

Python Data Preprocessing Tasks

1. Handle Missing Data:

```
num.py > ...
1  import pandas as pd
2  data = {'A': [1, 2, 2], 'B': [4, 5, 5]}
3  df = pd.DataFrame(data)
4
5  df_unique = df.drop_duplicates()
6
7  print("DataFrame after removing duplicates:")
8  print(df_unique)
9
```

2. Remove Duplicate Rows:

Code:

```
num.py > ...
1  import pandas as pd
2  data = {'A': [1, 2, 2], 'B': [4, 5, 5]}
3  df = pd.DataFrame(data)
4
5  df_unique = df.drop_duplicates()
6
7  print("DataFrame after removing duplicates:")
8  print(df_unique)
9
```

3. Detect and Remove Outliers (Z-Score Method):

Code:

```
num.py > ...
1  from scipy.stats import zscore
2  import pandas as pd
3
4  data = {'A': [10, 12, 11, 100], 'B': [5, 6, 5, 50]}
5  df = pd.DataFrame(data)
6
7
8  z_scores = zscore(df)
9  df_no_outliers = df[(abs(z_scores) < 3).all(axis=1)]
10
11 print("DataFrame after removing outliers:")
12 print(df_no_outliers)
13
```

4. Convert Categorical Data to Numerical (One-Hot Encoding):

Code:

```
num.py > ...
1 import pandas as pd
2 data = {'Category': ['A', 'B', 'A']}
3 df = pd.DataFrame(data)
4
5
6 df_encoded = pd.get_dummies(df, columns=['Category'])
7
8 print("DataFrame after one-hot encoding:")
9 print(df_encoded)
10
```

5. Normalize Data Using Min-Max Normalization:

Code:

```
num.py > ...
1 import pandas as pd
2 data = {'A': [10, 20, 30], 'B': [100, 200, 300]}
3 df = pd.DataFrame(data)
4
5 df_normalized = (df - df.min()) / (df.max() - df.min())
6
7 print("DataFrame after min-max normalization:")
8 print(df_normalized)
9
```

6. Clean Text Data:

Code:

```
num.py > ...
1 import re
2
3 text = "Hello, World! Welcome to Data Science 101."
4
5 cleaned_text = re.sub(r'^a-zA-Z0-9\s', '', text).lower()
6
7 print("Cleaned text:")
8 print(cleaned_text)
9
```

What I Learned:

Through these exercises, I deepened my understanding of essential data preprocessing techniques. Handling missing data using the mean improves dataset quality for analysis. Removing duplicates ensures data consistency. Detecting and eliminating outliers with Z-scores maintains data integrity for statistical

analysis. Transforming categorical variables with one-hot encoding allows machine learning models to process non-numerical data. Min-max normalization standardizes data ranges, enhancing comparability.

Finally, text cleaning prepares raw text for natural language processing by removing noise and standardizing case. These techniques are fundamental for building reliable and accurate machine learning models.