

Phase 3

Development part 1

This section will involve loading and pre-processing the dataset as you start to construct your project.

Import datasets and the necessary libraries:

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.preprocessing import StandardScaler, LabelEncoder
4
5 data = pd.read_csv('F:\Electricity\Electricity1.csv')
```

Exploring data analysis:

Investigate the dataset to learn about its properties and structure.

```
8 print("Original Dataset:")
9 print(data.head())
```

Handling the missing data:

When dealing with missing data, data scientists can use two primary methods to solve the error: Imputation or the removal

of data. The imputation method develops reasonable guesses for missing data. It's most useful when the percentage of missing data is low

```
10  
11 data.dropna(inplace=True)  
12
```

Encoding Categorical Data:

Data Encoding is an important pre-processing step in Machine Learning. It refers to the process of converting categorical or textual data into numerical format, so that it can be used as input for algorithms to process.

```
13 label_encoders = {}  
14 categorical_columns = ['categorical_column1', 'categorical_column2']  
15  
16 for col in categorical_columns:  
17     label_encoders[col] = LabelEncoder()  
18     data[col] = label_encoders[col].fit_transform(data[col])
```

Splitting the dataset:

Data Splitting is when data is divided into two or more subsets. Typically, with a two-part split, one part is used to evaluate or

test the data and the other to train the model. Data splitting is an important aspect of data science, particularly for creating models based on data.

```
20 X = data.drop('target_column', axis=1)
21 y = data['target_column']
22
23 scaler = StandardScaler()
24 X = scaler.fit_transform(X)
25
26 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Feature Scaling:

Feature scaling is a method used to normalize the range of independent variables or features of data. In data processing, it is also known as data normalization and is generally performed during the data preprocessing step.

```
27
28 print("\nPreprocessed Dataset:")
29 print("X_train shape:", X_train.shape)
30 print("X_test shape:", X_test.shape)
31 print("y_train shape:", y_train.shape)
32 print("y_test shape:", y_test.shape)
33
```

Output:

The screenshot displays a Jupyter Notebook environment with two tabs: 'main.py' and 'Electricity.py'. The 'Electricity.py' tab is active, showing a Python script that imports pandas and sklearn, reads a CSV file, and prints the first five rows of the dataset. Below the script, the console output shows the first five rows of the dataset, which include columns for DateTime, Holiday, HolidayFlag, ActualWindProduction, SystemLoadEP2, and SMPEP2. The output is formatted as a table with 5 rows and 18 columns.

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.preprocessing import StandardScaler, LabelEncoder
4 data = pd.read_csv('Electricity.csv')
5 print("Original Dataset:")
6 print(data.head())
7 data.dropna(inplace=True)
8 label_encoders = {}
9 categorical_columns = ['categorical_column1', 'categorical_column2']
10
11 for col in categorical_columns:
12     label_encoders[col] = LabelEncoder()
```

Original Dataset:

	DateTime	Holiday	HolidayFlag	...	ActualWindProduction	SystemLoadEP2	SMPEP2
0	01/11/2011 00:00	NaN	0	...	356.00	3159.60	54.32
1	01/11/2011 00:30	NaN	0	...	317.00	2973.01	54.23
2	01/11/2011 01:00	NaN	0	...	311.00	2834.00	54.23
3	01/11/2011 01:30	NaN	0	...	313.00	2725.99	53.47
4	01/11/2011 02:00	NaN	0	...	346.00	2655.64	39.87

[5 rows x 18 columns]

Team Leader: Ismail Hamdan K

Team Member: Hashir Kaamran K.A

Team Member: Izhan Ur Rahman H

Team Member: Mohammed Faheem G

Team Member: Mohammed Zakwan P