

SCHOOL OF COMPUTER SCIENCE AND ARTIFICIAL INTELLIGENCE		DEPARTMENT OF COMPUTER SCIENCE ENGINEERING	
Program Name: B. Tech		Assignment Type: Lab	
Course Coordinator Name		Venkataramana Veeramsetty	
Instructor(s) Name		Dr. V. Venkataramana (Co-ordinator) Dr. T. Sampath Kumar Dr. Pramoda Patro Dr. Brij Kishor Tiwari Dr.J.Ravichander Dr. Mohammand Ali Shaik Dr. Anirodh Kumar Mr. S.Naresh Kumar Dr. RAJESH VELPULA Mr. Kundhan Kumar Ms. Ch.Rajitha Mr. M Prakash Mr. B.Raju Intern 1 (Dharma teja) Intern 2 (Sai Prasad) Intern 3 (Sowmya) NS_2 (Mounika)	
Course Code	24CS002PC215	Course Title	AI Assisted Coding
Year/Sem	II/I	Regulation	R24
Date and Day of Assignment	Week9 - Thursday	Time(s)	
Duration	2 Hours	Applicable to Batches	
AssignmentNumber: 17.4(Present assignment number)/ 24 (Total number of assignments)			
Q.No.	Question		Expected Time to complete
1	Lab 17 – AI for Data Processing: Data Cleaning and Preprocessing Scripts Lab Objectives: <ul style="list-style-type: none"> Learn how to clean raw datasets using AI-assisted Python scripting. 		Week9 - Thursday

	<ul style="list-style-type: none"> Apply preprocessing techniques such as handling missing values, encoding categorical data, and normalization. Automate repetitive data-cleaning tasks with AI-generated code. Understand how preprocessing impacts model performance. 	
	<p>Task 1 – Employee Data Preprocessing</p> <p>Task: Use AI to generate a Python script for cleaning an employee dataset.</p> <p>Instructions:</p> <ul style="list-style-type: none"> 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). <p>Expected Output:</p> <ul style="list-style-type: none"> A cleaned Pandas DataFrame with consistent departments, proper dates, and encoded features. <p>CODE:</p> <pre>Task1.py > ... 1 import pandas as pd 2 import numpy as np 3 4 # Create a sample employee dataset 5 data = { 6 'employee_id': [101, 102, 103, 104, 105, 106, 107], 7 'name': ['Anita', 'Rahul', 'Suresh', 'Divya', 'Neha', 'Ravi', 'Meena'], 8 'salary': [55000, 62000, None, 58000, None, 60000, 75000], 9 'department': ['HR', 'human resources', 'hr', 'Sales & Marketing', 'Finance', None, 'Human Resources'], 10 'job_role': ['Manager', 'Executive', 'HR Officer', 'Sales Rep', 'Accountant', 'Clerk', 'Recruiter'], 11 'joining_date': ['2020-01-15', None, '2019-03-10', '2021-07-22', '2020-12-05', '2018-09-14', None] 12 } 13 14 15 df = pd.DataFrame(data) 16 print("Sample Employee Dataset (Raw):\n") 17 print(df) 18 df.to_csv('employee_data.csv', index=False) 19 print("\nDataset saved as 'employee_data.csv'")</pre> <p>OUTPUT:</p> <pre>Sample Employee Dataset (Raw): employee_id name salary department job_role joining_date 0 101 Anita 55000.0 HR Manager 2020-01-15 1 102 Rahul 62000.0 human resources Executive None 2 103 Suresh NaN hr HR Officer 2019-03-10 3 104 Divya 58000.0 Sales & Marketing Sales Rep 2021-07-22 4 105 Neha NaN Finance Accountant 2020-12-05 5 106 Ravi 60000.0 None Clerk 2018-09-14 6 107 Meena 75000.0 Human Resources Recruiter None Dataset saved as 'employee_data.csv'</pre>	

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.

Expected Output:

- A preprocessed DataFrame with valid dates, normalized amounts, and no invalid records.

CODE:

```
Task2.py > [+] data
1 import pandas as pd
2 from sklearn.preprocessing import MinMaxScaler
3 data = {
4     'transaction_id': [1,2,3,4,5,6,7,8,9,10],
5     'customer_id': [101,102,103,104,105,106,107,108,109,110],
6     'transaction_date': ['2025-10-01','2025-10-03','2025-10-05','2025-11-01',
7     '2025-11-02','2025-11-03','invalid_date','2025-11-05',
8     '2025-11-06','2025-11-07'],
9     'transaction_amount': [250,0,-50,450,300,500,200,0,150,600],
10    'product': ['Notebook','Pen','Marker','Printer','Paper','Laptop','Mouse','Keyboard','Stapler','Desk']
11 }
12 df = pd.DataFrame(data)
13 print("Original Dataset:")
14 print(df)
15 df['transaction_date'] = pd.to_datetime(df['transaction_date'], errors='coerce')
16 df = df.dropna(subset=['transaction_date'])
17 df['month_year'] = df['transaction_date'].dt.to_period('M')
18 df = df[df['transaction_amount'] > 0]
19 scaler = MinMaxScaler()
20 df['transaction_amount_normalized'] = scaler.fit_transform(df[['transaction_amount']])
21 df.reset_index(drop=True, inplace=True)
22 print("Uncleaned Dataset:")
23 print(df)
```

OUTPUT:

```
PS C:\Users\DEEKSHA\Desktop\AIAC\Lab-17> & C:/Users/DEEKSHA/AppData/Local/Microsoft/Windows/PowerShell/6.2/Task2.py
Original Dataset:
   transaction_id  customer_id transaction_date  transaction_amount      product
0              1          101    2025-10-01            250  Notebook
1              2          102    2025-10-03             0       Pen
2              3          103    2025-10-05            -50    Marker
3              4          104    2025-11-01            450  Printer
4              5          105    2025-11-02            300     Paper
5              6          106    2025-11-03            500   Laptop
6              7          107  invalid_date            200     Mouse
7              8          108    2025-11-05             0  Keyboard
8              9          109    2025-11-06            150  Stapler
9             10          110    2025-11-07            600     Desk

Cleaned Dataset:
   transaction_id  customer_id transaction_date  transaction_amount      product month_year  transaction_amount_normalized
0              1          101    2025-10-01            250  Notebook  2025-10           0.222222
1              4          104    2025-11-01            450  Printer  2025-11           0.666667
2              5          105    2025-11-02            300     Paper  2025-11           0.333333
3              6          106    2025-11-03            500   Laptop  2025-11           0.777778
4              9          109    2025-11-06            150  Stapler  2025-11           0.000000
5             10          110    2025-11-07            600     Desk  2025-11           1.000000
```

PS C:\Users\DEEKSHA\Desktop\AIAC\Lab-17>

	<p>Task 3 – Healthcare Patient Records Cleaning</p> <p>Task:</p> <p>Use AI to generate a script for cleaning healthcare patient records.</p> <p>Instructions:</p> <ul style="list-style-type: none"> 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. <p>Expected Output:</p> <ul style="list-style-type: none"> A cleaned healthcare dataset suitable for ML model training. <p>CODE:</p> <pre>import pandas as pd import numpy as np data = { 'patient_id': [1,2,3,4,5], 'name': ['Alice', 'Bob', 'Charlie', 'David', 'Eva'], 'gender': ['F', 'M', 'male', 'Female', 'M'], 'age': [25, 30, np.nan, 45, 50], 'height_cm': [160, 175, 180, np.nan, 165], 'weight_kg': [55, 70, 80, 90, np.nan], 'blood_pressure': [120, 130, np.nan, 140, 135], 'heart_rate': [80, np.nan, 78, 85, 90] } df = pd.DataFrame(data) print("Original Dataset:") print(df) numeric_cols = ['age', 'height_cm', 'weight_kg', 'blood_pressure', 'heart_rate'] for col in numeric_cols: df[col].fillna(df[col].mean(), inplace=True) df['height_m'] = df['height_cm'] / 100 df.drop(columns=['height_cm'], inplace=True) # drop old column df['gender'] = df['gender'].replace({'M':'Male', 'male':'Male', 'F':'Female', 'Female':'Female'}) df.drop(columns=['patient_id'], inplace=True) print("\nCleaned Healthcare Dataset:") print(df) df.to_csv('healthcare_data_cleaned.csv', index=False)</pre> <p>OUTPUT:</p> <pre>Cleaned Healthcare Dataset: name gender age weight_kg blood_pressure heart_rate height_m 0 Alice Female 25.0 55.00 120.00 80.00 1.60 1 Bob Male 30.0 70.00 130.00 83.25 1.75 2 Charlie Male 37.5 80.00 131.25 78.00 1.80 3 David Female 45.0 90.00 140.00 85.00 1.70 4 Eva Male 50.0 73.75 135.00 90.00 1.65 PS C:\Users\DEEKSHA\Desktop\AIAC\Lab-17></pre>	
	<p>Task 4 – Social Media Sentiment Dataset Preparation</p> <p>Task:</p>	

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.

Expected Output:

- A processed dataset with clean text, ready for NLP sentiment analysis.

CODE:

```
• Task4.py > ...
1  import pandas as pd
2  import re
3  import nltk
4  from nltk.corpus import stopwords
5  from nltk.stem import WordNetLemmatizer
6
7  nltk.download('punkt')
8  nltk.download('stopwords')
9  nltk.download('wordnet')
10
11 data = {
12     'post_id': [1,2,3,4,5],
13     'text': [
14         "I love this product! Check it out: https://example.com",
15         "Worst service ever!! #fail",
16         "Not bad, could be better... ",
17         "Totally amazing experience!!! @brand",
18         "Meh... it's okay. http://test.com "
19     ]
20 }
21 df = pd.DataFrame(data)
22 print("Original Dataset:")
23 print(df)
24
25 def clean_text(text):
26     # Remove URLs
```

```

text = re.sub(r'http\S+|www.\S+', '', text)
# Remove mentions and hashtags
text = re.sub(r'@\w+|\#\w+', '', text)
# Remove emojis and special characters
text = re.sub(r'[^\\w\\s]', '', text)
# Convert to lowercase
text = text.lower()
return text

df['clean_text'] = df['text'].apply(clean_text)

stop_words = set(str(stopwords.words('english')))
df['tokens'] = df['clean_text'].apply(lambda x: [word for word in x.split() if word not in stop_words])

lemmatizer = WordNetLemmatizer()
df['lemmatized'] = df['tokens'].apply(lambda x: [lemmatizer.lemmatize(word) for word in x])

df['processed_text'] = df['lemmatized'].apply(lambda x: ' '.join(x))

print("\nProcessed Dataset:")
print(df[['post_id','processed_text']])

df.to_csv('social_media_cleaned.csv', index=False)

```

OUTPUT:

Original Dataset:

	post_id	text
0	1	I love this product! Check it out: https://ex...
1	2	Worst service ever!! #fail

[nltk_data] C:\Users\DEEKSHA\AppData\Roaming\nltk_data...

Original Dataset:

	post_id	text
0	1	I love this product! Check it out: https://ex...
1	2	Worst service ever!! #fail
	post_id	text
0	1	I love this product! Check it out: https://ex...
1	2	Worst service ever!! #fail
0	1	I love this product! Check it out: https://ex...
1	2	Worst service ever!! #fail
1	2	Worst service ever!! #fail
2	3	Not bad, could be better...
2	3	Not bad, could be better...
3	4	Totally amazing experience!!! @brand
3	4	Totally amazing experience!!! @brand
4	5	Meh... it's okay. http://test.com
4	5	Meh... it's okay. http://test.com

Processed Dataset:

	post_id	processed_text
0	1	love product check
1	2	worst service ever
2	3	bad could better
3	4	totally amazing experience
4	5	meh okay

PS C:\Users\DEEKSHA\Desktop\AIAC\Lab-17>

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).

Expected Output:

- A feature-engineered DataFrame with new indicators and normalized values for ML tasks.

CODE:

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import StandardScaler, OneHotEncoder

# Example: Load your financial dataset
#df = pd.read_csv('financial_data.csv')

# For demonstration, create a mock financial dataset
data = {
    'date': pd.date_range("2023-01-01", periods=60),
    'company_name': ['AlphaTech']*20 + ['BetaBank']*20 + ['GammaEnergy']*20,
    'sector': ['Technology']*20 + ['Finance']*20 + ['Energy']*20,
    'stock_price': np.random.normal(100, 15, 60),
    'volume': np.random.normal(10000, 2000, 60)
}
df = pd.DataFrame(data)

# Randomly set some stock_price and volume as missing
np.random.seed(42)
missing_idx = np.random.choice(df.index, size=8, replace=False)
df.loc[missing_idx[:4], 'stock_price'] = np.nan
df.loc[missing_idx[4:], 'volume'] = np.nan

# 1. Handle missing values in stock_price and volume
df['stock_price'] = df['stock_price'].fillna(df['stock_price'].mean())
df['volume'] = df['volume'].fillna(df['volume'].median())

# 2. Create new features: moving averages (7-day, 30-day) by company
df.sort_values(['company_name', 'date'], inplace=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 using StandardScaler
continuous_columns = ['stock_price', 'volume', 'ma_7', 'ma_30']
scaler = StandardScaler()
df[continuous_columns] = scaler.fit_transform(df[continuous_columns])

# 4. Encode categorical columns
categorical_cols = ['sector', 'company_name']
encoder = OneHotEncoder(sparse_output=False, drop='first')
encoded = encoder.fit_transform(df[categorical_cols])
encoded_cols = encoder.get_feature_names_out(categorical_cols)
df_encoded = pd.DataFrame(encoded, columns=encoded_cols, index=df.index)

# Merge encoded columns and drop originals
df = df.drop(columns=categorical_cols).join(df_encoded)

# Final feature-engineered DataFrame ready for ML tasks
print(df.head())

# Optionally save to new CSV
df.to_csv('financial_data_feature_engineered.csv', index=False)
```

OUTPUT:

```
PS C:\Users\DEEKSHA\Desktop\AIAC\Lab-1> & c:/Users/DEEKSHA/AppData/Local/Microsoft/WindowsApps/python3.11.exe c:/Users/DEEKSHA/Desktop/AIAC\Lab-1\predictor.py
      date stock_price volume   ma_7   ma_30 sector_Finance sector_Technology company_name_BetaBank company_name_GammaEnergy
0 2023-01-01     0.000000 -1.388626 -0.023359  0.073102      0.0        1.0          0.0            0.0               0.0
1 2023-01-02     1.141799 -1.870012  1.138647  1.642710      0.0        1.0          0.0            0.0               0.0
2 2023-01-03    -0.107236  1.585444  0.678556  1.021231      0.0        1.0          0.0            0.0               0.0
3 2023-01-04     0.076179 -0.697816  0.541841  0.836560      0.0        1.0          0.0            0.0               0.0
4 2023-01-05    -1.296451 -0.532675 -0.098957 -0.029014      0.0        1.0          0.0            0.0               0.0
```

 Deliverables (For All Tasks)

1. AI-generated prompts for code and test case generation.
2. At least 3 assert test cases for each task.
3. AI-generated initial code and execution screenshots.
4. Analysis of whether code passes all tests.
5. Improved final version with inline comments and explanation.
6. Compiled report (Word/PDF) with prompts, test cases, assertions, code, and output.