



**RAJALAKSHMI ENGINEERING COLLEGE**

*Approved by AICTE | Affiliated to Anna University | Accredited by NAAC*

Department of Computer Science and Engineering

CS23334 Fundamentals of Data Science Lab

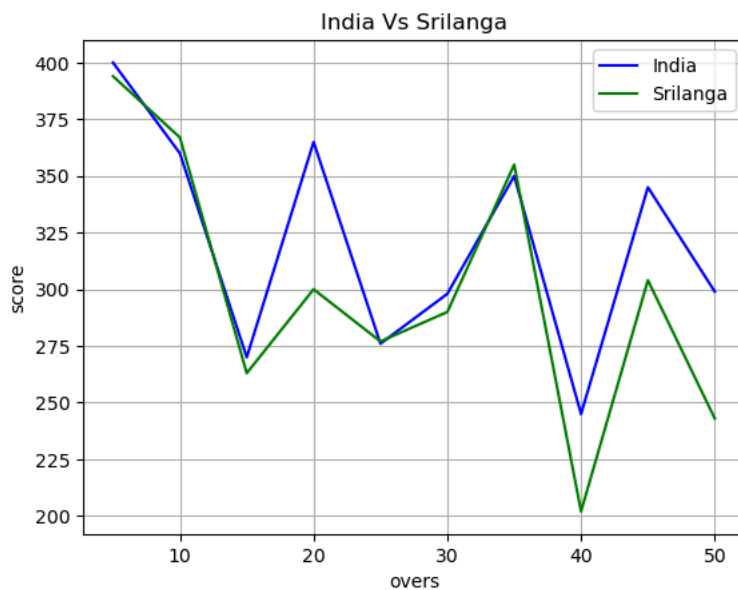
III semester II Year (2023R)

Name of the Student : VISHAAL S

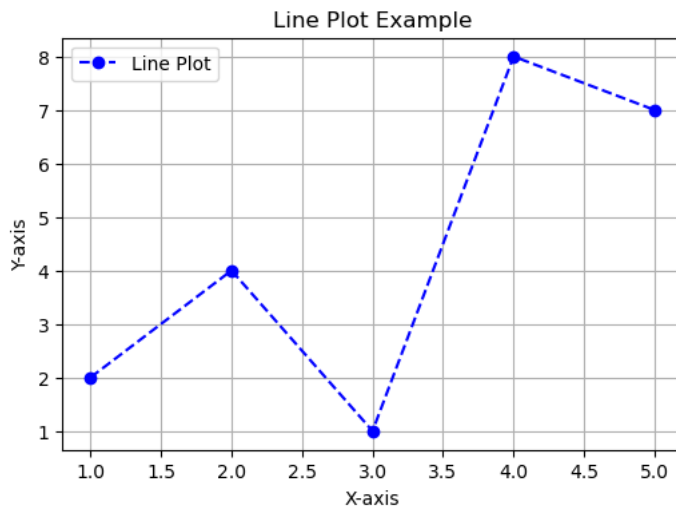
Register Number : 240701594

```
#Vishaal S
#240701594
#Fundamentals of Data Science
#17.07.2025
#LinePlot,Bargraph,Piechart,Histogram and Scatter plot using matplotlib
```

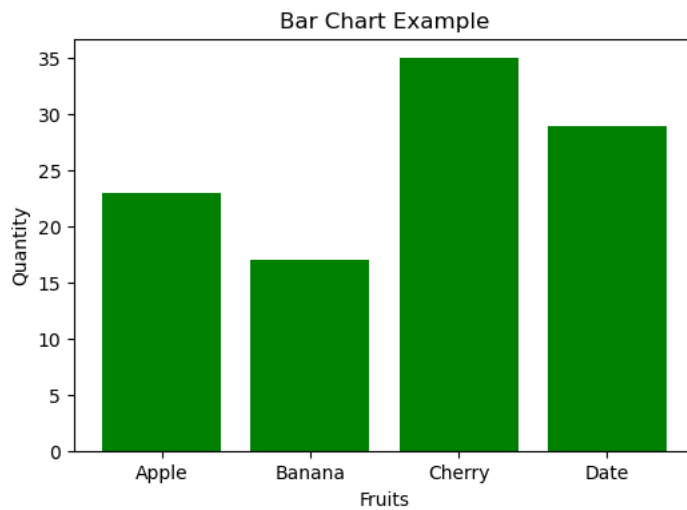
```
import matplotlib.pyplot as plt
overs = list(range(5, 51, 5))
India_Score = [400, 360, 270, 365, 276, 298, 350, 245, 345, 299]
Srilanga_Score = [394, 367, 263, 300, 277, 290, 355, 202, 304, 243]
plt.plot(overs, India_Score, color="blue", label="India")
plt.plot(overs, Srilanga_Score, color="green", label="Srilanga")
plt.title("India Vs Srilanga")
plt.xlabel("overs")
plt.ylabel("score")
plt.legend()
plt.grid(True)
plt.show()
```



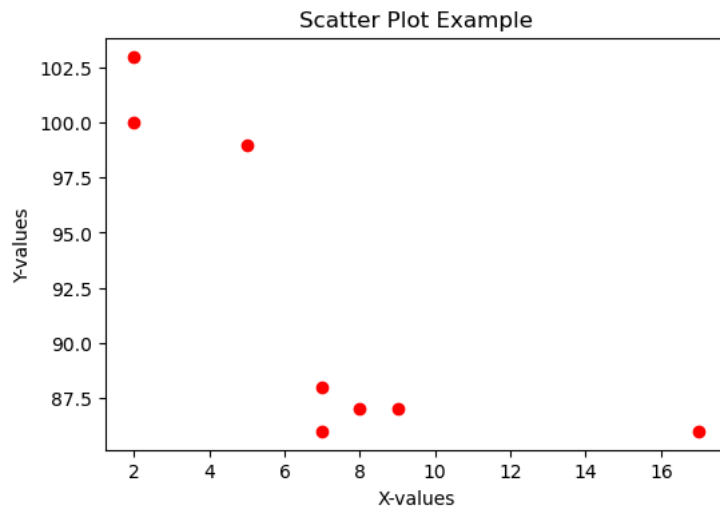
```
import matplotlib.pyplot as plt
x = [1, 2, 3, 4, 5]
y = [2, 4, 1, 8, 7]
plt.figure(figsize=(6, 4))
plt.plot(x, y, color='blue', marker='o', linestyle='--', label='Line Plot')
plt.title("Line Plot Example")
plt.xlabel("X-axis")
plt.ylabel("Y-axis")
plt.legend()
plt.grid(True)
plt.show()
```



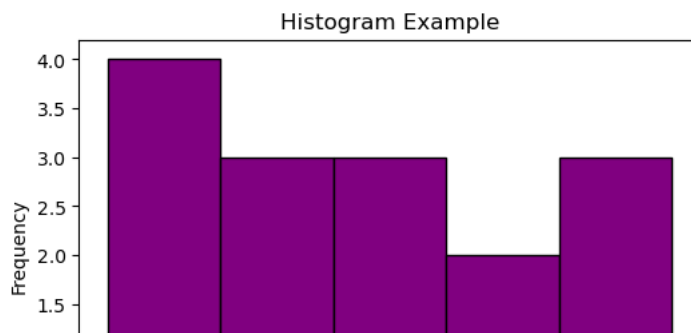
```
import matplotlib.pyplot as plt
categories = ['Apple', 'Banana', 'Cherry', 'Date']
values = [23, 17, 35, 29]
plt.figure(figsize=(6, 4))
plt.bar(categories, values, color='green')
plt.title("Bar Chart Example")
plt.xlabel("Fruits")
plt.ylabel("Quantity")
plt.show()
```



```
x_scatter = [5, 7, 8, 7, 2, 17, 2, 9]
y_scatter = [99, 86, 87, 88, 100, 86, 103, 87]
plt.figure(figsize=(6, 4))
plt.scatter(x_scatter, y_scatter, color='red')
plt.title("Scatter Plot Example")
plt.xlabel("X-values")
plt.ylabel("Y-values")
plt.show()
```



```
data = [22, 87, 5, 43, 56, 73, 55, 54, 11, 20, 51, 5, 79, 31, 27]
plt.figure(figsize=(6, 4))
plt.hist(data, bins=5, color='purple', edgecolor='black')
plt.title("Histogram Example")
plt.xlabel("Value Range")
plt.ylabel("Frequency")
plt.show()
```



```
# Vishaal S
# 240701594
# 24.7.25
# Data preprocessing
```

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns # Added this import
file_path='sales_data.csv'
df = pd.read_csv(file_path)
df
```

	Date	Product	Sales	Quantity	Region	
0	01-01-2023	Product A	200	4	North	
1	02-01-2023	Product B	150	3	South	
2	03-01-2023	Product A	220	5	North	
3	04-01-2023	Product C	300	6	East	
4	05-01-2023	Product B	180	4	West	
5	06-01-2023	Product A	210	5	North	
6	07-01-2023	Product C	320	7	East	
7	08-01-2023	Product B	160	3	South	
8	09-01-2023	Product A	230	6	North	
9	10-01-2023	Product C	310	7	East	
10	11-01-2023	Product B	190	4	West	
11	12-01-2023	Product A	240	6	North	
12	13-01-2023	Product C	330	8	East	
13	14-01-2023	Product B	170	3	South	
14	15-01-2023	Product A	250	7	North	
15	16-01-2023	Product C	340	8	East	

Next steps: [Generate code with df](#) [New interactive sheet](#)

```
df['Sales'] = df['Sales'].fillna(df['Sales'].mean())
df = df.dropna(subset=['Product', 'Quantity', 'Region'])
# Convert 'Date' column to datetime objects
df['Date'] = pd.to_datetime(df['Date'], dayfirst=True, errors='coerce')
```

```
df.describe()
```

	Sales	Quantity	
count	16.000000	16.000000	
mean	237.500000	5.375000	
std	64.031242	1.746425	
min	150.000000	3.000000	
25%	187.500000	4.000000	
50%	225.000000	5.500000	
75%	302.500000	7.000000	
max	340.000000	8.000000	

```
product_summary = df.groupby('Product').agg({
'Sales': 'sum',
'Quantity': 'sum'
```

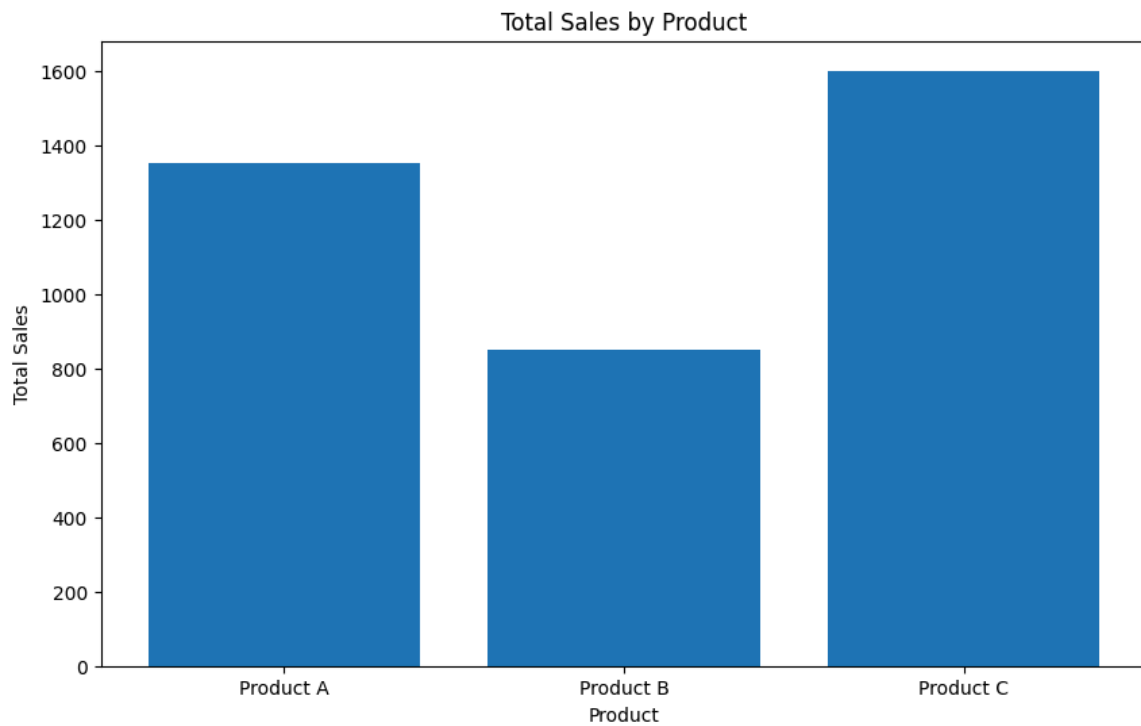
```
).reset_index()  
product_summary
```

	Product	Sales	Quantity	
0	Product A	1350	33	
1	Product B	850	17	
2	Product C	1600	36	

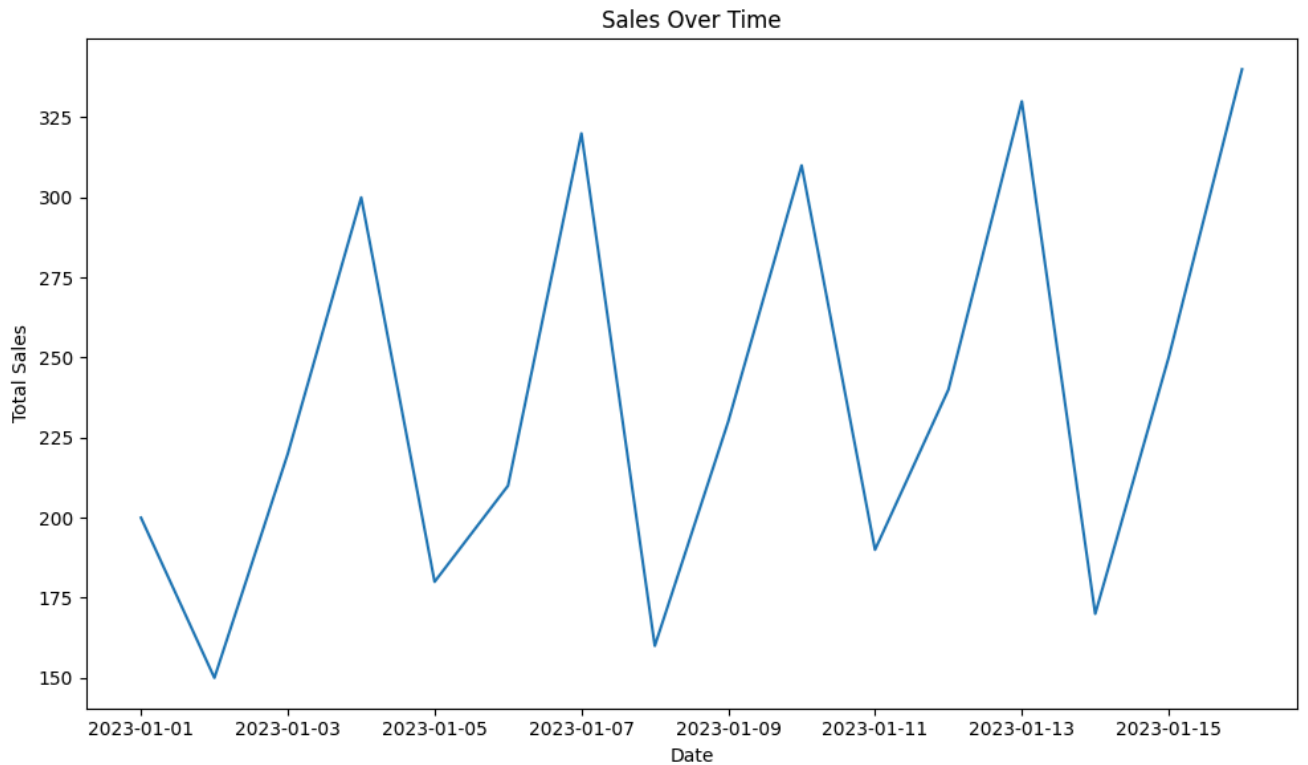
Next steps: [Generate code with product\\_summary](#) [New interactive sheet](#)

```
plt.figure(figsize=(10, 6))  
plt.bar(product_summary['Product'], product_summary['Sales'])  
plt.xlabel('Product')  
plt.ylabel('Total Sales')  
plt.title('Total Sales by Product')
```

Text(0.5, 1.0, 'Total Sales by Product')



```
sales_over_time = df.groupby('Date').agg({'Sales': 'sum'}).reset_index()  
plt.figure(figsize=(10, 6))  
plt.plot(sales_over_time['Date'], sales_over_time['Sales'])  
plt.xlabel('Date')  
plt.ylabel('Total Sales')  
plt.title('Sales Over Time')  
plt.tight_layout()  
plt.show()
```

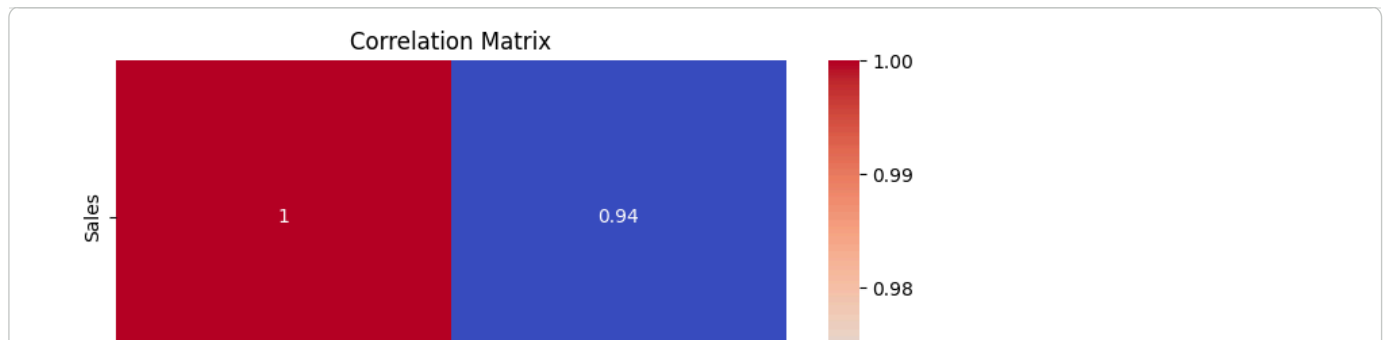


```
pivot_table = df.pivot_table(values='Sales', index='Region', columns='Product',aggfunc="sum", fill_value=0)
pivot_table
```

Product	Product A	Product B	Product C
Region			
East	0	0	1600
North	1350	0	0
South	0	480	0
West	0	370	0

Next steps: [Generate code with pivot\\_table](#) [New interactive sheet](#)

```
correlation_matrix = df.corr(numeric_only=True)
import seaborn as sns
plt.figure(figsize=(8, 6))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
plt.title('Correlation Matrix')
plt.show()
```





```
# Vishaal S
# 240701594
# 7.8.25
# Data preprocessing
```

```
import numpy as np
import pandas as pd
df=pd.read_csv("pre_process_datasample.csv")
df
```

	Country	Age	Salary	Purchased
0	France	44.0	72000.0	No
1	Spain	27.0	48000.0	Yes
2	Germany	30.0	54000.0	No
3	Spain	38.0	61000.0	No
4	Germany	40.0	NaN	Yes
5	France	35.0	58000.0	Yes
6	Spain	NaN	52000.0	No
7	France	48.0	79000.0	Yes
8	Germany	50.0	83000.0	No
9	France	37.0	67000.0	Yes

```
df.info()
df.Country.mode()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Country     10 non-null     object
1   Age         9 non-null      float64
2   Salary      9 non-null      float64
3   Purchased   10 non-null     object
dtypes: float64(2), object(2)
memory usage: 452.0+ bytes
0   France
Name: Country, dtype: object
```

```
df.Country.mode()[0]
```

```
'France'
```

```
type(df.Country.mode())
```

```
pandas.core.series.Series
```

```
df['Country'] = df['Country'].fillna(df['Country'].mode()[0])
df['Age'] = df['Age'].fillna(df['Age'].median())
df['Salary'] = df['Salary'].fillna(round(df['Salary'].mean()))
pd.get_dummies(df.Country)
```

	France	Germany	Spain
0	True	False	False
1	False	False	True
2	False	True	False
3	False	False	True
4	False	True	False
5	True	False	False
6	False	False	True
7	True	False	False
8	False	True	False
9	True	False	False

```
updated_dataset=pd.concat([pd.get_dummies(df.Country),df.iloc[:,[1,2,3]],axis=1)
updated_dataset
```

	France	Germany	Spain	Age	Salary	Purchased
0	True	False	False	44.0	72000.0	No
1	False	False	True	27.0	48000.0	Yes
2	False	True	False	30.0	54000.0	No
3	False	False	True	38.0	61000.0	No
4	False	True	False	40.0	63778.0	Yes
5	True	False	False	35.0	58000.0	Yes
6	False	False	True	38.0	52000.0	No
7	True	False	False	48.0	79000.0	Yes
8	False	True	False	50.0	83000.0	No
9	True	False	False	37.0	67000.0	Yes

```
updated_dataset.Purchased.replace(['No','Yes'],[0,1])
```

C:\Users\HP\AppData\Local\Temp\ipykernel\_8624\826484493.py:1: FutureWarning: Downcasting behavior in `replace` is deprecated and will be removed in a future version. To suppress this warning, please call `updated\_dataset.Purchased.replace(['No','Yes'],[0,1])`.

```
0    0
```

```
1    1
```

```
2    0
```

```
3    0
```

```
4    1
```

```
5    1
```

```
6    0
```

```
7    1
```

```
8    0
```

```
9    1
```

```
Name: Purchased, dtype: int64
```

Start coding or [generate](#) with AI.

```
# Vishaal S
# 240701594
# 7.8.25
# Data preprocessing
```

```
import numpy as np
import pandas as pd
df=pd.read_csv("Hotel_Dataset.csv")
df.duplicated()
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: 924.0+ bytes
```

```
df.drop_duplicates(inplace=True)
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
10	10	30-35	5	RedFox	non-Veg	-6755	4	87777	30-35

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

	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

```
df.loc[df["CustomerID"] < 0, "CustomerID"] = np.nan
df.loc[df["Bill"] < 0, "Bill"] = np.nan
df.loc[df["EstimatedSalary"] < 0, "EstimatedSalary"] = np.nan
df
```

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

```
df.loc[~df['NoOfPax'].between(1, 20), 'NoOfPax'] = np.nan
df
```

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

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

```
<bound method Series.unique of 0          veg
1      Non-Veg
2      Veg
3      Veg
4      Vegetarian
5      Non-Veg
6      Vegetarian
7      Veg
8      Non-Veg
9      Non-Veg
10     non-Veg
Name: FoodPreference, dtype: object>
```

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

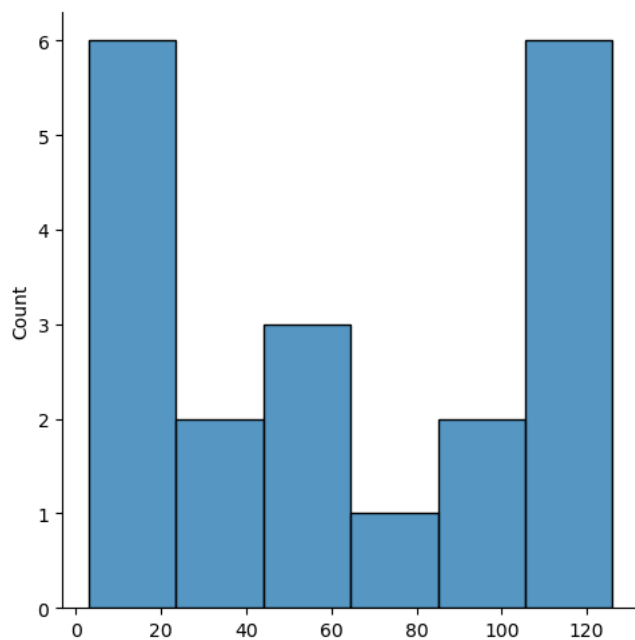
	CustomerID	Age_Group	Rating(1-5)	Hotel	FoodPreference	Bill	NoOfPax	EstimatedSalary	Age_Group.1
0	1.0	20-25	4	Ibis	veg	1300.0	2.0	40000.0	20-25
1	2.0	30-35	5	LemonTree	Non-Veg	2000.0	3.0	59000.0	30-35
2	3.0	25-30	6	RedFox	Veg	1322.0	2.0	30000.0	25-30
3	4.0	20-25	-1	LemonTree	Veg	1234.0	2.0	120000.0	20-25
4	5.0	35+	3	Ibis	Vegetarian	989.0	2.0	45000.0	35+
5	6.0	35+	3	Ibis	Non-Veg	1909.0	2.0	122220.0	35+
6	7.0	35+	4	RedFox	Vegetarian	1000.0	2.0	21122.0	35+
7	8.0	20-25	7	LemonTree	Veg	2999.0	2.0	345673.0	20-25
8	9.0	25-30	2	Ibis	Non-Veg	3456.0	3.0	96755.0	25-30
9	9.0	25-30	2	Ibis	Non-Veg	3456.0	3.0	96755.0	25-30
10	10.0	30-35	5	RedFox	non-Veg	1966.0	4.0	87777.0	30-35

Start coding or [generate](#) with AI.

```
# Vishaal S  
# 240701594  
# 14.8.25  
# Outliers
```

```
import numpy as np  
array=np.random.randint(1,150,20)  
def outDetection(array):  
    sorted(array)  
    Q1,Q3=np.percentile(array,[25,75])  
    IQR=Q3-Q1  
    lr=Q1-(1.5*IQR)  
    ur=Q3+(1.5*IQR)  
    return lr,ur  
lr,ur=outDetection(array)
```

```
import seaborn as sns  
import matplotlib.pyplot as plt  
%matplotlib inline  
sns.displot(array)  
plt.show()
```



```
sns.distplot(array)  
plt.show()
```

```
C:\Users\HP\AppData\Local\Temp\ipykernel_13032\2106234926.py:1: UserWarning:
```

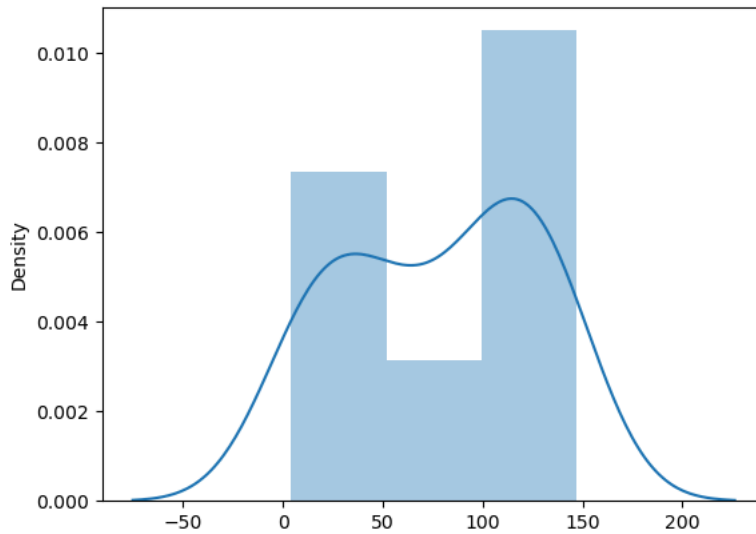
```
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
```

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

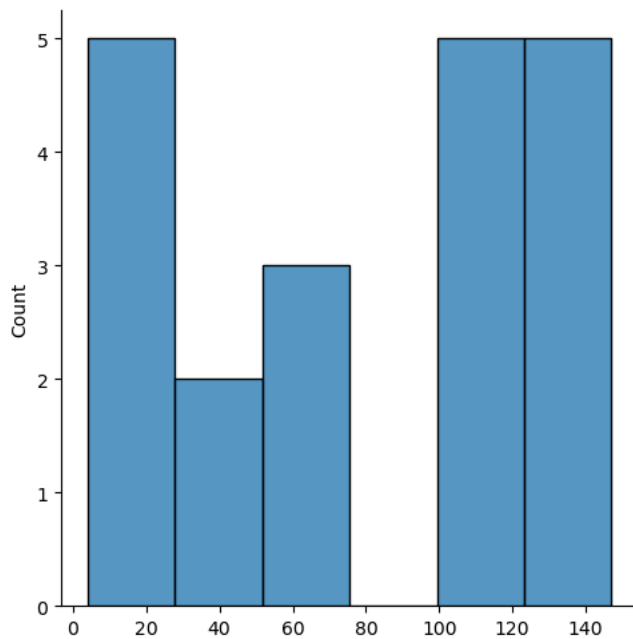
For a guide to updating your code to use the new functions, please see

<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(array)
```



```
new_array=array[(array>10) & (array<150)]  
sns.distplot(new_array);  
plt.show()
```



```
sns.distplot(new_array);  
plt.show()
```

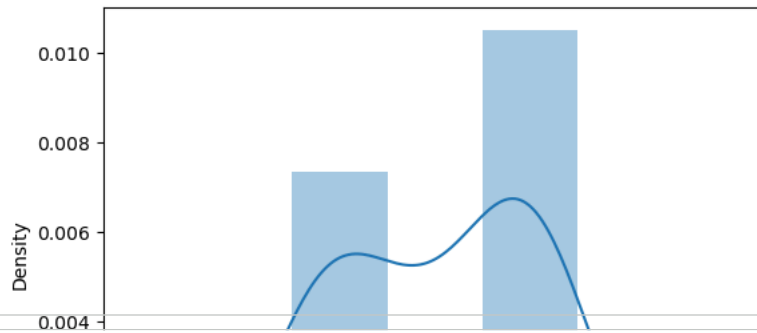
```
C:\Users\HP\AppData\Local\Temp\ipykernel_13032\1814760654.py:1: UserWarning:
```

```
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
```

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(new_array);
```



Start coding or [generate](#) with AI.



```
# Vishaal S
# 240701594
# 21.8.25
# Feature SCaling
```

```
import numpy as np
import pandas as pd
df=pd.read_csv('pre_process_datasample.csv')
df
```

	Country	Age	Salary	Purchased
0	France	44.0	72000.0	No
1	Spain	27.0	48000.0	Yes
2	Germany	30.0	54000.0	No
3	Spain	38.0	61000.0	No
4	Germany	40.0	NaN	Yes
5	France	35.0	58000.0	Yes
6	Spain	NaN	52000.0	No
7	France	48.0	79000.0	Yes
8	Germany	50.0	83000.0	No
9	France	37.0	67000.0	Yes

```
df.Country.fillna(df.Country.mode()[0])
features=df.iloc[:, :-1].values
label=df.iloc[:, -1].values
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]])
```

▼ SimpleImputer ⓘ ?

SimpleImputer()

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

▼ SimpleImputer ⓘ ?

SimpleImputer()

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

```
array([[ 'France', 44.0, 72000.0],
       [ 'Spain', 27.0, 48000.0],
       [ 'Germany', 30.0, 54000.0],
       [ 'Spain', 38.0, 61000.0],
       [ 'Germany', 40.0, 63777.77777777778],
       [ 'France', 35.0, 58000.0],
       [ 'Spain', 38.77777777777778, 52000.0],
       [ 'France', 48.0, 79000.0],
       [ 'Germany', 50.0, 83000.0],
       [ 'France', 37.0, 67000.0]], dtype=object)
```

```
from sklearn.preprocessing import OneHotEncoder
oh = OneHotEncoder(sparse_output=False)
Country=oh.fit_transform(features[:, [0]])
Country
```

```
array([[1., 0., 0.],
       [0., 0., 1.],
       [0., 1., 0.],
       [0., 0., 1.]])
```

```
[0., 1., 0.],  
[1., 0., 0.],  
[0., 0., 1.],  
[1., 0., 0.],  
[0., 1., 0.],  
[1., 0., 0.]])
```

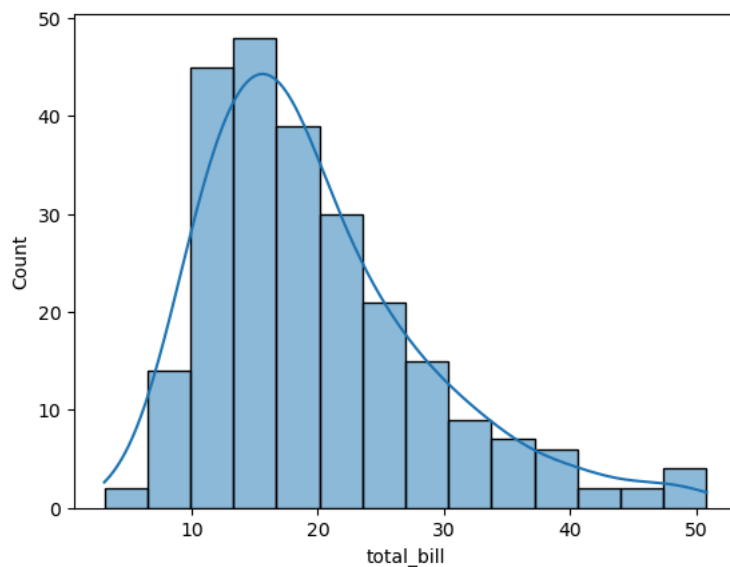
Start coding or [generate](#) with AI.

```
# Vishaal S  
# 240701594  
# 28.8.25  
# EDA
```

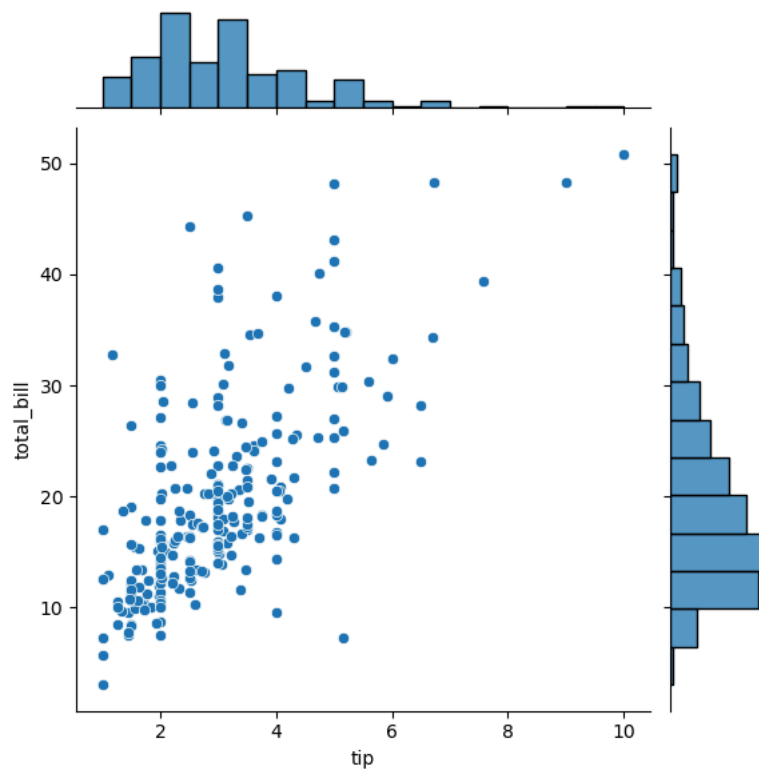
```
import seaborn as sns  
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
%matplotlib inline  
tips=sns.load_dataset('tips')  
tips.head()
```

	total_bill	tip	sex	smoker	day	time	size
0	16.99	1.01	Female	No	Sun	Dinner	2
1	10.34	1.66	Male	No	Sun	Dinner	3
2	21.01	3.50	Male	No	Sun	Dinner	3
3	23.68	3.31	Male	No	Sun	Dinner	2
4	24.59	3.61	Female	No	Sun	Dinner	4

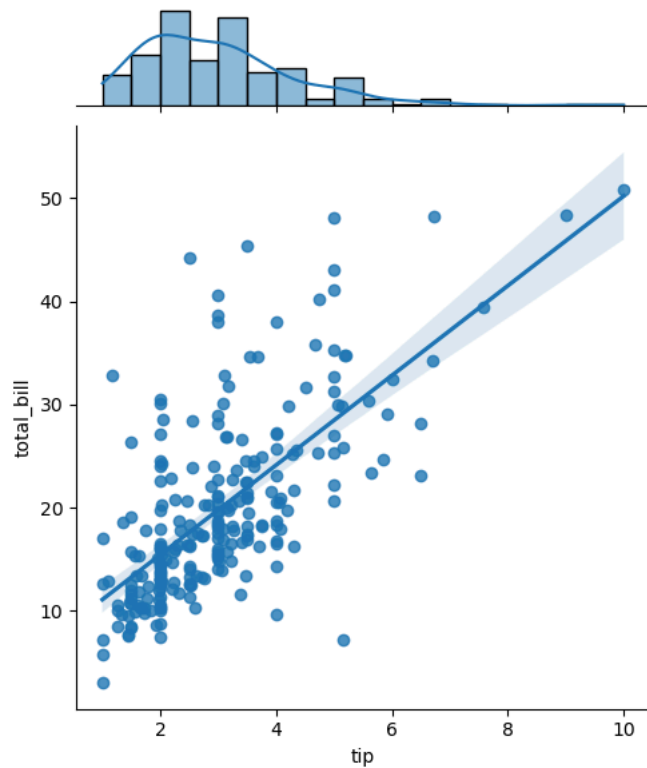
```
sns.histplot(tips.total_bill,kde=True)  
plt.show()
```



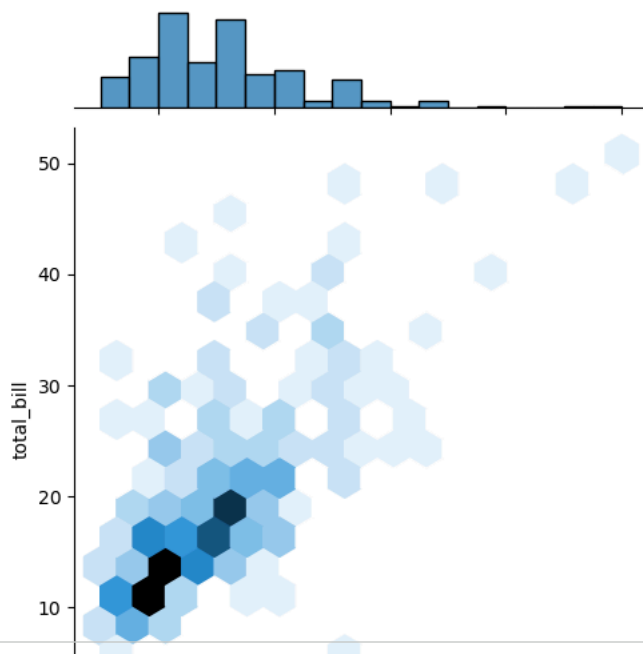
```
sns.jointplot(x=tips.tip,y=tips.total_bill)  
plt.show()
```



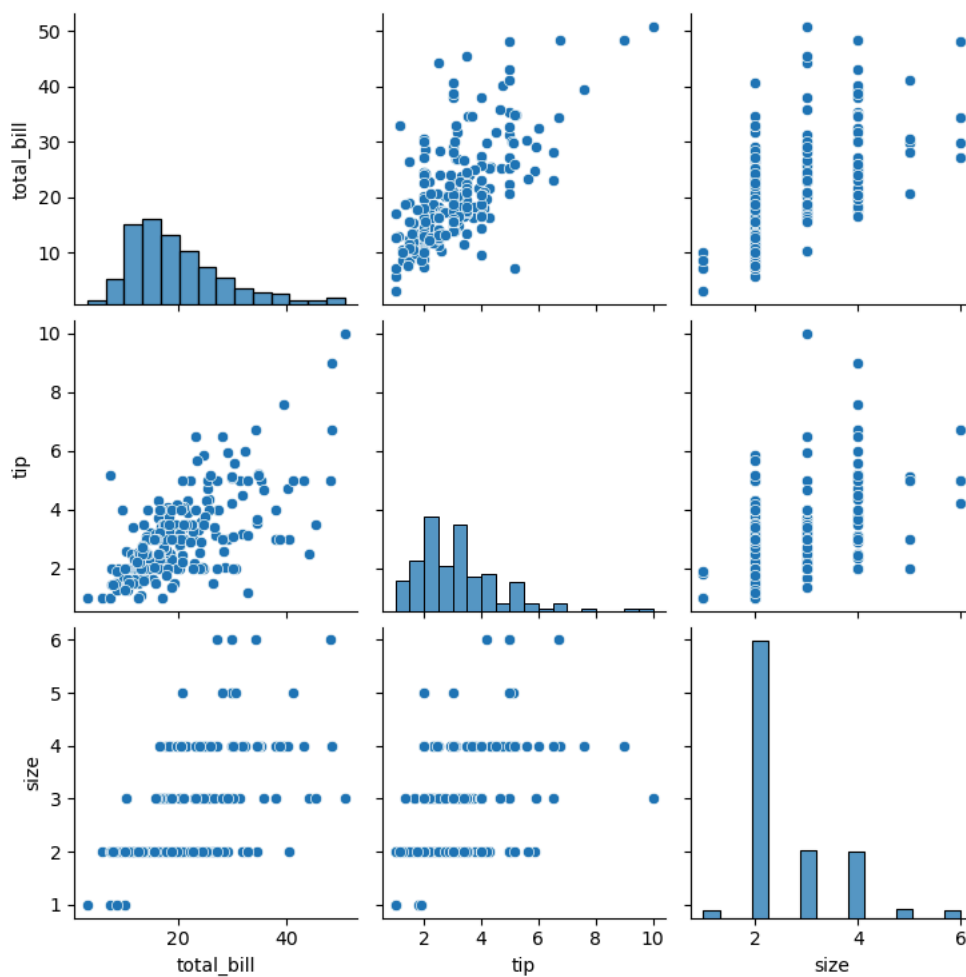
```
sns.jointplot(x=tips.tip,y=tips.total_bill,kind="reg")  
plt.show()
```



```
sns.jointplot(x=tips.tip,y=tips.total_bill,kind="hex")  
plt.show()
```



```
sns.pairplot(tips)
plt.show()
```



Start coding or [generate](#) with AI.

```
# Vishaal S
# 240701594
# 17.9.25
# Linear Regression
```

```
import numpy as np
import pandas as pd
df=pd.read_csv('Salary_data.csv')
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: 612.0 bytes
```

```
df.dropna(inplace=True)
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: 612.0 bytes
```

```
df.describe()
```

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

```
features=df.iloc[:,[0]].values
label=df.iloc[:,[1]].values
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=25)
from sklearn.linear_model import LinearRegression
model=LinearRegression()
model.fit(x_train,y_train)
```

```
LinearRegression
```

```
model.score(x_train,y_train)
```

```
0.9577907749872991
```

```
model.score(x_test,y_test)
```

```
0.9531732818427658
```

```
model.coef_
```

```
array([[9339.90339715]])
```

```
model.intercept_
```

```
array([26561.50676243])
```

```
import pickle
pickle.dump(model,open('SalaryPred.model','wb'))
model=pickle.load(open('SalaryPred.model','rb'))
yr_of_exp=float(input("Enter Years of Experience: "))
y=yr_of_exp
yr_of_exp_NP=np.array([[yr_of_exp]])
Salary=model.predict(yr_of_exp_NP)

print("Estimated salary for {} year of exp is{} ".format(yr_of_exp,Salary))
```

Start coding or [generate](#) with AI.

```
# Vishaal S  
# 240701594  
# 17.9.25  
# Logistic Regression
```

```
import numpy as np  
import pandas as pd  
df=pd.read_csv('Social_Network_Ads.csv')  
df
```

	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
398	15755018	Male	36	33000	0
399	15594041	Female	49	36000	1

400 rows × 5 columns

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



```
[ 49, 28000],
[ 57, 33000],
[ 56, 60000],
[ 49, 39000],
[ 39, 71000],
[ 47, 34000],
[ 48, 35000],
[ 48, 33000],
[ 47, 23000],
[ 45, 45000],
[ 60, 42000],
[ 39, 59000],
[ 46, 41000],
[ 51, 23000],
[ 50, 20000],
[ 36, 33000],
[ 49, 36000]])
```

label

```
array([0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 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, 0, 0, 0, 1, 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, 1, 0, 0,
       0, 0, 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, 0, 0, 0, 0, 0, 0, 0, 0, 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, 1, 0, 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, 0, 0, 1, 0, 0, 0, 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, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1,
       0, 0, 0, 0, 1, 0, 0, 1, 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, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1,
       0, 1, 0, 1, 1, 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, 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, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1,
       0, 1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1,
       1, 1, 0, 1])
```

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
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)

x_train,x_test,y_train,y_test=train_test_split(features,label,test_size=0.2,random_state=25)
finalModel=LogisticRegression()
finalModel.fit(x_train,y_train)
```

▼ LogisticRegression ⓘ ?  
LogisticRegression()

```
print(finalModel.score(x_train,y_train))
print(finalModel.score(x_test,y_test))
```

```
0.846875
0.8
```

```
from sklearn.metrics import classification_report
print(classification_report(label,finalModel.predict(features)))
```

	precision	recall	f1-score	support
0	0.84	0.92	0.88	257
1	0.82	0.69	0.75	143
accuracy			0.84	400
macro avg	0.83	0.81	0.82	400
weighted avg	0.84	0.84	0.83	400

Start coding or [generate](#) with AI.

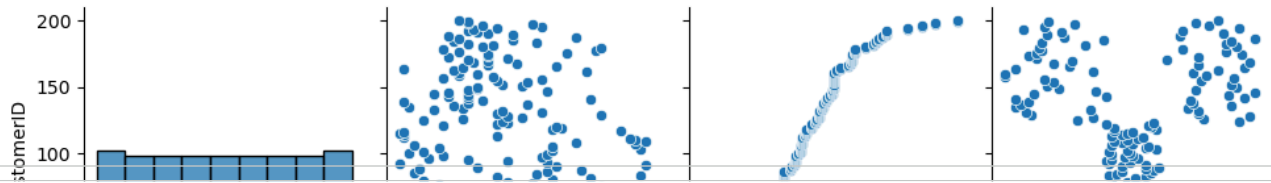


```
# Vishaal S  
# 240701594  
# 24.10.25  
# K means Clustering
```

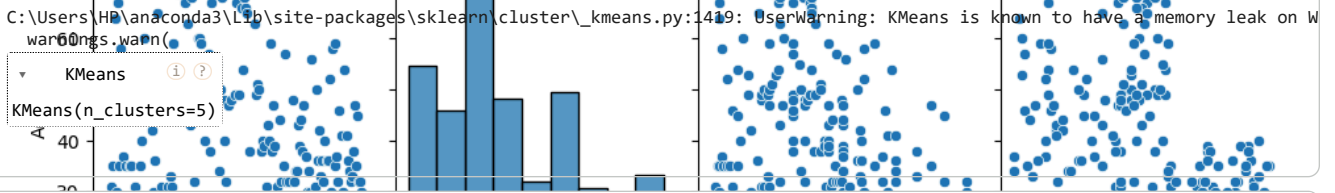
```
import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt  
import seaborn as sns  
%matplotlib inline  
df=pd.read_csv('Mall_Customers.csv')  
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 200 entries, 0 to 199  
Data columns (total 5 columns):  
#   Column                Non-Null Count  Dtype  
---  ---  
0   CustomerID            200 non-null   int64  
1   Gender                200 non-null   object  
2   Age                   200 non-null   int64  
3   Annual Income (k$)    200 non-null   int64  
4   Spending Score (1-100) 200 non-null   int64  
dtypes: int64(4), object(1)  
memory usage: 7.9+ KB
```

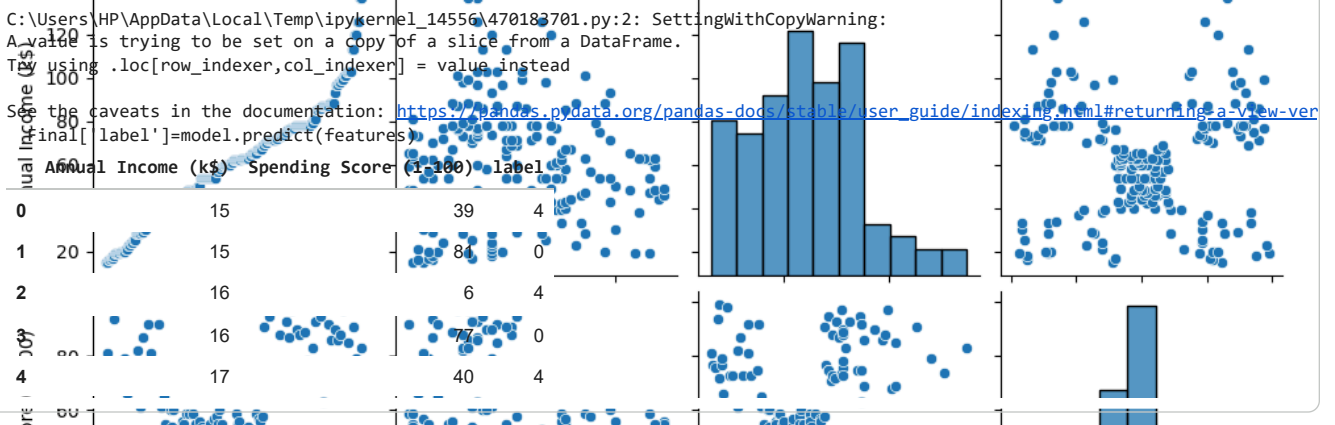
```
sns.pairplot(df)  
plt.show()
```



```
features=df.iloc[:,[3,4]].values
from sklearn.cluster import KMeans
model=KMeans(n_clusters=5)
model.fit(features)
KMeans(n_clusters=5)
```

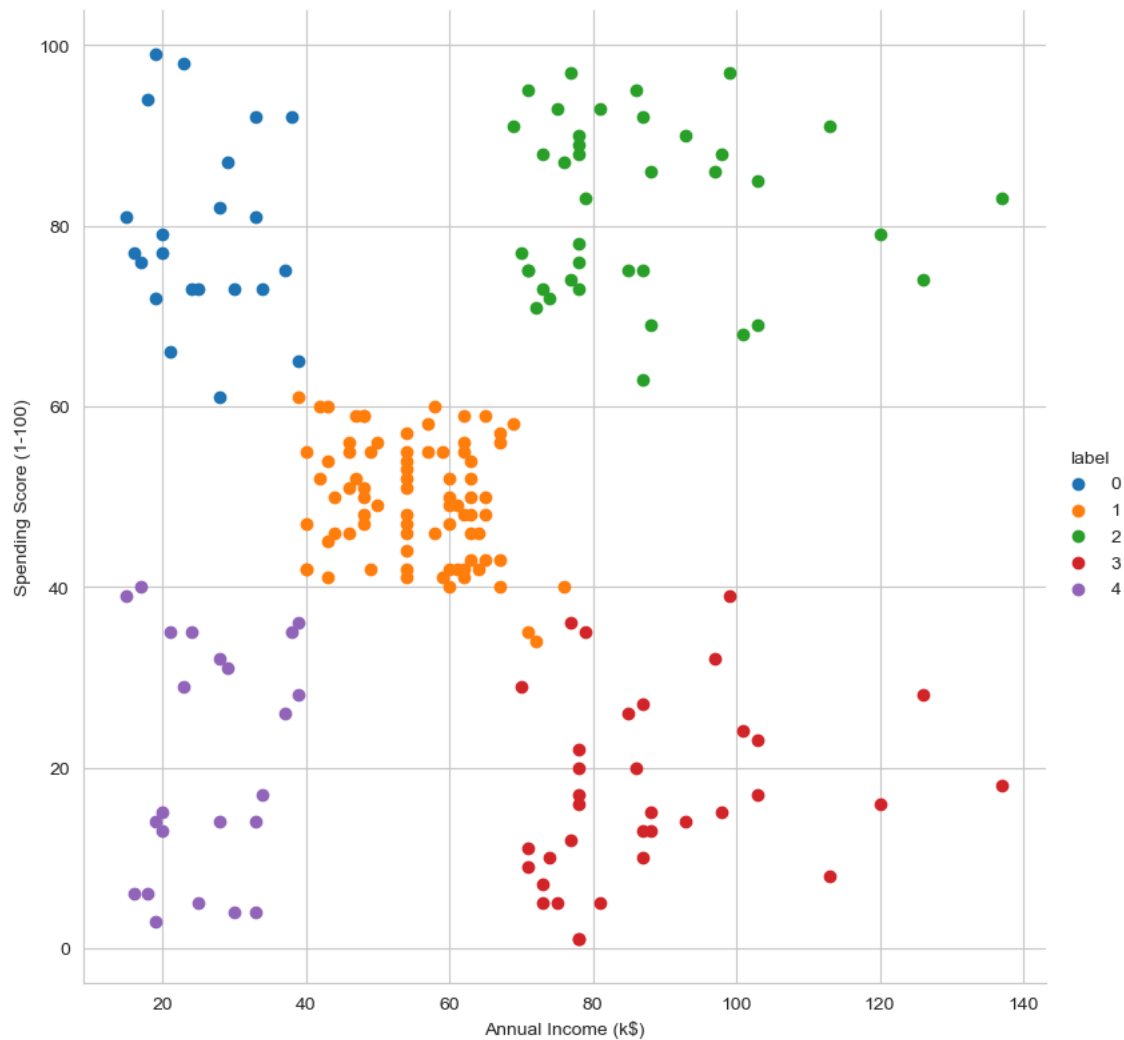


```
Final=df.iloc[:,[3,4]]
Final['label']=model.predict(features)
Final.head()
```



```
sns.set_style("whitegrid")
sns.FacetGrid(Final,hue="label",height=8) \
.map(plt.scatter,"Annual Income (k$)", "Spending Score (1-100)") \
.add_legend();
plt.show()
```

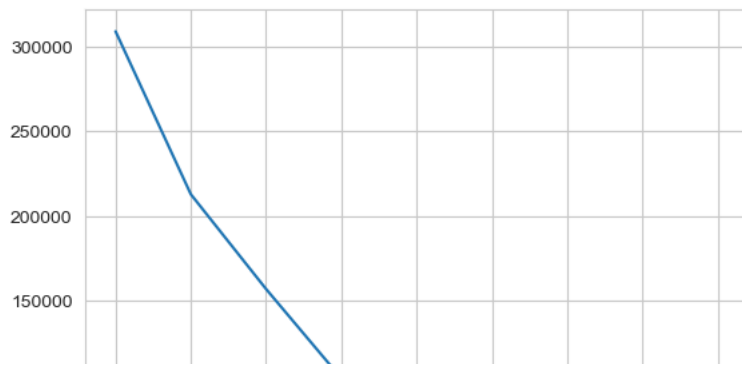




```
features_el=df.iloc[:,[2,3,4]].values
from sklearn.cluster import KMeans
wcss=[]
for i in range(1,10):
    model=KMeans(n_clusters=i)
    model.fit(features_el)
    wcss.append(model.inertia_)
plt.plot(range(1,10),wcss)
```

C:\Users\HP\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1419: UserWarning: KMeans is known to have a memory leak on W  
warnings.warn(  
C:\Users\HP\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1419: UserWarning: KMeans is known to have a memory leak on W  
warnings.warn(  
C:\Users\HP\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1419: UserWarning: KMeans is known to have a memory leak on W  
warnings.warn(  
C:\Users\HP\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1419: UserWarning: KMeans is known to have a memory leak on W  
warnings.warn(  
C:\Users\HP\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1419: UserWarning: KMeans is known to have a memory leak on W  
warnings.warn(  
C:\Users\HP\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1419: UserWarning: KMeans is known to have a memory leak on W  
warnings.warn(  
C:\Users\HP\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1419: UserWarning: KMeans is known to have a memory leak on W  
warnings.warn(  
C:\Users\HP\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1419: UserWarning: KMeans is known to have a memory leak on W  
warnings.warn(  
C:\Users\HP\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1419: UserWarning: KMeans is known to have a memory leak on W  
warnings.warn(  
[<matplotlib.lines.Line2D at 0x20725ab1310>]

```
plt.show()
```



```
# Vishaal S
# 240701594
# 1.10.25
# KNN
```

```
import numpy as np
import pandas as pd
df=pd.read_csv('Iris (1).csv')
df.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
```

```
df.variety.value_counts()
df
```

	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

Next steps: [Generate code with df](#) [New interactive sheet](#)

```
features = df.iloc[:, :-1].values
label = df.iloc[:, -1].values
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
xtrain, xtest, ytrain, ytest = train_test_split(
    features, label, test_size=0.2, random_state=20
)
model_KNN = KNeighborsClassifier(n_neighbors=5)
model_KNN.fit(xtrain, ytrain)
```

▼ KNeighborsClassifier ⓘ ?

KNeighborsClassifier()

```
print(model_KNN.score(xtrain,ytrain))
print(model_KNN.score(xtest,ytest))
```

```
0.9833333333333333
0.9666666666666667
```

```
from sklearn.metrics import confusion_matrix
confusion_matrix(label,model_KNN.predict(features))
```

```
array([[50,  0,  0],
       [ 0, 48,  2],
       [ 0,  1, 49]])
```

```
from sklearn.metrics import classification_report
print(classification_report(label,model_KNN.predict(features)))
```

	precision	recall	f1-score	support
Setosa	1.00	1.00	1.00	50
Versicolor	0.98	0.96	0.97	50
Virginica	0.96	0.98	0.97	50
accuracy			0.98	150
macro avg	0.98	0.98	0.98	150
weighted avg	0.98	0.98	0.98	150



```
#Vishaal S  
#240701594  
# 8.10.25  
#T-test
```

```
import numpy as np  
from scipy import stats  
  
marks = np.array([72, 68, 75, 70, 74, 69, 71, 73, 70, 72])  
  
mu_0 = 70
```

```
# Vishaal S  
# 240701594  
# 8.10.25  
# Z-test
```

```
import numpy as np  
from math import sqrt  
from scipy.stats import norm  
x_bar = 51.2 # sample mean  
mu_0 = 50 # population mean  
sigma = 3 # population standard deviation
```



```
#Vishaal S  
#240701594  
#Anova test
```

```
import numpy as np  
from scipy import stats  
A = [20, 22, 23]  
B = [19, 20, 18]  
C = [25, 27, 26]  
f_stat, p_value = stats.f_oneway(A, B, C)  
print(f"F-statistic: {f_stat:.3f}")
```