

❖ Day-8:-Task:-Data Preprocessing

❖ Data Preprocessing:

Data preprocessing is a crucial step in the data analysis and machine learning pipeline. It involves cleaning and transforming raw data into a format that is suitable for analysis or model training.

❖ Data processing working involves:

- I. Handling Missing Data: Identify and deal with missing values.
- II. Handling Duplicate Data: Detect and remove duplicate records.
- III. Handling outliers: Address outliers through methods like z-score or IQR.
- IV. Feature Scaling: standardize or normalize numerical features.
- V. Categorical Variable Encoding: convert categorical variable into numerical format.
- VI. Handling Text data: Transform text data for analysis or modeling.
- VII. Feature Engineering: create or transform features for improved performance.
- VIII. Data splitting: divide the dataset into training and testing sets.
- IX. Data Transformation: Apply necessary transformations for modeling.
- X. Handling imbalanced data: Address class imbalance if needed.
- XI. Data Visualization: visualize data for insights.

❖ Description:

I'm thrilled to apply this newfound knowledge to real-world scenarios, and the prospect of further honing these skills through the Skill Boost Internship Program adds an extra layer of excitement. Here's to the journey ahead and the exciting challenges awaiting in the Skill Boost Internship Program(www.Batweb.com).

day-8 data preprocessing - Jupyter Notebook

localhost:8890/notebooks/day-8%20data%20preprocessing.ipynb

File Edit View Insert Cell Kernel Widgets Help

Python 3 (pykernel)

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

In [4]: dataset=pd.read_csv('employee_data.csv')
dataset.head()
```

Out[4]:

	First Name	Last Name	Gender	Age	Salary	Expenditure	Savings	Expenditure Percentage	Savings Percentage
0	Aditi	Mishra	Female	35	INR 88,000.00	INR 65,000.00	INR 15,000.00	81.25%	18%
1	Sona	Singh	Female	35	INR 78,000.00	INR 58,000.00	INR 12,000.00	82.86%	17%
2	Devi	Thakur	Female	36	INR 65,000.00	INR 60,000.00	INR 5,000.00	92.31%	8%
3	Angal	Kaur	Female	42	INR 62,000.00	INR 58,000.00	INR 12,000.00	88.65%	19%
4	Sahil	Deshmukh	Male	29	INR 62,000.00	INR 57,000.00	INR 5,000.00	91.94%	8%

```
In [8]: dataset.tail()
```

Out[8]:

	First Name	Last Name	Gender	Age	Salary	Expenditure	Savings	Expenditure Percentage	Savings Percentage
30	Shaan	Pandey	Male	37	INR 25,200.00	INR 22,000.00	INR 3,200.00	87.30%	13%
31	Navi	Gupta	Female	51	INR 25,000.00	INR 20,000.00	INR 5,000.00	80.00%	20%
32	Ambar	Khatri	Male	36	INR 25,000.00	INR 22,000.00	INR 3,000.00	88.00%	12%
33	Lata	Agarwal	Female	28	INR 25,000.00	INR 24,000.00	INR 1,000.00	96.00%	4%
34	Anaya	Agarwal	Female	25	INR 24,100.00	INR 21,000.00	INR 3,100.00	87.14%	13%

```
In [9]: dataset.describe()
```

Out[9]:

	Age
count	35.000000
mean	37.628571
std	8.106228
min	25.000000
25%	30.500000
50%	36.000000
75%	42.000000
max	52.000000

```
In [ ]: dataset.isnull().sum()
```

day-8 data preprocessing - Jupyter Notebook

localhost:8891/notebooks/day-8%20data%20preprocessing.ipynb

File Edit View Insert Cell Kernel Widgets Help

Python 3 (pykernel)

```
#Identify numeric columns
numeric_columns = dataset.select_dtypes(include=np.number).columns.tolist()
numeric_columns

Out[6]: ['Age']

In [7]: #filling values with Mean
dataset[numeric_columns] = dataset[numeric_columns].fillna(dataset[numeric_columns].mean())

In [8]: #fill values with Mode (Most Frequent)
#Identify categorical columns
categorical_columns = ['Gender', 'Salary', 'Expenditure']
categorical_columns

Out[8]: ['Gender', 'Salary', 'Expenditure']

In [10]: ##filling values with Mode
dataset[categorical_columns] = dataset[categorical_columns].fillna(dataset[categorical_columns].mode().iloc[0])

In [12]: #Drop non required rows
#Drop rows for NaN values in 'Smoking'
dataset.dropna(subset=['Savings Percentage'], axis=0, inplace=True)
dataset

Out[12]:
```

	First Name	Last Name	Gender	Age	Salary	Expenditure	Savings	Expenditure Percentage	Savings Percentage
0	Aditi	Mishra	Female	35	INR 88,000.00	INR 65,000.00	INR 15,000.00	81.25%	18%
1	Sona	Singh	Female	35	INR 78,000.00	INR 58,000.00	INR 12,000.00	82.86%	17%
2	Devi	Thakur	Female	36	INR 65,000.00	INR 60,000.00	INR 5,000.00	92.31%	8%
3	Angal	Kaur	Female	42	INR 62,000.00	INR 58,000.00	INR 12,000.00	88.65%	19%
4	Sahil	Deshmukh	Male	29	INR 62,000.00	INR 57,000.00	INR 5,000.00	91.94%	8%
5	Adhira	Das	Male	42	INR 42,000.00	INR 59,000.00	INR 3,000.00	95.16%	5%
6	Dhaya	Bhatt	Female	42	INR 52,500.00	INR 47,000.00	INR 5,500.00	89.52%	10%

The image shows a Jupyter Notebook interface with a light gray theme. At the top, there's a browser address bar showing 'localhost:8891/notebooks/day-8%20data%20preprocessing.ipynb'. The notebook title is 'jupyter day-8 data preprocessing' with a subtitle 'Last Checkpoint 6 minutes ago (unsaved changes)'. Below the title bar, there's a menu bar with 'File', 'Edit', 'View', 'Insert', 'Cell', 'Kernel', 'Widgets', and 'Help'. To the right of the menu bar are buttons for 'Not Trusted' and 'Python 3 (pykernel)'. The main area of the notebook contains three code cells. The first cell has a comment '#data encoding - Handle/encode categorical data' and imports 'ColumnTransformer' and 'OneHotEncoder'. The second cell has a comment '#label encoding' and imports 'LabelEncoder'. The third cell has a comment '#split the dataset for training and testing' and imports 'train_test_split'. The output of the third cell shows the first few rows of the dataset, including columns for name, gender, age, income, and a target variable. The notebook is running on a local host, and the system tray at the bottom shows the time as 12:55 AM on 12/23/2023.

day-8 data preprocessing - Jupyter

```
['Anjali', 'Kaur', 'Female', 42, 'INR 62,000.00', 'INR 50,000.00', 'INR 12,000.00', '80.65%', 'INR 20,000.00', 'INR 5,500.00', '78.43%']
['Lata', 'Agarwal', 'Female', 28, 'INR 25,000.00', 'INR 24,000.00', 'INR 1,000.00', '96.00%', 'INR 1,000.00']
['Nikhil', 'Shah', 'Male', 30, 'INR 30,000.00', 'INR 28,000.00', 'INR 2,000.00', '93.33%', 'INR 2,000.00']
['Shaan', 'Sharma', 'Male', 35, 'INR 35,000.00', 'INR 25,000.00', 'INR 10,000.00', '71.43%', 'INR 10,000.00']
['Darsh', 'Kulkarni', 'Male', 41, 'INR 30,000.00', 'INR 28,000.00', 'INR 2,000.00', '93.33%', 'INR 2,000.00']
['Inaya', 'Kumari', 'Female', 52, 'INR 30,000.00', 'INR 20,000.00', 'INR 10,000.00', '66.67%', 'INR 10,000.00']
[ 2  0  1  10  10  15  15  0  8  4  3  1  4  0  9  2  3  5  0  6  1  17  16  11
  2  7]
[ 2  3  6  10  16  18  12  10  14]
```

In [23]: #Reset Index

```
dataset.reset_index(drop=True, inplace=True)
dataset
```

Out[23]:

	First Name	Last Name	Gender	Age	Salary	Expenditure	Savings	Expenditure Percentage	Savings Percentage
0	Anjali	Kaur	Female	35	INR 60,000.00	INR 65,000.00	INR 15,000.00	81.25%	18%
1	Sona	Singh	Female	35	INR 70,000.00	INR 58,000.00	INR 12,000.00	82.86%	17%
2	Devi	Thakur	Female	36	INR 65,000.00	INR 60,000.00	INR 5,000.00	92.31%	8%
3	Anjali	Kaur	Female	42	INR 62,000.00	INR 50,000.00	INR 12,000.00	80.65%	19%
4	Sahil	Deshmukh	Male	29	INR 42,000.00	INR 57,000.00	INR 5,000.00	91.94%	8%
5	Aditya	Das	Male	42	INR 42,000.00	INR 59,000.00	INR 3,000.00	95.19%	5%
6	Divya	Bhatt	Female	42	INR 52,500.00	INR 47,000.00	INR 5,500.00	89.52%	10%
7	Amar	Rao	Male	52	INR 52,100.00	INR 46,000.00	INR 6,100.00	88.29%	12%
8	Jaya	Mehta	Female	28	INR 52,000.00	INR 46,000.00	INR 6,000.00	88.46%	12%
9	Candice	Patel	Female	35	INR 52,000.00	INR 40,000.00	INR 12,000.00	76.92%	23%
10	Sudhir	Mishra	Male	35	INR 52,000.00	INR 47,000.00	INR 5,000.00	90.38%	10%
11	Raghar	Verma	Male	28	INR 51,000.00	INR 49,000.00	INR 2,000.00	96.08%	4%
12	Amitabh	Joshi	Male	36	INR 51,000.00	INR 50,000.00	INR 1,000.00	98.04%	2%
13	Karish	Patel	Male	45	INR 50,500.00	INR 40,000.00	INR 10,500.00	79.21%	21%
14	Tulsi	Acharya	Female	40	INR 45,200.00	INR 40,000.00	INR 5,200.00	88.50%	12%
15	Jasleen	Kaur	Female	50	INR 47,000.00	INR 40,000.00	INR 7,000.00	85.11%	11%
16	Nit	Kumar	Male	44	INR 47,000.00	INR 50,000.00	INR 4,000.00	96.40%	10%

day-8 data preprocessing - Jupyter

```
29 Agastya Joshi Male 50 INR 25,500.00 INR 20,000.00 INR 5,500.00 78.43% 22%
30 Shaan Pandey Male 37 INR 25,200.00 INR 22,000.00 INR 3,200.00 87.30% 13%
31 Navi Gupta Female 51 INR 25,000.00 INR 20,000.00 INR 5,000.00 80.00% 20%
32 Ambar Khatri Male 36 INR 25,000.00 INR 22,000.00 INR 3,000.00 88.00% 12%
33 Lata Agarwal Female 28 INR 25,000.00 INR 24,000.00 INR 1,000.00 96.00% 4%
34 Anaya Agarwal Female 25 INR 24,100.00 INR 21,000.00 INR 3,100.00 87.14% 13%
```

In [25]: #Check info and NULL again

```
dataset.shape
dataset.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 35 entries, 0 to 34
Data columns (total 9 columns):
 # Column          Non-Null Count  Dtype
---  --
 0 First Name       35 non-null    object
 1 Last Name        35 non-null    object
 2 Gender           35 non-null    object
 3 Age              35 non-null    int64
 4 Salary           35 non-null    object
 5 Expenditure       35 non-null    object
 6 Savings           35 non-null    object
 7 Expenditure Percentage 35 non-null    object
 8 Savings Percentage 35 non-null    object
dtypes: int64(1), object(8)
memory usage: 2.6+ KB
```

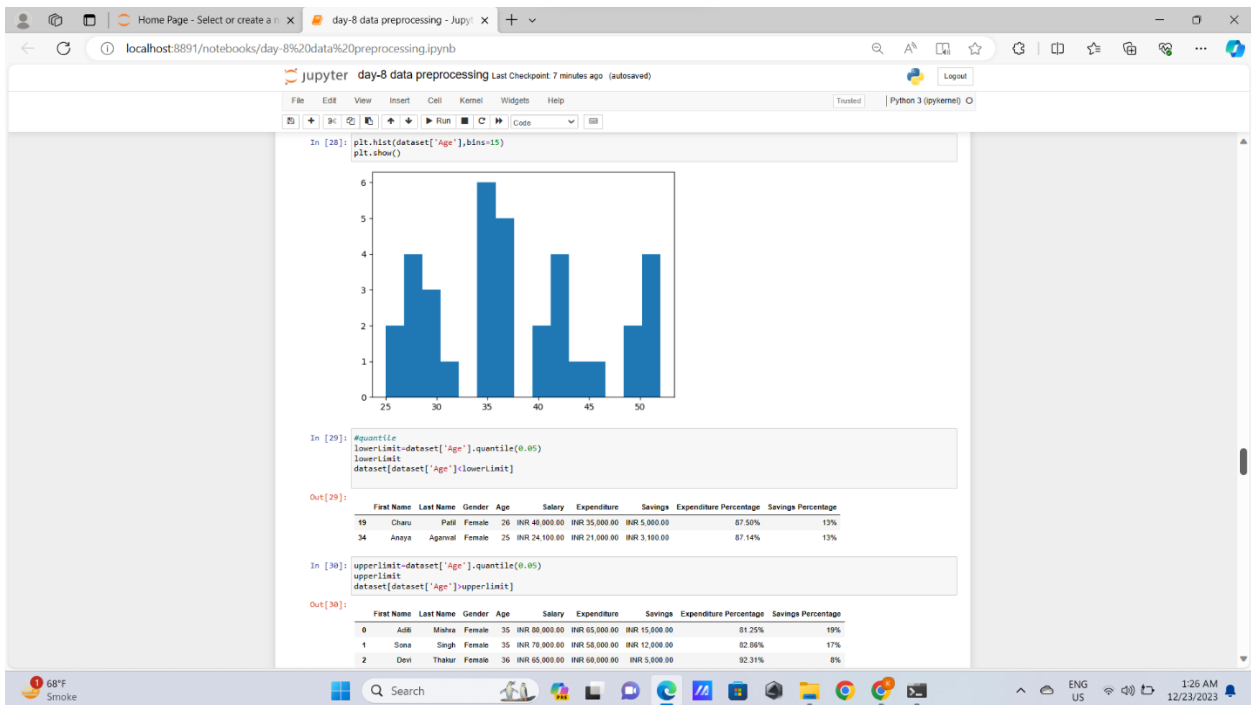
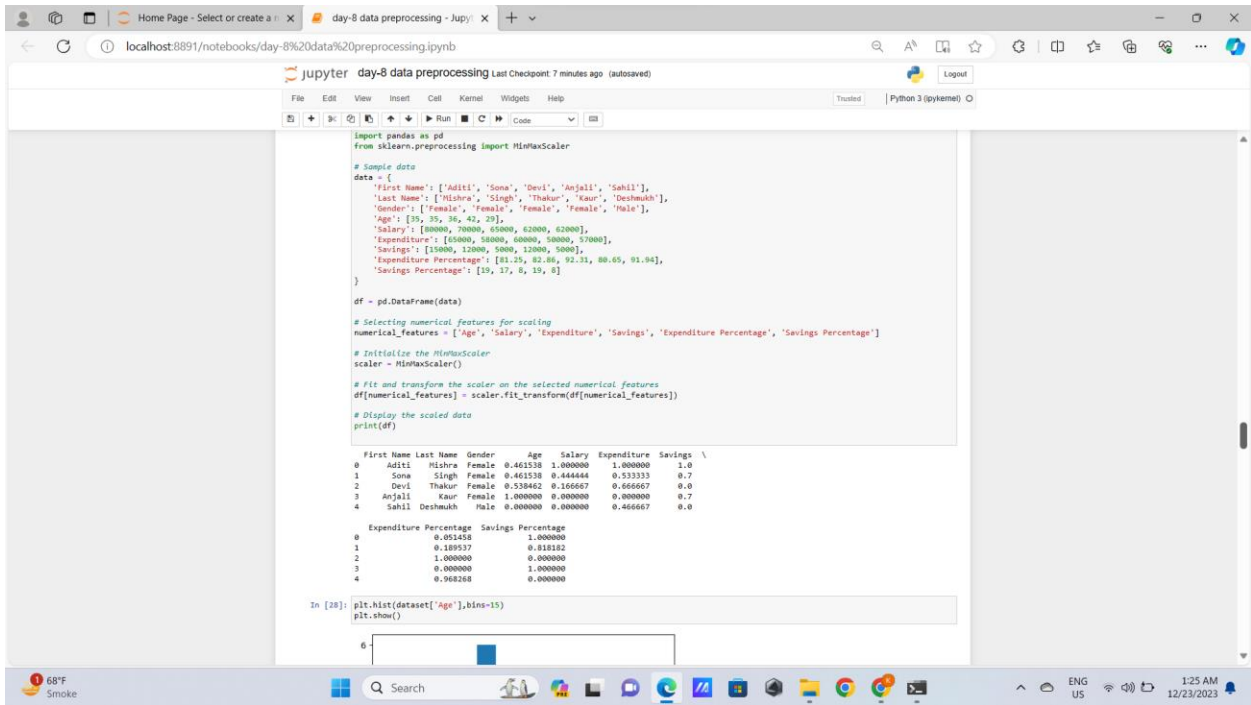
In [26]: dataset.isnull().sum()

Out[26]:

```
First Name      0
Last Name       0
Gender          0
Age             0
Salary          0
Expenditure     0
Savings         0
Expenditure Percentage 0
Savings Percentage 0
dtype: int64
```

In [27]: #Feature scaling

```
import pandas as pd
from sklearn.preprocessing importMinMaxScaler
```



day-8 data preprocessing - Jupyter

localhost:8891/notebooks/day-8%20data%20preprocessing.ipynb

jupyter day-8 data preprocessing Last Checkpoint: 7 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (pykernel)

4	Sahil	Deshmukh	Male	29	INR 42,000.00	INR 57,000.00	INR 5,000.00	91.94%	8%
5	Adithi	Das	Male	42	INR 42,000.00	INR 59,000.00	INR 3,000.00	95.16%	5%
6	Divya	Bhatt	Female	42	INR 52,500.00	INR 47,000.00	INR 5,500.00	89.52%	10%
7	Amar	Rao	Male	52	INR 52,100.00	INR 46,000.00	INR 6,100.00	88.29%	12%
8	Jaya	Mehra	Female	28	INR 52,000.00	INR 46,000.00	INR 6,000.00	88.46%	12%
9	Candly	Patel	Female	35	INR 52,000.00	INR 40,000.00	INR 12,000.00	76.92%	23%
10	Sudhar	Mishra	Male	35	INR 52,000.00	INR 47,000.00	INR 5,000.00	90.38%	10%
11	Raghar	Verma	Male	28	INR 51,000.00	INR 49,000.00	INR 2,000.00	96.00%	4%
12	Amalash	Joshi	Male	36	INR 51,000.00	INR 50,000.00	INR 1,000.00	96.04%	2%
13	Karak	Patel	Male	45	INR 50,500.00	INR 40,000.00	INR 10,500.00	79.21%	21%
14	Tulsi	Acharya	Female	40	INR 45,200.00	INR 40,000.00	INR 5,200.00	88.50%	12%
15	Jasleen	Kaur	Female	50	INR 45,000.00	INR 40,000.00	INR 5,000.00	88.89%	11%
16	Nili	Kumar	Male	44	INR 42,000.00	INR 59,000.00	INR 3,000.00	90.48%	10%
17	Jasand	Yadav	Male	25	INR 42,000.00	INR 38,000.00	INR 4,000.00	90.48%	10%
18	Atharva	Das	Male	29	INR 41,000.00	INR 35,000.00	INR 6,000.00	85.37%	15%
20	Ajay	Shah	Male	35	INR 36,000.00	INR 32,000.00	INR 4,000.00	88.89%	11%
21	Jai	Pati	Male	36	INR 35,500.00	INR 30,000.00	INR 5,500.00	84.51%	15%
22	Ishaan	Sharma	Male	35	INR 35,000.00	INR 25,000.00	INR 10,000.00	71.43%	29%
23	Dhruva	Yadav	Male	42	INR 35,000.00	INR 22,000.00	INR 13,000.00	62.86%	37%
24	Artha	Sharma	Male	31	INR 35,000.00	INR 31,000.00	INR 4,000.00	88.57%	11%
25	Shaila	Gupta	Female	51	INR 35,000.00	INR 34,000.00	INR 1,000.00	97.14%	3%
26	Itanya	Kumari	Female	52	INR 30,000.00	INR 20,000.00	INR 10,000.00	66.67%	33%
27	Darsh	Kulkarni	Male	41	INR 30,000.00	INR 28,000.00	INR 2,000.00	93.33%	7%
28	Mihir	Shah	Male	30	INR 30,000.00	INR 28,000.00	INR 2,000.00	93.33%	7%
29	Agastya	Joshi	Male	50	INR 25,500.00	INR 20,000.00	INR 5,500.00	78.43%	22%
30	Shaan	Pandey	Male	37	INR 25,200.00	INR 22,000.00	INR 3,200.00	87.30%	13%
31	Nani	Gupta	Female	51	INR 25,000.00	INR 20,000.00	INR 5,000.00	80.00%	20%
32	Ambar	Khatri	Male	36	INR 25,000.00	INR 22,000.00	INR 3,000.00	88.00%	12%
33	Lata	Agarwal	Female	28	INR 25,000.00	INR 24,000.00	INR 1,000.00	96.00%	4%

```
In [31]: dataset=dataset[(dataset['Age']>lowerlimit)& (dataset['Age']<upperlimit)]
dataset
Out[31]:
```

First Name	Last Name	Gender	Age	Salary	Expenditure	Savings	Expenditure Percentage	Savings Percentage
------------	-----------	--------	-----	--------	-------------	---------	------------------------	--------------------

```
In [32]:
```

day-8 data preprocessing - Jupyter

localhost:8891/notebooks/day-8%20data%20preprocessing.ipynb

jupyter day-8 data preprocessing Last Checkpoint: 7 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 (pykernel)

```
Out[31]:
```

First Name	Last Name	Gender	Age	Salary	Expenditure	Savings	Expenditure Percentage	Savings Percentage
------------	-----------	--------	-----	--------	-------------	---------	------------------------	--------------------

```
In [32]:
```

```
# Sample data
names = ['Aditi', 'Sona', 'Devi', 'Anjali', 'Sahil']
expenditure_percentages = [81.25, 82.86, 92.31, 89.65, 91.94]

# Create a pie chart
plt.pie(expenditure_percentages, labels=names, autopct='%1.1f%%', startangle=90)

# Add a title
plt.title('Expenditure Percentage Distribution')

# Display the pie chart
plt.show()
```

Expenditure Percentage Distribution

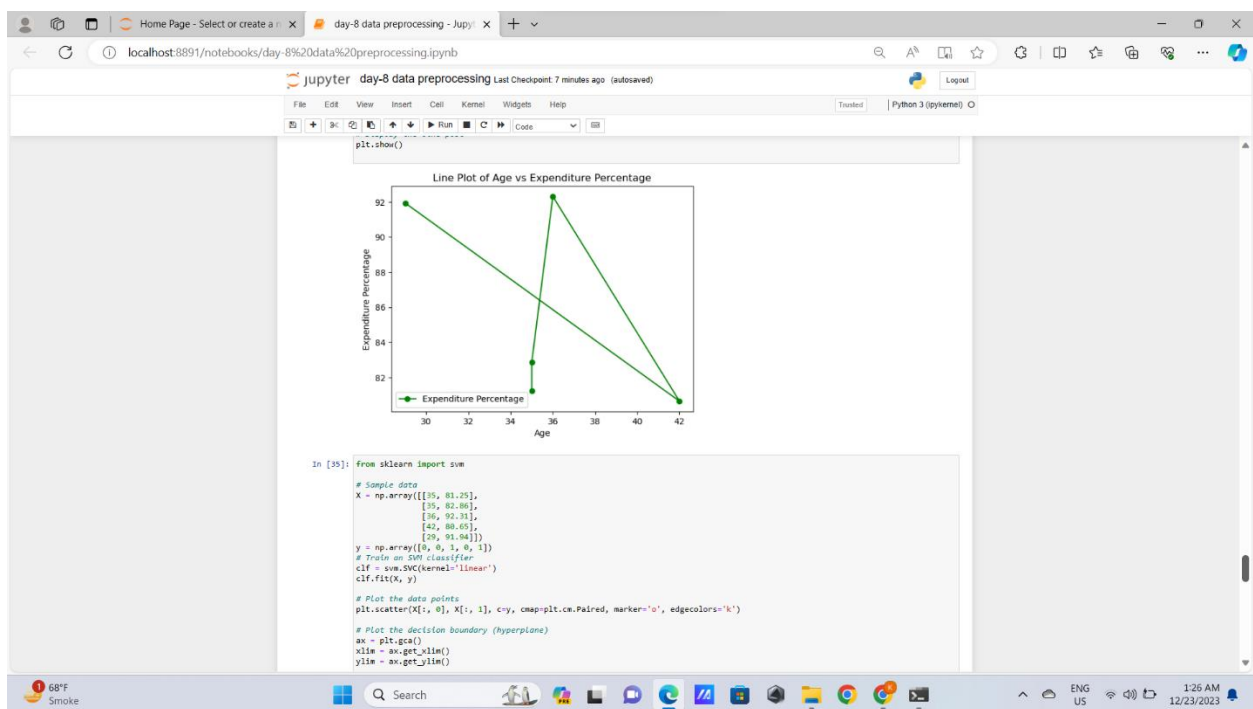
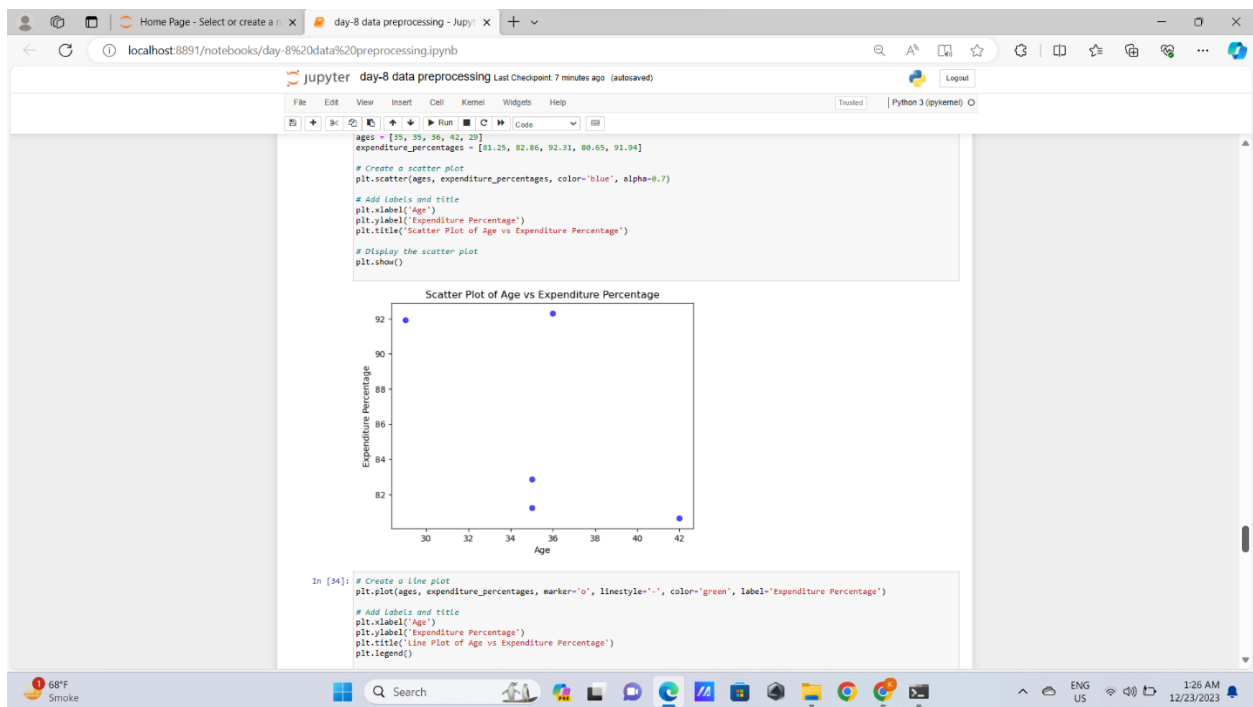
First Name	Last Name	Gender	Age	Salary	Expenditure	Savings	Expenditure Percentage	Savings Percentage
Aditi							18.9%	
Sona							19.3%	
Devi							21.5%	
Anjali							18.8%	
Sahil							21.4%	

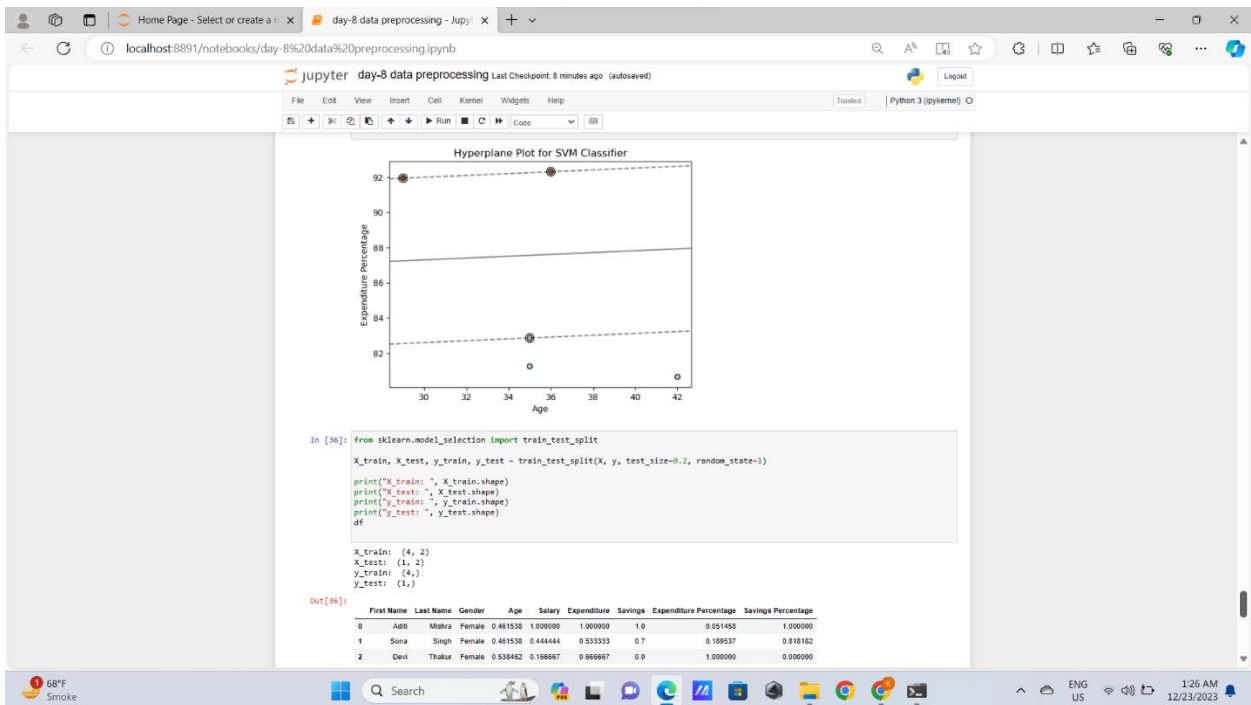
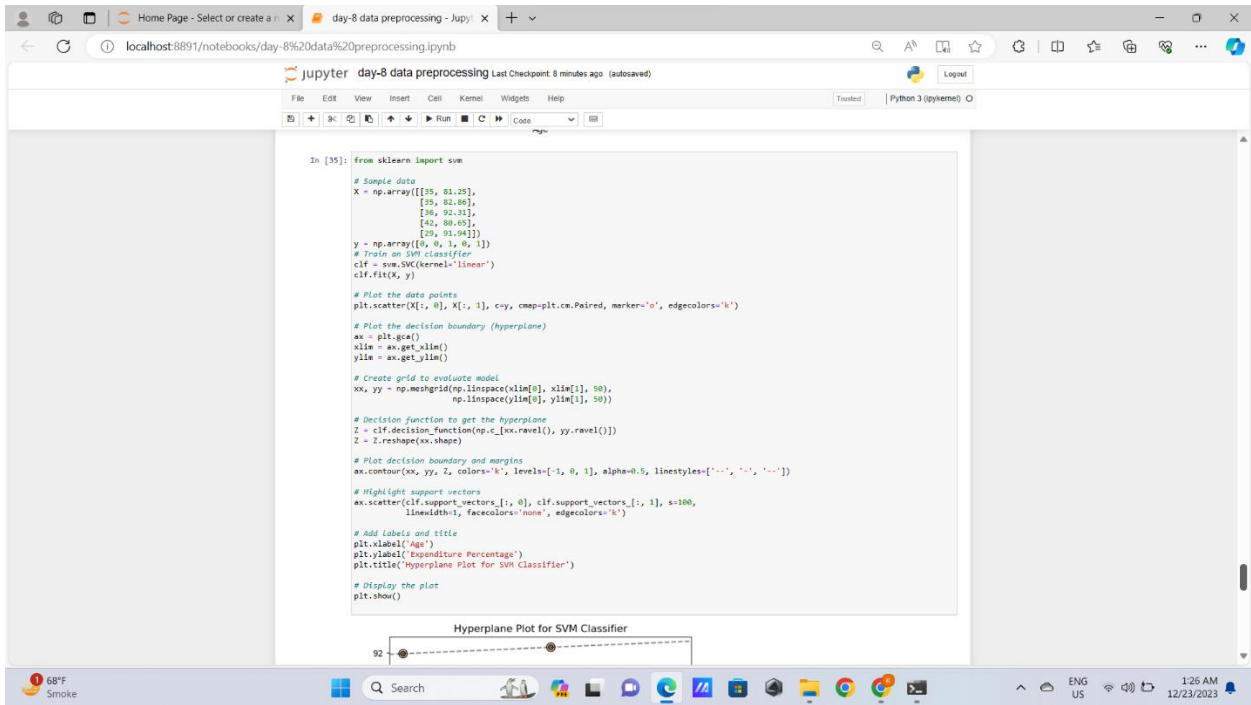
```
In [33]:
```

```
# Sample data
ages = [35, 35, 36, 42, 29]
expenditure_percentages = [81.25, 82.86, 92.31, 89.65, 91.94]

# Create a scatter plot
plt.scatter(ages, expenditure_percentages, color='blue', alpha=0.7)

# Add labels and title
```





day-8 data preprocessing - Jupyter

localhost:8891/notebooks/day-8%20data%20preprocessing.ipynb

Out[36]:

	First Name	Last Name	Gender	Age	Salary	Expenditure	Savings	Expenditure Percentage	Savings Percentage
0	Aditi	Mishra	Female	0.461538	1.000000	1.000000	1.0	0.051458	1.000000
1	Sona	Singh	Female	0.461538	0.444444	0.533333	0.7	0.109537	0.818182
2	Devi	Thakur	Female	0.538462	0.166667	0.666667	0.0	1.000000	0.000000
3	Anjali	Kaur	Female	1.000000	0.000000	0.000000	0.7	0.000000	1.000000
4	Sahil	Deshmukh	Male	0.000000	0.000000	0.466667	0.0	0.962268	0.000000

In [42]:

```
from sklearn.svm import SVC
model_linear = SVC(kernel='linear')
model_linear.fit(X_train[:500], y_train[:500])
```

Out[42]:

```
SVC
SVC(kernel='linear')
```

In [43]:

```
from sklearn import metrics
# predict
y_pred = model_linear.predict(X_test)
# accuracy
print("accuracy:", metrics.accuracy_score(y_test, y_pred), "\n")
```

accuracy: 1.0

In [44]:

```
#Add Dummies to code
#Creating Dummies
convert_to_dummies = ['Salary', 'Age', 'Salary']
field_dummies = pd.get_dummies(df[convert_to_dummies])
field_dummies
dataset = pd.concat([df, field_dummies], axis = 1)
dataset.drop(convert_to_dummies, axis=1, inplace=True)
dataset
```

Out[44]:

	First Name	Last Name	Gender	Expenditure	Savings	Expenditure Percentage	Savings Percentage
0	Aditi	Mishra	Female	1.000000	1.0	0.051458	1.000000
1	Sona	Singh	Female	0.533333	0.7	0.109537	0.818182

day-8 data preprocessing - Jupyter

localhost:8891/notebooks/day-8%20data%20preprocessing.ipynb

Out[42]:

```
SVC
SVC(kernel='linear')
```

In [43]:

```
from sklearn import metrics
# predict
y_pred = model_linear.predict(X_test)
# accuracy
print("accuracy:", metrics.accuracy_score(y_test, y_pred), "\n")
```

accuracy: 1.0

In [44]:

```
#Add Dummies to code
#Creating Dummies
convert_to_dummies = ['Salary', 'Age', 'Salary']
field_dummies = pd.get_dummies(df[convert_to_dummies])
field_dummies
dataset = pd.concat([df, field_dummies], axis = 1)
dataset.drop(convert_to_dummies, axis=1, inplace=True)
dataset
```

Out[44]:

	First Name	Last Name	Gender	Expenditure	Savings	Expenditure Percentage	Savings Percentage
0	Aditi	Mishra	Female	1.000000	1.0	0.051458	1.000000
1	Sona	Singh	Female	0.533333	0.7	0.109537	0.818182
2	Devi	Thakur	Female	0.666667	0.0	1.000000	0.000000
3	Anjali	Kaur	Female	0.000000	0.7	0.000000	1.000000
4	Sahil	Deshmukh	Male	0.466667	0.0	0.962268	0.000000

In []:

In []: