# AI LAB 17

**Name: K. Jashuva**

**Roll No: 2503A51L16**

**Batch No: 19**

## TASK 1 - Employee Data Preprocessing

- Use AI to generate a Python script for cleaning an employee dataset

## Instructions:

• Handle missing values in columns (salary, department, joining_date).

• Convert the "joining_date" column into proper datetime format.

• Standardize department names (e.g., "HR", "hr", "Human Resources" → "HR").

• Encode categorical variables (department, job_role)

## Prompt:

Clean the employee dataset by filling missing values in salary, department, and joining_date, convert joining_date to datetime, standardize department names (e.g., unify "HR" variants), and encode categorical variables (department, job_role) for modeling.

## Code :

```python
import pandas as pd
import numpy as np
from sklearn.preprocessing import LabelEncoder

# Sample employee data (replace this with your actual dataset)
data = {
    'employee_id': [1, 2, 3, 4, 5],
    'salary': [50000, np.nan, 55000, 60000, np.nan],
    'department': ['HR', 'human resources', 'IT', 'hr', 'Finance'],
    'joining_date': ['2020-05-01', '2021-06-15', np.nan, '2022-07-20', '2021-08-30'],
    'job_role': ['Manager', 'Analyst', 'Developer', 'Manager', 'Analyst']
}

# Create DataFrame
df = pd.DataFrame(data)

# 1. Handle missing values
# Fill missing salary with median salary (without inplace=True)
df['salary'] = df['salary'].fillna(df['salary'].median())

# Fill missing joining_date with a placeholder (e.g., 'Unknown') or with a default date
df['joining_date'] = df['joining_date'].fillna('2020-01-01')

# 2. Convert 'joining_date' to proper datetime format
df['joining_date'] = pd.to_datetime(df['joining_date'], format='%Y-%m-%d')

# 3. Standardize department names
df['department'] = df['department'].replace({
    'human resources': 'HR',
    'hr': 'HR',
    'Human Resources': 'HR'
})

# 4. Encode categorical variables (department, job_role)
label_encoder = LabelEncoder()

# Encode 'department'
df['department_encoded'] = label_encoder.fit_transform(df['department'])

# Encode 'job_role'
df['job_role_encoded'] = label_encoder.fit_transform(df['job_role'])

# Final cleaned DataFrame
print(df)
```

## Output:



```
PS C:\Users\HP> & C:/Users/HP/AppData/Local/Microsoft/WindowsApps/python3.13.exe "c:/Users/HP/Documents/TASK 1.py"
   employee_id  salary department joining_date  job_role  department_encoded  job_role_encoded
PS C:\Users\HP> & C:/Users/HP/AppData/Local/Microsoft/WindowsApps/python3.13.exe "c:/Users/HP/Documents/TASK 1.py"
   employee_id  salary department joining_date  job_role  department_encoded  job_role_encoded
0            1  50000.0        HR   2020-05-01   Manager                   1                 2
1            2  55000.0        HR   2021-06-15   Analyst                   1                 0
2            3  55000.0        IT   2020-01-01  Developer                  2                 1
3            4  60000.0        HR   2022-07-20   Manager                   1                 2
4            5  55000.0   Finance   2021-08-30   Analyst                   0                 0
PS C:\Users\HP>
```

## Observations:

The dataset contains missing values and inconsistent department names that need cleaning for accuracy. The joining_date column requires conversion to a proper datetime format, and categorical fields like department and job_role should be encoded. After preprocessing, the dataset will be clean, consistent, and ready for analysis or modeling.

# Task 2 – Sales Transaction Data Preprocessing

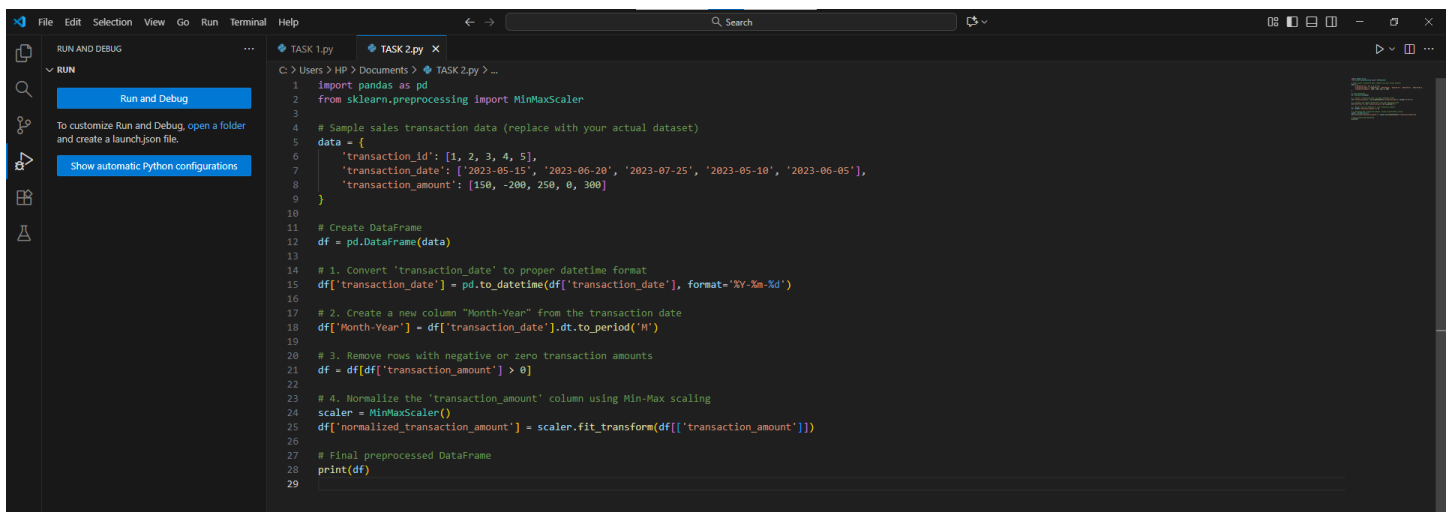**Task:** Use AI to generate a script for preprocessing a sales transaction dataset.

**Instructions:**

- Convert transaction dates to proper datetime format.

- Create a new column for "Month-Year" from the transaction date.

- Remove rows with negative or zero transaction amounts.

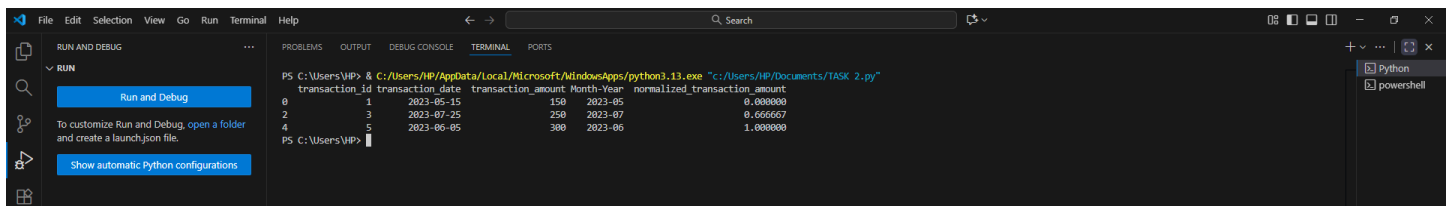- Normalize the "transaction_amount" column using Min-Max scaling

**Prompt :**

Create a Python script to preprocess a sales transaction dataset by converting transaction dates to datetime format, adding a "Month-Year" column, and removing rows with zero or negative amounts. Normalize the "transaction_amount" column using Min-Max scaling.

**Code:**



**Output:**

## Observations:

The sales transaction dataset requires cleaning and transformation to ensure data accuracy. Transaction dates need to be standardized, and a new Month-Year column will help in time-based analysis. Removing invalid transaction amounts and normalizing values will prepare the data for consistent and reliable analysis or modeling.

## Task 3 – Healthcare Patient Records Cleaning

## Task:

Use AI to generate a script for cleaning healthcare patient records.

## Instructions:

• Fill missing values in numeric columns (e.g., blood_pressure, heart_rate) with column mean.

• Standardize units (convert height from cm to meters).

• Correct inconsistent categorical labels (e.g., "M", "Male", "male"→ "Male").

• Drop irrelevant columns such as patient_id after cleaning.

## Prompt:

Create a Python script to clean healthcare patient records by filling missing numeric values with the column mean, converting height from centimeters to meters, standardizing gender labels (e.g., "M", "Male", "male" → "Male"), and dropping irrelevant columns like patient_id after cleaning.

## Code:

```python
import pandas as pd
import numpy as np

# Sample healthcare patient records data (replace with your actual dataset)
data = {
    'patient_id': [101, 102, 103, 104, 105],
    'age': [45, 60, 38, np.nan, 55],
    'blood_pressure': [120, np.nan, 130, 140, 115],
    'heart_rate': [80, 75, np.nan, 72, 85],
    'height_cm': [170, 180, 160, 175, 165],
    'gender': ['M', 'Male', 'male', 'M', 'Female']
}

# Create DataFrame
df = pd.DataFrame(data)

# 1. Fill missing values in numeric columns with the column mean
df['blood_pressure'] = df['blood_pressure'].fillna(df['blood_pressure'].mean())
df['heart_rate'] = df['heart_rate'].fillna(df['heart_rate'].mean())
df['age'] = df['age'].fillna(df['age'].mean())

# 2. Standardize units - Convert height from cm to meters
df['height_m'] = df['height_cm'] / 100

# 3. Correct inconsistent categorical labels (e.g., 'M', 'Male', 'male' → 'Male')
df['gender'] = df['gender'].replace({'M': 'Male', 'male': 'Male'})

# 4. Drop irrelevant columns (e.g., patient_id)
df.drop(columns=['patient_id', 'height_cm'], inplace=True)

# Final cleaned DataFrame
print(df)
```

**Output:**



```
PS C:\Users\HP> & C:/Users/HP/AppData/Local/Microsoft/WindowsApps/python3.13.exe "c:/Users/HP/Documents/TASK 3.py"
    age  blood_pressure  heart_rate  gender  height_m
0  45.0          120.00        80.0    Male      1.70
1  60.0          126.25        75.0    Male      1.80
2  38.0          130.00        78.0    Male      1.60
3  49.5          140.00        72.0    Male      1.75
4  55.0          115.00        85.0  Female      1.65
PS C:\Users\HP>
```

**Observations:**

The healthcare dataset contains missing values and inconsistent data formats that need correction for accuracy. Numeric columns require imputation with mean values, and units like height must be standardized for consistency. Cleaning categorical labels and removing irrelevant columns will result in a reliable and well-structured dataset for analysis.

## Task 4 – Social Media Sentiment Dataset Preparation

**Task:**

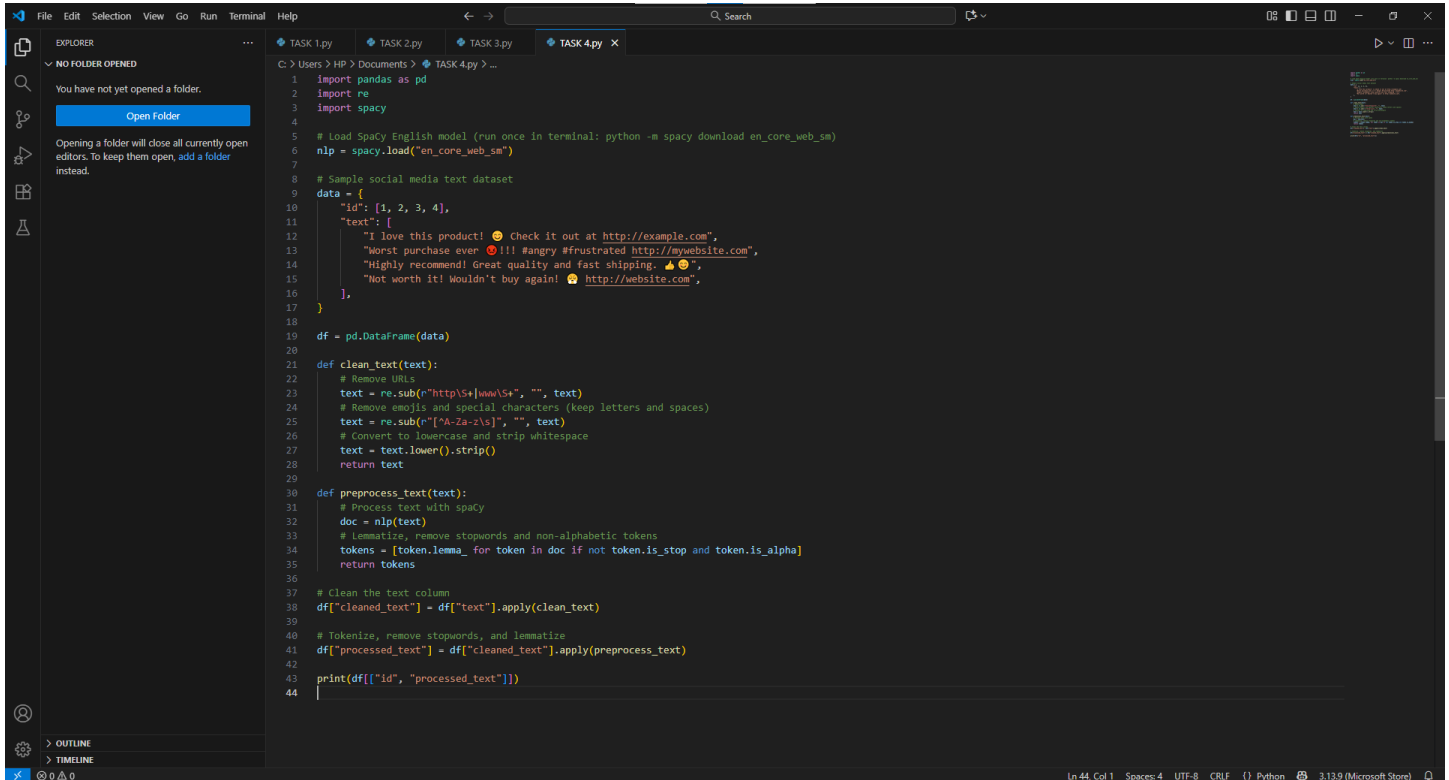Use AI to write a script to preprocess a social media text dataset.

**Instructions:**

• Remove special characters, URLs, and emojis from text.

• Convert all text to lowercase.

• Tokenize and remove stopwords.

• Apply lemmatization for standardizing words
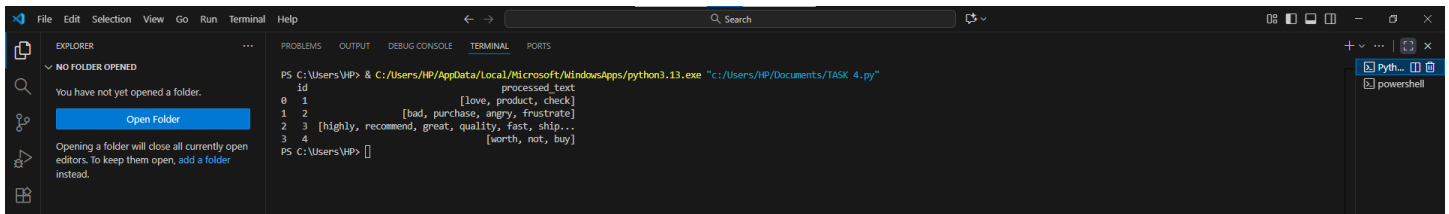
## Prompt:

Create a Python script to preprocess a social media sentiment dataset by removing special characters, URLs, and emojis from text. Convert all text to lowercase, tokenize it, remove stopwords, and apply lemmatization to standardize words for accurate sentiment analysis.

## Code:

```python
import pandas as pd
import re
import spacy

# Load SpaCy English model (run once in terminal: python -m spacy download en_core_web_sm)
nlp = spacy.load("en_core_web_sm")

# Sample social media text dataset
data = {
    "id": [1, 2, 3, 4],
    "text": [
        "I love this product! 😊 Check it out at http://example.com",
        "Worst purchase ever 😡!!! #angry #frustrated http://mywebsite.com",
        "Highly recommend! Great quality and fast shipping. 👍😊",
        "Not worth it! Wouldn't buy again! 😠 http://website.com",
    ],
}

df = pd.DataFrame(data)

def clean_text(text):
    # Remove URLs
    text = re.sub(r"http\S+|www\S+", "", text)
    # Remove emojis and special characters (keep letters and spaces)
    text = re.sub(r"[^A-Za-z\s]", "", text)
    # Convert to lowercase and strip whitespace
    text = text.lower().strip()
    return text

def preprocess_text(text):
    # Process text with spaCy
    doc = nlp(text)
    # Lemmatize, remove stopwords and non-alphabetic tokens
    tokens = [token.lemma_ for token in doc if not token.is_stop and token.is_alpha]
    return tokens

# Clean the text column
df["cleaned_text"] = df["text"].apply(clean_text)

# Tokenize, remove stopwords, and lemmatize
df["processed_text"] = df["cleaned_text"].apply(preprocess_text)

print(df[["id", "processed_text"]])
```

## Output:

```
PS C:\Users\HP> & C:/Users/HP/AppData/Local/Microsoft/WindowsApps/python3.13.exe "c:/Users/HP/Documents/TASK 4.py"
   id                                     processed_text
0   1                               [love, product, check]
1   2                      [bad, purchase, angry, frustrate]
2   3        [highly, recommend, great, quality, fast, ship...
3   4                                    [worth, not, buy]
PS C:\Users\HP>
```

## Observations:

The social media text data contains noise such as emojis, URLs, and special characters that can affect sentiment accuracy. Standardizing text through lowercasing, tokenization, and lemmatization ensures consistency. After preprocessing, the dataset will be cleaner and more suitable for effective sentiment analysis and model training.

## Task 5 – Financial Dataset Feature Engineering

## Task:

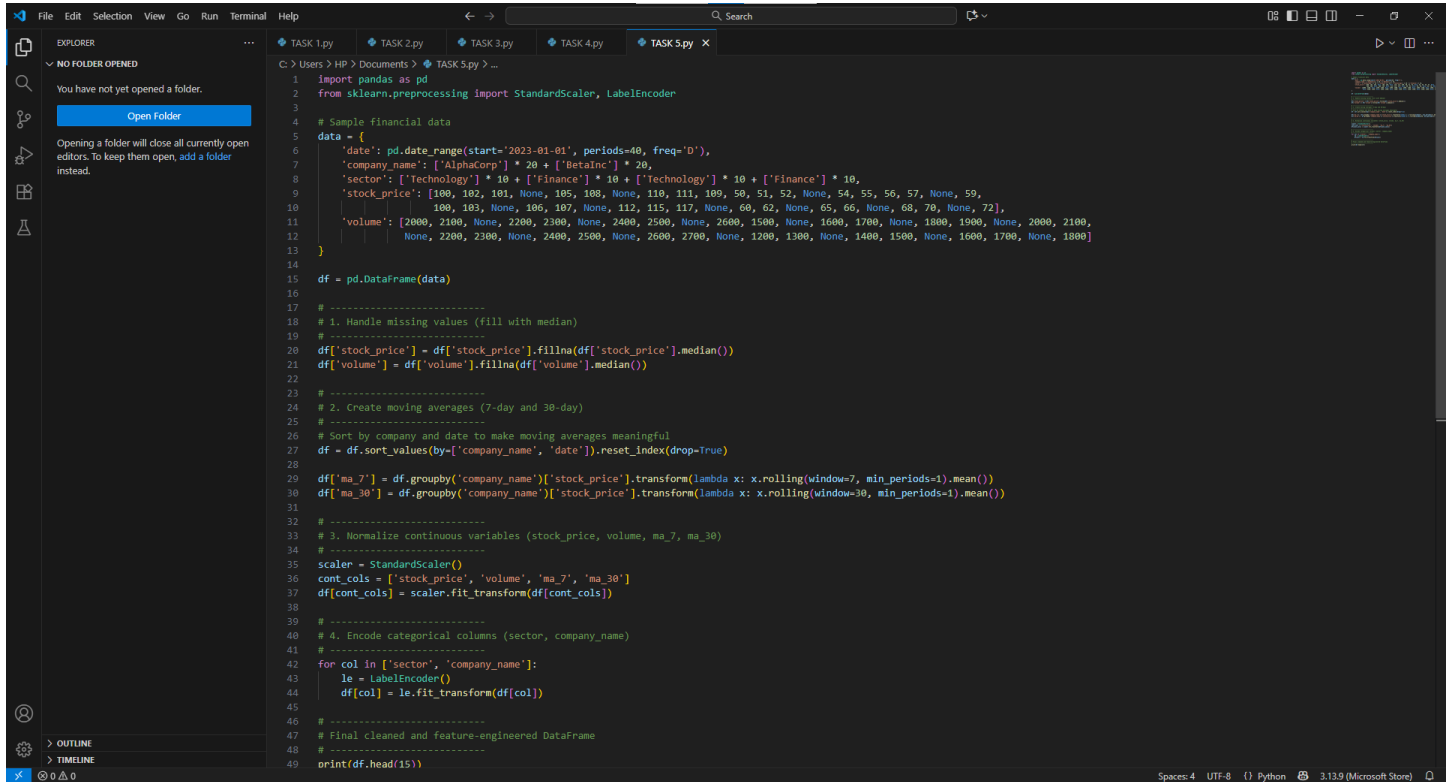Use AI to create a preprocessing script for a financial dataset.

## Instructions:

• Handle missing values in stock price and volume.

• Create new features such as moving average (7-day, 30-day).

• Normalize continuous variables using StandardScaler.

• Encode categorical columns (sector, company_name).

## Prompt:

Create a Python script to preprocess a financial dataset by handling missing values in stock price and volume, generating new features like 7-day and 30-day moving averages, normalizing continuous variables using StandardScaler, and encoding categorical columns such as sector and company_name.

## Code:



```python
import pandas as pd
from sklearn.preprocessing import StandardScaler, LabelEncoder

# Sample financial data
data = {
    'date': pd.date_range(start='2023-01-01', periods=40, freq='D'),
    'company_name': ['AlphaCorp'] * 20 + ['BetaInc'] * 20,
    'sector': ['Technology'] * 10 + ['Finance'] * 10 + ['Technology'] * 10 + ['Finance'] * 10,
    'stock_price': [100, 102, 101, None, 105, 108, None, 110, 111, 109, 50, 51, 52, None, 54, 55, 56, 57, None, 59,
                    100, 103, None, 106, 107, None, 112, 115, 117, None, 60, 62, None, 65, 66, None, 68, 70, None, 72],
    'volume': [2000, 2100, None, 2200, 2300, None, 2400, 2500, None, 2600, 1500, None, 1600, 1700, None, 1800, 1900, None, 2000, 2100,
               None, 2200, 2300, None, 2400, 2500, None, 2600, 2700, None, 1200, 1300, None, 1400, 1500, None, 1600, 1700, None, 1800]
}

df = pd.DataFrame(data)

# ----------------------------
# 1. Handle missing values (fill with median)
# ----------------------------
df['stock_price'] = df['stock_price'].fillna(df['stock_price'].median())
df['volume'] = df['volume'].fillna(df['volume'].median())

# ----------------------------
# 2. Create moving averages (7-day and 30-day)
# ----------------------------
# Sort by company and date to make moving averages meaningful
df = df.sort_values(by=['company_name', 'date']).reset_index(drop=True)

df['ma_7'] = df.groupby('company_name')['stock_price'].transform(lambda x: x.rolling(window=7, min_periods=1).mean())
df['ma_30'] = df.groupby('company_name')['stock_price'].transform(lambda x: x.rolling(window=30, min_periods=1).mean())

# ----------------------------
# 3. Normalize continuous variables (stock_price, volume, ma_7, ma_30)
# ----------------------------
scaler = StandardScaler()
cont_cols = ['stock_price', 'volume', 'ma_7', 'ma_30']
df[cont_cols] = scaler.fit_transform(df[cont_cols])

# ----------------------------
# 4. Encode categorical columns (sector, company_name)
# ----------------------------
for col in ['sector', 'company_name']:
    le = LabelEncoder()
    df[col] = le.fit_transform(df[col])

# ----------------------------
# Final cleaned and feature-engineered DataFrame
# ----------------------------
print(df.head(15))
```

## Output:



## Observations:

The financial dataset requires handling missing values to maintain data reliability. Creating moving average features will help capture market trends, while normalization ensures balanced scaling for modeling. Encoding categorical variables makes the dataset suitable for machine learning and predictive analytics.