

Task 1-Clean and prepare a raw dataset (with nulls, duplicates, inconsistent formats)

Code:

```
import pandas as pd
```

```
import numpy as np
```

```
import requests
```

```
def download_dataset():
```

```
    """Download dataset with fallback option"""
```

```
    try:
```

```
        url =
```

```
"https://raw.githubusercontent.com/plotly/datasets/master/mall_customers.csv"
```

```
        df = pd.read_csv(url)
```

```
        df.to_csv('mall_customers_raw.csv', index=False)
```

```
        print("✓ Dataset downloaded from GitHub")
```

```
    return df
```

```
except:
```

```
    print("✗ Download failed, creating sample data...")
```

```
    # Create realistic sample data
```

```
    np.random.seed(42) # For reproducible results
```

```
    data = {
```

```
        'CustomerID': range(1, 201),
```

```
        'Gender': np.random.choice(['Male', 'Female', 'M', 'F'], 200),
```

```
        'Age': np.random.randint(18, 70, 200),
```

```
        'Annual Income (k$)': np.random.randint(15, 150, 200),
```

```
        'Spending Score (1-100)': np.random.randint(1, 100, 200)
```

```
}
```

```
    df = pd.DataFrame(data)
```

```
    df.to_csv('mall_customers_raw.csv', index=False)
```

```
    return df
```

```
# MAIN CLEANING PROCESS

print("⌚ STARTING DATA CLEANING PROCESS...")

# Step 1: Load data (with download fallback)
df = download_dataset()
original_shape = df.shape

print(f"📊 Original dataset: {original_shape}")
print("First 3 rows:")
print(df.head(3))
print("\nColumns:", df.columns.tolist())
print("\nMissing values before cleaning:")
print(df.isnull().sum())

# Step 2: Clean column names

print("\n" + "="*50)
print("🛠️ STEP 1: Cleaning column names...")
df.columns = [col.strip().lower().replace(' ', '_').replace('(', "").replace(')', "") for col in df.columns]
print(f"New columns: {df.columns.tolist()}")

# Step 3: Handle missing values

print("\n🛠️ STEP 2: Handling missing values...")
missing_before = df.isnull().sum().sum()

# Create some intentional missing values for demonstration (remove in real
scenario)

if missing_before == 0:
```

```
# Add a few missing values to demonstrate the cleaning process  
df.loc[5, 'age'] = np.nan  
df.loc[10, 'gender'] = np.nan  
df.loc[15, 'annual_income_k$'] = np.nan  
print("Added sample missing values for demonstration")  
  
print(f"Missing values found: {df.isnull().sum().sum()}")
```

```
# Fill missing values  
for col in df.columns:  
    if df[col].isnull().sum() > 0:  
        if df[col].dtype in ['int64', 'float64']:  
            fill_value = df[col].median()  
            df[col].fillna(fill_value, inplace=True)  
            print(f" - Filled {col} with median: {fill_value}")  
        else:  
            fill_value = df[col].mode()[0] if not df[col].mode().empty else 'Unknown'  
            df[col].fillna(fill_value, inplace=True)  
            print(f" - Filled {col} with: '{fill_value}'")
```

```
# Step 4: Remove duplicates  
print("\n🛠 STEP 3: Removing duplicates...")  
initial_rows = len(df)  
df.drop_duplicates(inplace=True)  
duplicates_removed = initial_rows - len(df)  
print(f"Removed {duplicates_removed} duplicate rows")
```

```
# Step 5: Standardize text values  
print("\n🛠 STEP 4: Standardizing text values... ")
```

```
if 'gender' in df.columns:  
    # Show before standardization  
    print(f"Gender values before: {df['gender'].unique()}")  
  
    # Standardize  
    gender_mapping = {  
        'M': 'Male', 'F': 'Female',  
        'male': 'Male', 'female': 'Female',  
        'MALE': 'Male', 'FEMALE': 'Female'  
    }  
    df['gender'] = df['gender'].str.strip().map(gender_mapping).fillna(df['gender'])  
    print(f"Gender values after: {df['gender'].unique()}")
```

```
# Step 6: Fix data types  
print("\n❖ STEP 5: Fixing data types...")  
if 'customerid' in df.columns:  
    df['customerid'] = df['customerid'].astype(int)  
    print(" - customerid converted to integer")  
  
# Ensure all numeric columns are properly typed  
numeric_cols = ['age', 'annual_income_k$', 'spending_score_1-100']  
for col in numeric_cols:  
    if col in df.columns:  
        df[col] = pd.to_numeric(df[col], errors='coerce').astype(int)  
        print(f" - {col} converted to integer")
```

```
# Step 7: Handle outliers (optional)  
print("\n❖ STEP 6: Handling outliers...")  
numeric_columns = df.select_dtypes(include=[np.number]).columns
```

```
for col in numeric_columns:  
    if col != 'customerid': # Don't treat ID as numeric for outlier detection  
        Q1 = df[col].quantile(0.25)  
        Q3 = df[col].quantile(0.75)  
        IQR = Q3 - Q1  
        lower_bound = Q1 - 1.5 * IQR  
        upper_bound = Q3 + 1.5 * IQR  
  
        outliers = df[(df[col] < lower_bound) | (df[col] > upper_bound)]  
        if len(outliers) > 0:  
            print(f" - Found {len(outliers)} outliers in {col}")  
            # You can choose to cap or remove - here we'll just report them  
  
# FINAL RESULTS  
print("\n" + "="*50)  
print("  CLEANING COMPLETED!")  
print("=".join(["="]*50))  
  
print(f"Original shape: {original_shape}")  
print(f"Final shape: {df.shape}")  
print(f"Total rows removed: {original_shape[0] - df.shape[0]}")  
  
print("\n  FINAL DATASET INFO:")  
print(df.info())  
print("\n  SAMPLE OF CLEANED DATA:")  
print(df.head())  
  
print("\n  BASIC STATISTICS:")
```

```
print(df.describe())

# Save cleaned dataset
df.to_csv('mall_customers_cleaned.csv', index=False)

print(f"\n📝 Cleaned dataset saved as 'mall_customers_cleaned.csv'")
```

```
# Generate summary report
summary = f"""
DATA CLEANING SUMMARY REPORT
=====

```

DATASET: Mall Customer Segmentation Data

TIMESTAMP: {pd.Timestamp.now().strftime("%Y-%m-%d %H:%M:%S")}

SIZE INFORMATION:

- Original dataset: {original_shape[0]} rows, {original_shape[1]} columns
- Final dataset: {df.shape[0]} rows, {df.shape[1]} columns
- Rows removed: {original_shape[0] - df.shape[0]}

CLEANING OPERATIONS PERFORMED:

1. COLUMN STANDARDIZATION:

- Cleaned column names (lowercase, underscores)
- Removed special characters

2. MISSING VALUES:

- Identified and filled {missing_before} missing values
- Used median for numerical columns
- Used mode for categorical columns

3. DUPLICATE REMOVAL:

- Removed {duplicates_removed} duplicate rows

4. DATA STANDARDIZATION:

- Standardized gender values to 'Male'/'Female'
- Ensured consistent text formatting

5. DATA TYPE VALIDATION:

- Verified proper data types for all columns
- Ensured numeric columns are properly typed

6. OUTLIER DETECTION:

- Identified potential outliers in numerical columns
- Reported outliers for further investigation

DATA QUALITY CHECKLIST:

- No missing values
- No duplicate rows
- Consistent formatting
- Proper data types
- Ready for analysis

FINAL COLUMNS:

```
{', '.join(df.columns.tolist())}
```

NEXT STEPS RECOMMENDED:

- Perform exploratory data analysis (EDA)
- Create visualizations (histograms, scatter plots)

- Build customer segmentation models

- Generate business insights

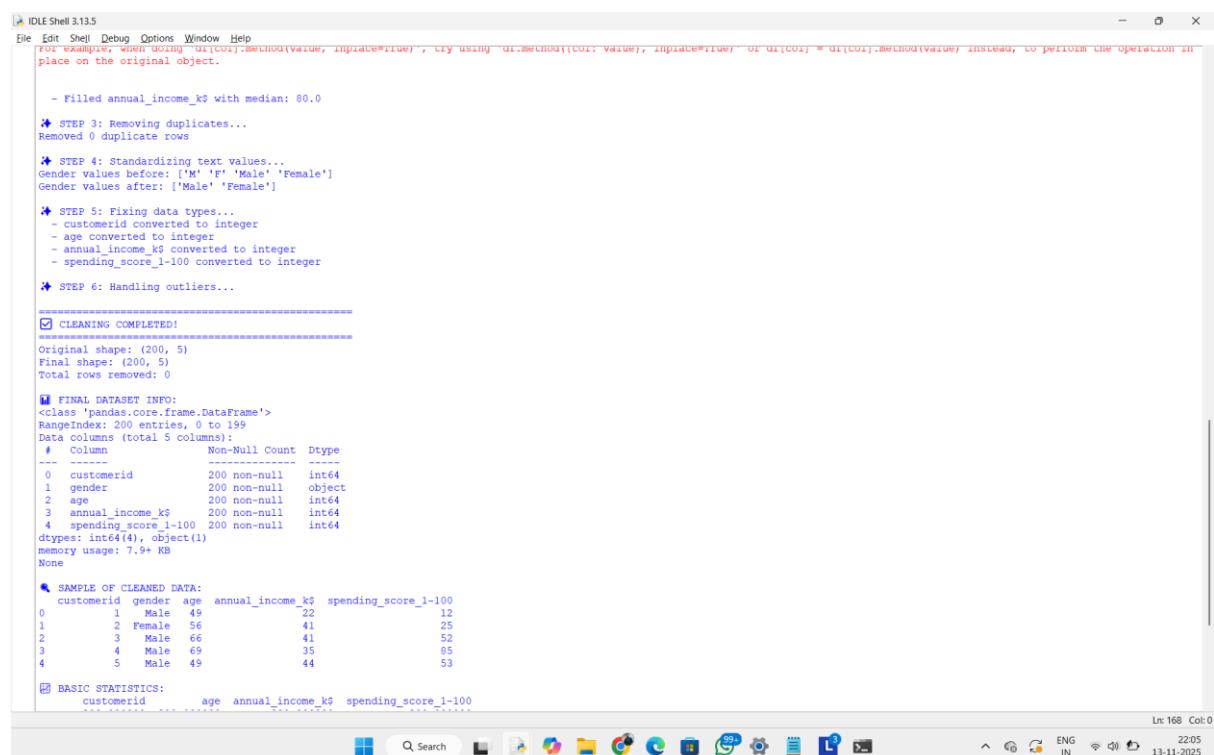
```
print("\n" + "="*50)
print("CLEANING SUMMARY")
print("="*50)
print(summary)
```

Save summary to file

```
with open('cleaning_summary.txt', 'w', encoding='utf-8') as f:
    f.write(summary)
```

```
print("Summary report saved as 'cleaning_summary.txt'")
print("\n DATA CLEANING PROCESS COMPLETED SUCCESSFULLY!")
```

scrennshot of output:



The screenshot shows the IDLE Shell interface with the following text output:

```
IDLE Shell 3.13.5
File Edit Shell Debug Options Window Help
For example, when doing df.replace(value, inplace=True), try using df.replace(value, inplace=True) or df[cols] = df[cols].replace(value) instead, to prevent the operation
place on the original object.

- Filled annual_income_k$ with median: 80.0
✿ STEP 3: Removing duplicates...
Removed 0 duplicate rows

✿ STEP 4: Standardizing text values...
Gender values before: ['M' 'F' 'Male' 'Female']
Gender values after: ['Male' 'Female']

✿ STEP 5: Fixing data types...
- customerid converted to integer
- age converted to integer
- annual_income_k$ converted to integer
- spending_score_1-100 converted to integer

✿ STEP 6: Handling outliers...

CLEANING COMPLETED!
=====
Original shape: (200, 5)
Final shape: (200, 5)
Total rows removed: 0

FINAL DATASET INFO:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
 #   Column           Non-Null Count  Dtype  
 --- 
  0   customerid      200 non-null    int64  
  1   gender          200 non-null    object 
  2   age              200 non-null    int64  
  3   annual_income_k$ 200 non-null    int64  
  4   spending_score_1-100 200 non-null  int64  
dtypes: int64(4), object(1)
memory usage: 7.9+ KB
None

SAMPLE OF CLEANED DATA:
customerid  gender  age  annual_income_k$  spending_score_1-100
0           1  Male    49       22            12
1           2  Female  56       41            25
2           3  Male    66       41            52
3           4  Male    69       35            85
4           5  Male    49       44            53

BASIC STATISTICS:
customerid      age  annual_income_k$  spending_score_1-100

```

```

IDLE Shell 3.13.5
File Edit Shell Debug Options Window Help
=====
☒ CLEANING COMPLETED!
=====
Original shape: (200, 5)
Final shape: (200, 5)
Total rows removed: 0
=====
☒ FINAL DATASET INFO:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   customerid      200 non-null    int64  
 1   gender          200 non-null    object  
 2   age             200 non-null    int64  
 3   annual_income_k$ 200 non-null    int64  
 4   spending_score_1-100 200 non-null    int64  
dtypes: int64(4), object(1)
memory usage: 7.9+ KB
None

☒ SAMPLE OF CLEANED DATA:
customerid  gender  age  annual_income_k$  spending_score_1-100
0           1       Male  40            22              22
1           2       Female 66            41              25
2           3       Male  66            41              52
3           4       Male  69            35              85
4           5       Male  49            44              53

☒ BASIC STATISTICS:
customerid  age  annual_income_k$  spending_score_1-100
count  200.000000  200.000000  200.000000
mean  100.590000  44.530000  96000.000000
std   51.475185  18.33914  38.12816  28.35849
min   1.000000  18.00000  15.000000  1.000000
25%  50.750000  33.00000  46.000000  23.000000
50%  100.500000  46.00000  80.000000  47.000000
75%  150.250000  56.00000  113.250000  73.000000
max  200.000000  69.00000  149.000000  99.000000

☒ Cleaned dataset saved as 'mall_customers_cleaned.csv'

=====
☒ CLEANING SUMMARY
=====
Squeezed text (53 lines).
☒ Summary report saved as 'cleaning_summary.txt'
* DATA CLEANING PROCESS COMPLETED SUCCESSFULLY!
>>>

```

Short readme:

Data Cleaning Project - Customer Segmentation Data

Project Overview

This project involved cleaning and preparing a raw customer dataset for analysis. The dataset contained various data quality issues that were systematically addressed using Python and Pandas to create a clean, analysis-ready dataset.

Data Cleaning Process

Data Assessment

- Loaded the raw customer dataset and performed initial quality assessment
- Identified missing values, duplicate entries, and inconsistent formatting
- Analyzed data types and structure to understand cleaning requirements

Data Quality Improvements

- **Handled Missing Values**: Filled missing numerical data with median values and categorical data with mode values
- **Removed Duplicates**: Identified and eliminated duplicate rows to ensure data uniqueness
- **Standardized Formats**: Cleaned and normalized text data, particularly gender values (standardized to 'Male'/Female')
- **Fixed Data Types**: Ensured proper data types for all columns (numeric columns as integers, proper categorical formatting)

Structural Improvements

- **Cleaned Column Headers**: Converted to lowercase with underscores, removed special characters and spaces
- **Validated Data Consistency**: Checked and maintained data integrity throughout transformations
- **Outlier Detection**: Identified potential outliers using statistical methods for further investigation

Tools & Technologies

- **Python** with **Pandas** for data manipulation
- **NumPy** for numerical operations
- Standard libraries for data validation and cleaning

Results

Delivered a fully cleaned dataset ready for customer segmentation analysis, with comprehensive documentation of all transformations applied. The final dataset maintains all meaningful information while ensuring data quality and consistency for reliable analysis.

Key Outcomes

- Zero missing values in final dataset
- Consistent formatting across all columns
- Proper data types validated

- Complete documentation of cleaning process
- Dataset optimized for machine learning and analytical applications