

LAB ASSIGNMENT-17

Task-1:

Prompt:

Create a Python script to clean an employee dataset: fill missing values (salary, department, joining_date), convert joining_date to datetime, standardize department names (e.g., 'HR', 'hr', 'Human Resources' → 'HR'), and encode categorical variables (department, job_role).

Code and Output:

```
1 import pandas as pd
2 import numpy as np
3
4 # -----
5 # Step 1: Generate a Sample Employee Dataset
6 #
7 data = {
8     'employee_id': [101, 102, 103, 104, 105, 106, 107, 108],
9     'salary': [50000, np.nan, 75000, 60000, None, 80000, 55000,
10    'department': ['HR', 'hr', 'Human Resources', 'IT', 'Finance', 'hr', 'Sales', None],
11    'joining_date': ['2020-05-10', None, '2019-07-22', '', '2021/03/15', '2022-11-30', '2020-08-01', '2021-12-12'],
12    'job_role': ['Manager', 'Analyst', 'Developer', 'Clerk', 'Manager', 'Analyst', 'Developer', 'Clerk']
13 }
14
15 # Create DataFrame with missing values and inconsistent department names
16 df = pd.DataFrame(data)
17
18 # -----
19 # Step 2: Data Cleaning and Preprocessing
20 #
21
22 # Handle missing values:
23 # For 'salary', fill missing values with the median salary calculated from non-missing entries
24 df['salary'] = df['salary'].fillna(df['salary'].median())
25
26 # For 'department', fill missing values with "Unknown"
27 df['department'] = df['department'].fillna("Unknown")
28
29 # For 'joining_date', fill missing or empty strings with a default date (e.g., '2000-01-01')
30 df['joining_date'] = df['joining_date'].replace("", None).fillna("2000-01-01")
31
32 # Convert 'joining_date' column to datetime format (errors='coerce' will convert invalid formats to NaT)
```



```
33 # Convert 'joining_date' column to datetime format (errors='coerce' will convert invalid formats to NaT)
34 df['joining_date'] = pd.to_datetime(df['joining_date'], errors='coerce')
35
36 # Standardize department names:
37 # Map various forms of HR to a single standardized value "HR"
38 def standardize_department(dept):
39     dept_lower = str(dept).strip().lower()
40     if dept_lower in ['hr', 'human resources']:
41         return "HR"
42     return dept # Return the original if no standardization rule applies
43
44 df['department'] = df['department'].apply(standardize_department)
45
46 # Encode categorical variables:
47 # Using one-hot encoding for 'department' and 'job_role'
48 cleaned_df = pd.get_dummies(df, columns=['department', 'job_role'], drop_first=True)
49
50 # -----
51 # Step 3: Output the Cleaned DataFrame
52 #
53 print(cleaned_df.head())
```

```
>>> %Run -c $EDITOR_CONTENT
    employee_id  salary ... job_role_Developer  job_role_Manager
0           101  50000.0 ...             False            True
1           102  60000.0 ...             False            False
2           103  75000.0 ...             True             False
3           104  60000.0 ...             False            False
4           105  60000.0 ...             False            True
[5 rows x 10 columns]
```

Explanation:

The code performs the following key operations:

- **Data Generation:**
Creates a sample employee dataset with various issues like missing values and inconsistent department names.
- **Missing Values Handling:**
 - Fills missing salaries with the median salary.
 - Replaces missing department entries with "Unknown".
 - Substitutes missing or empty `joining_date` values with a default date ("2000-01-01").
- **Data Type Conversion:**
Converts the `joining_date` column into datetime objects to ensure consistent date formatting.
- **Department Standardization:**
Uses a custom function to map different forms of "HR" (like "hr" or "Human Resources") to a single, standardized value ("HR").
- **Categorical Encoding:**
Applies one-hot encoding to `department` and `job_role`, while dropping the first category to avoid multicollinearity.
- **Output:**
Finally, prints the head of the cleaned DataFrame for a quick preview of the processed data.

Task 2:

Prompt:

Generate a Python script to preprocess a sales dataset: convert `transaction_date` to datetime, add Month-Year column, remove rows with non-positive `transaction_amount`, and apply Min-Max scaling to `transaction_amount`.

Code and Output:

```
1 import pandas as pd
2
3 # -----
4 # Step 1: Generate a Sample Sales Transaction Dataset
5 #
6 data = {
7     'transaction_id': [1, 2, 3, 4, 5, 6],
8     'transaction_date': ['2025-01-15', '2025/02/20', '15-03-2025', '2025-04-05', '2025-05-10', '2025-06-25'],
9     'transaction_amount': [250.0, -50.0, 0.0, 500.0, 750.0, 300.0]
10 }
11
12 df = pd.DataFrame(data)
13 # Step 2: Preprocess the Dataset
14
15 # Convert 'transaction_date' column to datetime format (invalid dates become NaT)
16 df['transaction_date'] = pd.to_datetime(df['transaction_date'], errors='coerce')
17
18 # Create a new column 'Month-Year' extracted from the 'transaction_date'
19 df['Month-Year'] = df['transaction_date'].dt.strftime('%b-%Y')
20
21 # Remove rows with negative or zero transaction amounts
22 df = df[df['transaction_amount'] > 0]
23
24 # Normalize the 'transaction_amount' column using Min-Max scaling without sklearn
25 min_amount = df['transaction_amount'].min()
26 max_amount = df['transaction_amount'].max()
27 df['transaction_amount'] = (df['transaction_amount'] - min_amount) / (max_amount - min_amount)
28
29 # Step 3: Output the Preprocessed DataFrame
30 print(df)
```

```
>>> %Run -c $EDITOR_CONTENT
transaction_id transaction_date transaction_amount Month-Year
0            1    2025-01-15          0.0   Jan-2025
3            4    2025-04-05          0.5   Apr-2025
4            5    2025-05-10          1.0   May-2025
5            6    2025-06-25          0.1   Jun-2025
```

Explanation:

1. **Dataset Creation**
 - Creates a dictionary with three keys:
 - `transaction_id`: Unique IDs for each transaction.
 - `transaction_date`: Dates of transactions in various formats.
 - `transaction_amount`: Monetary values for each transaction.
 - Converts this dictionary into a pandas DataFrame.
2. **Data Preprocessing**
 - **Convert Dates:**
 - Uses `pd.to_datetime` with `errors='coerce'` to convert the `transaction_date` column to datetime format.
 - Any invalid date strings are set to `NaT` (Not a Time).
 - **Extract Month-Year:**
 - Creates a new column `Month-Year` by extracting the month and year from the valid `transaction_date` values using the `dt.strftime('%b-%Y')` method.
 - **Filter Transactions:**
 - Removes rows where `transaction_amount` is less than or equal to zero, keeping only rows with a positive amount.
 - **Normalize Amounts:**
 - Applies Min-Max scaling to the `transaction_amount` column without using sklearn.
 - This is done by subtracting the minimum value and dividing by the range (max - min) so that all transaction amounts are scaled between 0 and 1.
3. **Output**
 - Finally, the cleaned and preprocessed DataFrame is printed.

Task 3:

Prompt:

Write a Python script to clean healthcare patient data: fill missing numeric values with column mean, convert height from cm to meters, standardize gender labels (e.g., 'M', 'Male', 'male' → 'Male'), and drop the patient_id column.

Code and Output:

```
1 import pandas as pd
2 import numpy as np
3
4 # -----
5 # Generate a Sample Healthcare Patient Dataset
6 #
7 data = {
8     'patient_id': [1, 2, 3, 4, 5, 6, 7, 8],
9     'blood_pressure': [120, 130, np.nan, 110, 140, np.nan, 125, 135],
10    'heart_rate': [80, np.nan, 70, 75, 90, 85, np.nan, 78],
11    'height': [170, 165, 180, np.nan, 160, 175, 168, 172], # Heights in cm
12    'gender': ['M', 'Male', 'male', 'F', 'Female', 'f', 'F', 'female'],
13    'age': [25, 35, 45, 55, np.nan, 65, 75, 85] # Additional numeric column
14 }
15 df = pd.DataFrame(data)
16
17 # -----
18 # Data Cleaning and Preprocessing
19 #
20
21 # 1. Fill missing values in numeric columns with the column mean
22 numeric_columns = ['blood_pressure', 'heart_rate', 'height', 'age']
23 for col in numeric_columns:
24     df[col].fillna(df[col].mean(), inplace=True)
25
26 # 2. Standardize units: Convert height from centimeters to meters
27 df['height'] = df['height'] / 100
28
29 # 3. Correct inconsistent categorical labels for 'gender'
30 def standardize_gender(g):
31     g_clean = str(g).strip().lower()
32
```

```
33     if g_clean in ['m', 'male']:
34         return 'Male'
35     elif g_clean in ['f', 'female']:
36         return 'Female'
37     return g # Return original if no rule applies
38
39 df['gender'] = df['gender'].apply(standardize_gender)
40
41 # 4. Drop irrelevant columns such as patient_id
42 df_cleaned = df.drop(columns=['patient_id'])
43
44 # -----
45 # Output the Cleaned Dataset
46 #
47 print("Cleaned Healthcare Patient Records:")
48 print(df_cleaned.head())

```

```
Cleaned Healthcare Patient Records:
  blood_pressure heart_rate height gender age
0      120.000000   80.000000  1.70  Male 25.0
1      130.000000   79.666667  1.65  Male 35.0
2      126.666667   70.000000  1.80  Male 45.0
3      110.000000   75.000000  1.70 Female 55.0
4      140.000000   90.000000  1.60 Female 55.0
```

```
>>>
```

Explanation:

```
1. Import Libraries
The code starts by importing the necessary libraries:
  o pandas for handling data in tabular form.
  o numpy for handling numerical data and to use the np.nan value as missing data.

2. Generate a Sample Healthcare Dataset
A dictionary called data is defined which contains sample patient records. Each key in the dictionary represents a column (like patient_id, blood_pressure, etc.) and the values are lists representing the data entries. This data is then converted into a pandas DataFrame called df.

3. Fill Missing Values in Numeric Columns
  o The code identifies numeric columns such as blood_pressure, heart_rate, height, and age.
  o It loops over these columns and fills any missing values (denoted by np.nan) with the mean value of that column. This ensures that there are no empty values in these fields.

4. Standardize Units for Height
The height column is stored in centimeters. The code divides the values by 100 to convert them into meters, which might be the preferred unit for analysis or modeling.

5. Correct Inconsistent Categorical Labels for Gender
  o A function called standardize_gender is defined. It takes a gender value, converts it to lowercase, and then maps it to either "Male" or "Female".
  o This function is applied to the gender column to standardize the labels, ensuring consistency in the data.

6. Drop the Irrelevant patient_id Column
After cleaning the records, the patient_id column is removed from the DataFrame since it is assumed not to be useful for further analysis (e.g., in a machine learning model).

7. Display the Cleaned Dataset
Finally, the cleaned DataFrame (df_cleaned) is printed to the console. This shows the first few rows of the cleaned data, which is now ready for further analysis or modeling.
```

Task 4:

Prompt:

Write a Python script to clean social media text data: remove special characters, URLs, and emojis; convert text to lowercase; tokenize and remove stopwords; and apply lemmatization to prepare for sentiment analysis.

Code & Output:

```
1 import re
2 import pandas as pd
3
4 # Sample dataset: a list of social media text posts
5 data = [
6     "text": [
7         "Loving this new product! 🌟 #amazing http://example.com",
8         "Worst service ever!! @@ Visit http://badservice.com for details!",
9         "Happy days! Enjoying life & sunshine ☀️ https://weather.com",
10        "OMG! Can't believe it!!! #shocked 🤯 http://trending.com",
11        "Feeling blessed & grateful. Life is beautiful 😊"
12    ]
13 }
14
15 df = pd.DataFrame(data)
16 print("Original Dataset:")
17 print(df)
18
19 # Define a simple stopword list (you can expand this list as needed)
20 stop_words = [
21     "the", "and", "is", "in", "it", "for", "to", "a",
22     "of", "this", "that", "with", "on", "at", "ever", "cant"
23 ]
24
25 def clean_text(text):
26     # Remove URLs
27     text = re.sub(r'http\S+', '', text)
28     # Remove emojis using an approximate Unicode regex pattern
29     emoji_pattern = re.compile(
30         "["
31         "\U0001F600-\U0001F64F" # Emoticons
32         "\U0001F300-\U0001F5FF" # Symbols & pictographs
33
34         "\U0001F680-\U0001F6FF" # Transport & map symbols
35         "\U0001F1E0-\U0001F1F" # Flags
36         "]",
37         flags=re.UNICODE)
38     text = emoji_pattern.sub('', text)
39     # Remove special characters (keep letters, numbers, and whitespace)
40     text = re.sub(r'[\W\s]', '', text)
41     # Convert to lowercase
42     text = text.lower()
43     return text
44
45 def simple_tokenize(text):
46     # A very basic tokenizer that splits on whitespace
47     return text.split()
48
49 def simple_lemmatize(word):
50     # A basic and naive lemmatization:
51     # Remove common suffixes: 'ing', 'ed', 'es' and trailing 's'
52     if word.endswith('ing') and len(word) > 4:
53         return word[:-3]
54     if word.endswith('ed') and len(word) > 3:
55         return word[:-2]
56     if word.endswith('es') and len(word) > 3:
57         return word[:-2]
58     if word.endswith('s') and len(word) > 2:
59         return word[:-1]
60     return word
61
62 def preprocess_text(text):
63     # Clean the text to remove URLs, emojis, special characters, and convert to lowercase
64     cleaned = clean_text(text)
65     # Tokenize the cleaned text
66     tokens = simple_tokenize(cleaned)
67
68     # Remove stopwords
69     filtered_tokens = [token for token in tokens if token not in stop_words]
70     # Apply simple lemmatization
71     lemmatized_tokens = [simple_lemmatize(token) for token in filtered_tokens]
72     # Return the processed text as a string
73     return ' '.join(lemmatized_tokens)
74
75 df['processed_text'] = df['text'].apply(preprocess_text)
76
77 print("\nProcessed Dataset:")
78 print(df[['processed_text']])
```

```
>>> %Run -c $EDITOR_CONTENT
Original Dataset:
text
0 Loving this new product! 🌟 #amazing http://exa...
1 Worst service ever!! @@ Visit http://badservice...
2 Happy days! Enjoying life & sunshine ☀️ https://...
3 OMG! Can't believe it!!! #shocked 🤯 http://trend...
4 Feeling blessed & grateful. Life is beautiful 😊

Processed Dataset:
processed_text
0      lov new product amaz
1      worst service visit detail
2      happy day enjoy life sunshine
3          omg believe shock
4  feel bless grateful life beautiful
>>>
```

Explanation:

- **Importing libraries and data setup:**
The script imports the required libraries:
 - `re` for working with regular expressions.
 - `pandas` for handling and manipulating the dataset.It then creates a sample dataset (a few social media posts) as a Python dictionary and converts it into a DataFrame. Finally, it prints the original dataset.
- **Defining stopwords:**
A basic set of stopwords (common words that are typically removed) is defined. This helps to filter out words that may not be important for sentiment analysis.
- **Cleaning the text:**
The `clean_text` function performs several cleaning steps:
 - **Remove URLs:** It uses a regex pattern to remove any web links.
 - **Remove emojis:** It applies a regex that targets emoji characters and removes them.
 - **Remove special characters:** It keeps only letters, numbers, and whitespace.
 - **Convert to lowercase:** It changes all characters to lowercase for uniformity.
- **Tokenization:**
The `simple_tokenize` function splits the cleaned text into words by splitting on whitespace. This breaks the text into manageable tokens.
- **Simple lemmatization:**
The `simple_lemmatize` function removes common suffixes like "ing", "ed", "es" and a trailing "s" from each word. This is a very basic method to reduce words to their base form (lemmatization).
- **Preprocessing text:**
The `preprocess_text` function puts all these steps together:
 - It cleans the text using `clean_text`.
 - It tokenizes the text using `simple_tokenize`.
 - It filters out any stopwords from these tokens.
 - It applies the simple lemmatization to the remaining tokens.
 - Finally, it joins these tokens back into a single string which represents the cleaned text.
- **Applying preprocessing and printing:**
The script applies the `preprocess_text` function on each entry in the original "text" column of the DataFrame, creating a new column called "processed_text." It then prints this processed dataset, which is ready for further NLP sentiment analysis.

Task 5:

Prompt:

Write a Python script to preprocess a financial dataset: fill missing values in stock price and volume, create 7-day and 30-day moving averages, normalize continuous columns with StandardScaler, and encode categorical columns like sector and company_name.

Code & Output:

```
1 import pandas as pd
2 import numpy as np
3
4 # Generate a simulated financial dataset
5 np.random.seed(42)
6 num_days = 60
7 dates = pd.date_range(start='2025-01-01', periods=num_days, freq='D')
8 sectors = ['Technology', 'Finance', 'Healthcare']
9 companies = ['CompanyA', 'CompanyB', 'CompanyC']
10
11 data = {
12     'date': dates,
13     'stock_price': np.random.uniform(100, 200, size=num_days),
14     'volume': np.random.randint(1000, 10000, size=num_days),
15     'sector': np.random.choice(sectors, size=num_days),
16     'company_name': np.random.choice(companies, size=num_days)
17 }
18 df = pd.DataFrame(data)
19
20 # Introduce missing values randomly in stock_price and volume
21 for col in ['stock_price', 'volume']:
22     df.loc[df.sample(frac=0.1, random_state=42).index, col] = np.nan
23
24 # Handle missing values: use median imputation
25 for col in ['stock_price', 'volume']:
26     median_val = df[col].median()
27     df[col].fillna(median_val, inplace=True)
28
29 # Create new features: moving averages (7-day and 30-day) of stock_price.
30 df['MA_7'] = df['stock_price'].rolling(window=7, min_periods=1).mean()
31 df['MA_30'] = df['stock_price'].rolling(window=30, min_periods=1).mean()
32
33 # List of continuous columns to normalize
34 cont_cols = ['stock_price', 'volume', 'MA_7', 'MA_30']
35
36 # Normalize continuous variables using manual StandardScaler approach
37 for col in cont_cols:
38     mean_val = df[col].mean()
39     std_val = df[col].std()
40     df[col + '_norm'] = (df[col] - mean_val) / std_val
41
42 # Encode categorical columns using dummy variables
43 cat_cols = ['sector', 'company_name']
44 df_encoded = pd.get_dummies(df, columns=cat_cols, drop_first=True)
45
46 # Final feature engineered DataFrame
47 print("Feature engineered dataset:")
48 print(df_encoded.head())
49
50 # Optionally, save to CSV file for further ML tasks
51 df_encoded.to_csv("feature_engineered_financial_data.csv", index=False)
52 print("\nData saved to feature_engineered_financial_data.csv")
```

```
Feature engineered dataset:
   date  stock_price ... company_name_CompanyB  company_name_CompanyC
0 2025-01-01    147.562345 ...                  True                  False
1 2025-01-02    195.071431 ...                  True                  False
2 2025-01-03    173.199394 ...                 False                  True
3 2025-01-04    159.865848 ...                 False                  False
4 2025-01-05    115.601864 ...                 False                  False
```

[5 rows x 13 columns]

Data saved to feature_engineered_financial_data.csv

>>>

Explanation:

1. Import Modules

The code imports necessary packages:

- o `pandas` for handling data as DataFrames
- o `numpy` for numerical operations

2. Generate a Simulated Financial Dataset

- o A date range of 60 days starting from January 1, 2025, is created.
- o Random data is generated for stock prices and trading volume.
- o Two categorical columns (`sector` and `company_name`) are generated from predefined lists.
- o All these columns are combined into a DataFrame.

3. Introduce Missing Values

- o The code randomly selects 10% of the rows for `stock_price` and `volume` columns, replacing those values with `NaN` (to simulate missing data).

4. Handle Missing Values

- o The missing values in `stock_price` and `volume` are filled with the median value of each respective column.

5. Calculate Moving Averages

- o Two new features are created:
 - `MA_7`: the 7-day moving average of the stock price.
 - `MA_30`: the 30-day moving average of the stock price.
- o The rolling mean is computed using a window of 7 or 30 days.

6. Normalize Continuous Variables Without Sklearn

- o For each continuous column (stock price, volume, and the moving averages), the script computes the mean and standard deviation.
- o A new column is added for each feature (appending `_norm`) that represents the standardized value, calculated as:
$$(\text{value} - \text{mean}) / \text{std}$$

7. Encode Categorical Columns

- o The categorical columns (`sector` and `company_name`) are encoded using one-hot encoding (creating dummy variables).
- o The `drop_first=True` option is used to avoid creating redundant columns.

8. Output the Final DataFrame

- o The code prints the first few rows of the processed DataFrame, showing the new features and normalized columns.
- o The final dataset is saved as a CSV file named "feature_engineered_financial_data.csv" for later machine learning tasks.