtype=object)

[4]: label=df.iloc[:, -1].values

[5]: from sklearn.impute import SimpleImputer
age=SimpleImputer(strategy="mean",missing_values=np.nan)
Salary=SimpleImputer(strategy="mean",missing_values=np.nan)
age.fit(features[:,1:])

[5]: SimpleImputer
SimpleImputer()

[6]: Salary.fit(features[:,2:])

[6]: SimpleImputer
SimpleImputer()

[7]: features[:,1]=age.transform(features[:,1])
features[:,2]=Salary.transform(features[:,2])
features
```

```
jupyter 230701521 EX 10 Last Checkpoint: 14 minutes ago

File Edit View Run Kernel Settings Help Trusted

JupyterLab Python 3 (ipykernel)

[12]: array([[ 7.58874362e-01,  7.49473254e-01],
          [-8.16496581e-01, -6.54653671e-01,  1.52752523e+00,
          -1.71150388e+00, -1.43817841e+00],
          [-8.16496581e-01,  1.52752523e+00, -6.54653671e-01,
          -1.27555478e+00, -8.91265492e-01],
          [-8.16496581e-01, -6.54653671e-01,  1.52752523e+00,
          -1.13023841e-01, -2.53200424e-01],
          [-8.16496581e-01,  1.52752523e+00, -6.54653671e-01,
          1.77608893e-01,  6.63219199e-16],
          [ 1.22474487e+00, -6.54653671e-01, -6.54653671e-01,
          -5.48972942e-01, -5.26656882e-01],
          [-8.16496581e-01, -6.54653671e-01,  1.52752523e+00,
          0.00000000e+00, -1.07356980e+00],
          [ 1.22474487e+00, -6.54653671e-01, -6.54653671e-01,
          1.34013983e+00,  1.38753832e+00],
          [-8.16496581e-01,  1.52752523e+00, -6.54653671e-01,
          1.63077256e+00,  1.75214693e+00],
          [ 1.22474487e+00, -6.54653671e-01, -6.54653671e-01,
          -2.58340208e-01,  2.93712492e-01]])

[12]: from sklearn.preprocessing import MinMaxScaler
mms=MinMaxScaler(feature_range=(0,1))
mms.fit(final_set)
feat_minmax_scaler=mms.transform(final_set)
feat_minmax_scaler

[12]: array([[1.      , 0.      , 0.      , 0.73913043, 0.68571429],
          [0.      , 0.      , 1.      , 0.      , 0.      ],
          [0.      , 1.      , 0.      , 0.13043478, 0.17142857],
          [0.      , 0.      , 1.      , 0.47826087, 0.37142857],
          [0.      , 1.      , 0.      , 0.56521739, 0.45079365],
          [1.      , 0.      , 0.      , 0.34782609, 0.28571429],
          [0.      , 0.      , 1.      , 0.51207729, 0.11428571],
          [1.      , 0.      , 0.      , 0.91304348, 0.88571429],
          [0.      , 1.      , 0.      , 1.      , 1.      ],
          [1.      , 0.      , 0.      , 0.43478261, 0.54285714]])
```

230701521

JABARAJ E

“A”

11. LINEAR REGRESSION

Jupyter 230701521 EX 11 Last Checkpoint: 18 minutes ago

File Edit View Run Kernel Settings Help

Trusted

JupyterLab Python 3 (ipykernel)

```
[1]: import numpy as np
import pandas as pd
df = pd.read_csv('Salary_data.csv')
df
```

C:\Users\selsva\AppData\Local\Temp\ipykernel_35980\978986724.py:2: DeprecationWarning:
Pyarrow will become a required dependency of pandas in the next major release of pandas (pandas 3.0),
(to allow more performant data types, such as the Arrow string type, and better interoperability with other libraries)
but was not found to be installed on your system.
If this would cause problems for you,
please provide us feedback at <https://github.com/pandas-dev/pandas/issues/54466>

import pandas as pd

```
[1]:
```

	YearsExperience	Salary
0	1.1	39343
1	1.3	46205
2	1.5	37731
3	2.0	43525
4	2.2	39891
5	2.9	56642
6	3.0	60150
7	3.2	54445
8	3.2	64445
9	3.7	57189

Jupyter 230701521 EX 11 Last Checkpoint: 20 minutes ago

File Edit View Run Kernel Settings Help

Trusted

JupyterLab Python 3 (ipykernel)

```
[4]: df.info()
```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30 entries, 0 to 29
Data columns (total 2 columns):
Column Non-Null Count Dtype
--- ---
0 YearsExperience 30 non-null float64
1 Salary 30 non-null int64
dtypes: float64(1), int64(1)
memory usage: 608.0 bytes

```
[5]: df.describe()
```

	YearsExperience	Salary
count	30.000000	30.000000
mean	5.313333	76003.000000
std	2.837888	27414.429785
min	1.100000	37731.000000
25%	3.200000	56720.750000
50%	4.700000	65237.000000
75%	7.700000	100544.750000
max	10.500000	122391.000000

```
[6]: features = df.iloc[:,[0]].values
label = df.iloc[:,[1]].values
features
```

230701521

JABARAJ E

“A”

Jupyter 230701521 EX 11 Last Checkpoint: 19 minutes ago

File Edit View Run Kernel Settings Help Trusted

JupyterLab Python 3 (ipykernel)

```
[2]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30 entries, 0 to 29
Data columns (total 2 columns):
 #   Column        Non-Null Count  Dtype  
---  --
 0   YearsExperience  30 non-null     float64
 1   Salary          30 non-null     int64   
dtypes: float64(1), int64(1)
memory usage: 608.0 bytes
```

```
[3]: df.dropna(inplace=True);
df
```

```
[3]:
```

	YearsExperience	Salary
0	1.1	39343
1	1.3	46205
2	1.5	37731
3	2.0	43525
4	2.2	39891
5	2.9	56642
6	3.0	60150
7	3.2	54445
8	3.2	64445

Jupyter 230701521 EX 11 Last Checkpoint: 20 minutes ago

File Edit View Run Kernel Settings Help Trusted

JupyterLab Python 3 (ipykernel)

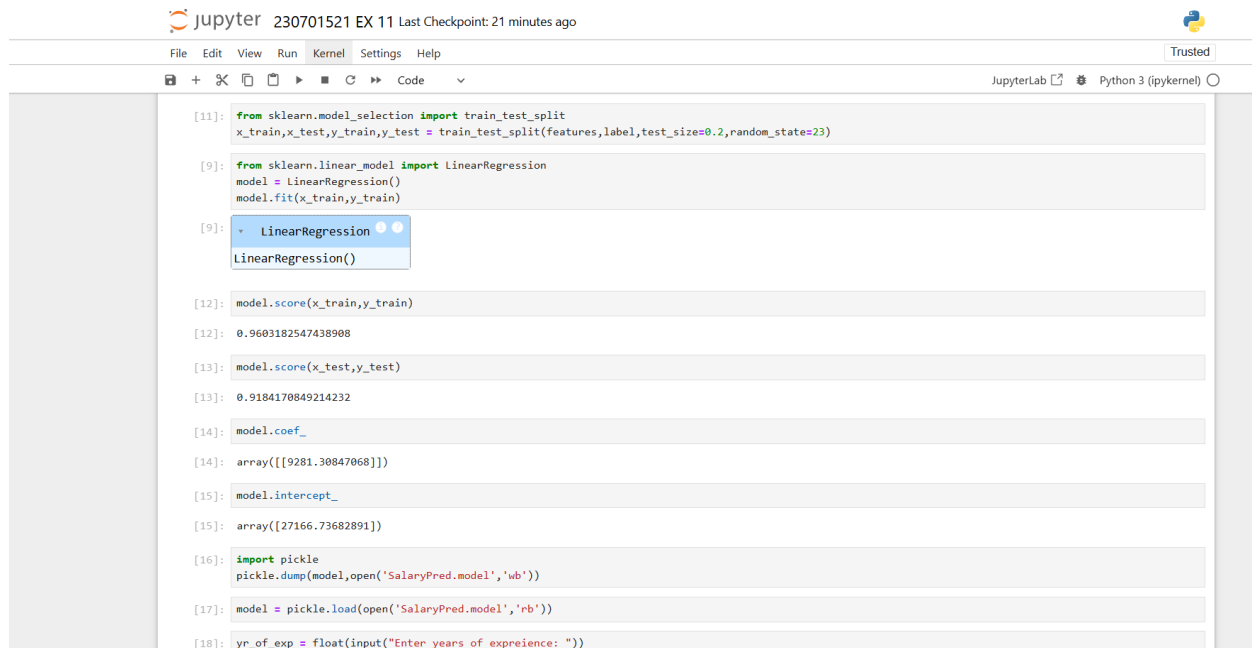
```
[7]: label
```

```
[7]: array([[ 39343],
 [ 46205],
 [ 37731],
 [ 43525],
 [ 39891],
 [ 56642],
 [ 60150],
 [ 54445],
 [ 64445],
 [ 57189],
 [ 63218],
 [ 55794],
 [ 56957],
 [ 57081],
 [ 61111],
```

230701521

JABARAJ E

“A”



The image shows a JupyterLab interface with a Python 3 (ipykernel) environment. The interface includes a top bar with the Jupyter logo, the text "jupyter 230701521 EX 11 Last Checkpoint: 21 minutes ago", and a "Trusted" status indicator. Below the top bar is a menu bar with "File", "Edit", "View", "Run", "Kernel", "Settings", and "Help". A toolbar with various icons is located below the menu bar. The main area displays a series of code cells and their outputs:

```
[11]: from sklearn.model_selection import train_test_split
      x_train,x_test,y_train,y_test = train_test_split(features,label,test_size=0.2,random_state=23)

[9]: from sklearn.linear_model import LinearRegression
      model = LinearRegression()
      model.fit(x_train,y_train)

[9]: LinearRegression
      LinearRegression()

[12]: model.score(x_train,y_train)

[12]: 0.9603182547438908

[13]: model.score(x_test,y_test)

[13]: 0.9184170849214232

[14]: model.coef_

[14]: array([[9281.30847068]])

[15]: model.intercept_

[15]: array([27166.73682891])

[16]: import pickle
      pickle.dump(model,open('SalaryPred.model','wb'))

[17]: model = pickle.load(open('SalaryPred.model','rb'))

[18]: yr_of_exp = float(input("Enter years of expreience: "))
```

230701521
JABARAJ E
“A”

12. LOGISTIC REGRESSION

The image displays two screenshots of a JupyterLab environment. The top screenshot shows the initial setup where a CSV file named 'Social_Network_Ads.csv' is loaded into a pandas DataFrame. A deprecation warning for Pyarrow is shown. The DataFrame is then printed, showing columns: User ID, Gender, Age, EstimatedSalary, and Purchased. The bottom screenshot shows further data manipulation: the last five rows are viewed with 'df.tail(5)', the first five rows with 'df.head(5)', and the first three rows are selected as 'features' while the 'Purchased' column is the 'label'.

Top Screenshot:

```
[1]: import numpy as np
import pandas as pd
import warnings
warnings.filterwarnings('ignore')
df = pd.read_csv('Social_Network_Ads.csv')
df
```

C:\Users\selva\AppData\Local\Temp\ipykernel_26828\2834319537.py:2: DeprecationWarning: Pyarrow will become a required dependency of pandas in the next major release of pandas (pandas 3.0), (to allow more performant data types, such as the Arrow string type, and better interoperability with other libraries) but was not found to be installed on your system. If this would cause problems for you, please provide us feedback at <https://github.com/pandas-dev/pandas/issues/54466>

```
import pandas as pd
```

```
[1]:
```

	User ID	Gender	Age	EstimatedSalary	Purchased
0	15624510	Male	19	19000	0
1	15810944	Male	35	20000	0
2	15668575	Female	26	43000	0
3	15603246	Female	27	57000	0
4	15804002	Male	19	76000	0
...
395	15691863	Female	46	41000	1
396	15706071	Male	51	23000	1
397	15654296	Female	50	20000	1

Bottom Screenshot:

```
[2]: df.tail(5)
```

```
[2]:
```

395	15691863	Female	46	41000	1
396	15706071	Male	51	23000	1
397	15654296	Female	50	20000	1
398	15755018	Male	36	33000	0
399	15594041	Female	49	36000	1

```
[3]: df.head(5)
```

```
[3]:
```

	User ID	Gender	Age	EstimatedSalary	Purchased
0	15624510	Male	19	19000	0
1	15810944	Male	35	20000	0
2	15668575	Female	26	43000	0
3	15603246	Female	27	57000	0
4	15804002	Male	19	76000	0

```
[4]: features = df.iloc[:,2:3].values
label = df.iloc[:,4].values
features
```

```
[4]: array([[ 19, 19000],
        [ 35, 20000],
        [ 26, 43000],
```

230701521
JABARAJ E
"A"

Jupyter

230701521 ex 12 Last Checkpoint: 24 minutes ago

File Edit View Run Kernel Settings Help

Trusted

+ ↻ 📄 🔍 ⏮️ ⏪ ⏩ ⏭ Code ▾

JupyterLab Python 3 (ipykernel)

[5]: label

[5]: array([[0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
0,
0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1,
0, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 0, 1, 1, 0,
1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 1, 0, 1, 1, 0,
1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 0, 1,
0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 0, 1,
1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 1, 1,
0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0,
1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0,
0, 1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1,
1, 1, 0, 1], dtype=int64)

[6]: from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression

[7]: for i in range(1, 401):
x_train, x_test, y_train, y_test = train_test_split(features, label, test_size=0.2, random_state=i)
model = LogisticRegression()
model.fit(x_train, y_train)

train_score = model.score(x_train, y_train)
test_score = model.score(x_test, y_test)

if test_score > train_score:
print(f"Test Score: {test_score:.4f} | Train Score: {train_score:.4f} | Random State: {i}")

Jupyter

230701521 ex 12 Last Checkpoint: 25 minutes ago

File Edit View Run Kernel Settings Help

Trusted

+ ↻ 📄 🔍 ⏮️ ⏪ ⏩ ⏭ Code ▾

JupyterLab Python 3 (ipykernel)

Test Score: 0.8750 | Train Score: 0.8344 | Random State: 18
Test Score: 0.8500 | Train Score: 0.8438 | Random State: 19
Test Score: 0.8750 | Train Score: 0.8438 | Random State: 20
Test Score: 0.8625 | Train Score: 0.8344 | Random State: 21
Test Score: 0.8750 | Train Score: 0.8406 | Random State: 22
Test Score: 0.8750 | Train Score: 0.8406 | Random State: 24
Test Score: 0.8500 | Train Score: 0.8344 | Random State: 26
Test Score: 0.8500 | Train Score: 0.8406 | Random State: 27
Test Score: 0.8625 | Train Score: 0.8344 | Random State: 30

[8]: x_train,x_test,y_train,y_test=train_test_split(features,label,test_size=0.2,random_states=209)
finalModel=LogisticRegression()
finalModel.fit(x_train,y_train)

[8]: LogisticRegression
LogisticRegression()

[9]: print(finalModel.score(x_train,y_train))
print(finalModel.score(x_test,y_train))

0.85
0.85

[10]: from sklearn.metrics import classification_report
print(classification_report(label,finalModel.predict(features)))

precision recall f1-score support

0 0.86 0.91 0.89 257
1 0.83 0.73 0.77 143

accuracy 0.84 0.82 0.85 400
macro avg 0.84 0.82 0.83 400
weighted avg 0.85 0.85 0.85 400