



Case Study of Data Mining Code-24CAT-611

Data Cleaning in E-Commerce Customer Analytics

Submitted to:

Dr.Arun kumar

Professor

Employee ID: E17568

Submitted by:

Pragati Gupta

UID:24MCI10255

Problem statement

An e-commerce company, XYZ, aims to enhance its customer experience and improve sales through data-driven insights. The company collects vast amounts of data from customer transactions, website interactions, and feedback surveys. However, as the data grows, inconsistencies, duplicates, and errors have emerged, leading to unreliable analytics.

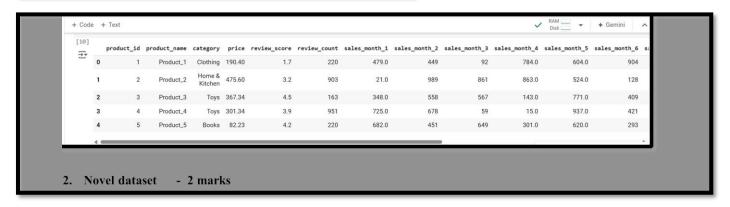
Take dataset from any source, clean it by applying different data cleaning techniques whatever required into that along with the explanation.

GitHub Link for this Case study:

Rubrics for evaluation:

1. Relevant dataset _ 2 marks

```
+ Code + Text
 D import numpy as np
      import pandas as pd
      import matplotlib.pyplot as plt
     import seaborn as sns
     # 1, Relevant dataset:
     df-pd.read_csv('/content/econnerce_analysis.csv')
     display(df.head())
     product id - df| 'product id' |
     product name - df['product name']
     category - dt["category"]
     price - df['price']
     review_score - df['review_score']
     review_count - df['review_count']
      sales_month_1 - df['sales_month_1']
      sales_month_2 - df['sales_month_2']
      sales_month_3 - df['sales_month_3']
      sales_month_4 = df['sales_month_4']
      sales_month_5 - df['sales_month_5']
      sales month 6 - df['sales month 6']
      sales month 7 - df['sales month 7']
      sales month 8 - df['sales_month_8']
      sales_month_9 - df['sales_month_9']
      sales month 10 - df['sales month 10']
      sales_month_11 - df['sales_month_11'
      sales month 12 - df['sales month 12'
```



2. Novel dataset - 2 marks

```
File Edit View Insert Runtime Tools Help Allchangessaved

+ Code + Text

** 2. Novel this dataset

** Explore basic statistics and data types

print(f"Overview of Data: \n (df.info())")  # Get an overview of the data, including column types and missing values

print(f"Descriptive Statistics of Columns: \n(df.describe())")  # Show descriptive statistics for numerical columns
```



```
Insert Runtime
                         Tools Help
  + Text
Descriptive Statistics of Columns:
        product_id
                          price review_score
                                                review_count
                                                               sales_month_1 \
count
       1000.000000
                    1000.000000
                                  1000,000000
                                                 1000.000000
                                                                 1000.000000
        500.500000
                     247.677138
                                      3.027600
                                                  526,586888
                                                                  498.814000
mean
                     144.607983
std
        288.819436
                                      1.171243
                                                  282.269932
                                                                  289.512243
          1.000000
                                                                    0.000000
min
                       7.290000
                                      1.000000
                                                    1.000000
        250.750000
                                      2.000000
                                                  283.750000
                                                                  245.000000
25%
                     121.810000
50%
        500,500000
                     258,920000
                                      3,100000
                                                  543,000000
                                                                  508,000000
75%
                     373.435000
                                      4,000000
                                                  772.000000
                                                                  740,758000
        750,250000
max
       1000.000000
                     499.860000
                                      5.886088
                                                  999.000000
                                                                 1000,000000
       sales_month_2
                      sales_month_3
                                      sales_month_4 sales_month_5 \
count
         1000.000000
                         1080,000000
                                         999.000000
                                                        999.000000
          507.661000
                         586.739886
                                         504.324324
                                                         487.674675
mean
          285.992689
                                         286.350194
                                                        287.586675
std
                         294.010873
            2.000000
                           0.000000
                                           0.000000
                                                          0.000000
min
25%
          262.500000
                         243,750000
                                         263.000000
                                                        222.000000
50%
          508,000000
                         493,000000
                                         502,600000
                                                        497,000000
75%
                                         750,000000
                                                        727,000000
          756,250000
                         777.250000
                                        1800,000000
max
         1000,000000
                         999.000000
                                                       1000,000000
       sales_month_6 sales_month_7
                                      sales_month_8
                                                      sales_month_9
         1000.000000
                         1000,000000
                                         999,000000
                                                       1000.000000
mean
          491.653000
                          507.011000
                                         505.022022
                                                        491.934000
          289.234018
                         291.047287
                                         289.736615
                                                        287,514731
std
            0.000000
                           0.000000
min
                                           5.000000
                                                          0.000000
          236.000000
                          254,886886
                                         241.500000
                                                         247.250000
25%
50%
          479,588688
                         522,566666
                                         500,000000
                                                        495,586688
          740,500000
                         757,250000
                                         762,500000
                                                        735,250000
75%
max
         1000,000000
                         1000.000000
                                        1000,000000
                                                       1000,000000
```

File Edit View Insert Runtime Tools Help All changes saved										
+	Cod	e + Te	ext							
/ [26]		sales_month_10	sales_month_11	sales_month_12					
	₹	count	1000.000000	1000.00000	1000.000000					
	Ξ,	mean	514.798000	505.83800	500.386000					
		std	288.710119	288.82451	278.509459					
		min	1.000000	0.00000	4.000000					
		25%	267.000000	251.25000	259.000000					
		50%	532.000000	502.00000	500.500000					
		75%	770.250000	761.00000	730.000000					
		max	1000.000000	1000.00000	1000.000000					

- 3. Data cleaning techniques along with the relevant explanation that why these techniques are being applied. -- 3 marks
 - a) **Identifying missing values:** The complete Dataset is checked if there is the presence of any null value or not.
 - b) **Checking for the duplicate rows:** The dataset is checked for any kind of duplicate row present in the dataset.

- c) Replaced Missing values with a central tendency, i.e., Median: One of the suitable method to handle Missing value is to replace it with the central tendencies like mean, median, mode, or standard deviation.
- d) **Detecting the outliers:** Outliers refers to value out of range with respect to the dataset available. We have detected outliers based on their IQR (Inter quartile Range).



e) **Replacing the outliers with Median:** Outliers detected by the IQR are replaced with the central tendency in order to clean the dataset.

```
Code + Text
    # 3. Data cleaning in the dataset
    # Identifying Missing Values in the dataset
    print(f"Checking for any missing value present:\n {df.isnull().sum()}") # Count missing values in each column
    # Checking for duplicate rows
    print(f"Checking for any duplicate rows present: {df.duplicated().sum()}")
    # Replacing missing value in "sales_month_1" with the median value
    median_sales_month_1 = df['sales_month_1'].median()
    df['sales_month_1'].fillna(median_sales_month_1, inplace=True)
    print(f"Replaced missing value: {median_sales_month_1}")
    # Detection of outliers in "price" column
    Q1 = df['price'].quantile(0.25)
    Q3 = df['price'].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers = df[(df['price'] < lower_bound) | (df['price'] > upper_bound)]
    print("Outliers in 'price' column:")
    print(outliers)
    # Replacing outliers with the median value
    median_price = df['price'].median()
    df.loc[df['price'] < lower_bound, 'price'] = median_price</pre>
    df.loc[df['price'] > upper_bound, 'price'] = median_price
    print(f"Replaced outliers: {median_price}")
```

4. Relevant graphs and tables depicting clean data –3 marks

```
File Edit View Insert Runtime Tools Help All changes saved
+ Code + Text
         # 4.Relevant graphs and tables depicting clean data
          # Table showing basic statistics of the cleaned data
print(f"Table of cleaned data:\n {df.describe()}")
           # a) Distribution of Product Prices
          plt.figure(figsize=(8, 6))
sns.histplot(df['price'], bins=20, kde=True, color='y')
plt.title('Distribution of Product Prices')
           plt.xlabel('Price')
           plt.ylabel('Frequency')
           plt.show()
           # b) Relationship between Price and Review Score
          # 0) Relationship between Price and Review Score
plt.figure(figsize=(8, 6))
sns.scatterplot(x='price', y='review_score', data=df,c=price, cmap='viridis')
plt.title('Relationship between Price and Review Score')
plt.xlabel('Price')
plt.ylabel('Review Score')
           plt.show()
          # c) Top 10 Categories by Sales in Month 1
top_10_categories_month_1 = df.groupby('category')['sales_month_1'].sum().nlargest(10)
plt.figure(figsize=(10, 6))
sns.barplot(x=top_10_categories_month_1.values, y=top_10_categories_month_1.index, color='red')
plt.title('Top 10_categories by Sales in Month 1')
           plt.xlabel('Total Sales')
           plt.ylabel('Category')
           plt.show()
```

```
File Edit View Insert Runtime Tools Help All changes saved
+ Code + Text
[25] # d) Average Review Score by Category
      average_review_score_by_category = df.groupby('category')['review_score'].mean()
      plt.figure(figsize=(10, 6))
      sns.barplot(x=average review_score_by_category.values, y=average review_score_by_category.index, color='k')
      plt.title('Average Review Score by Category')
      plt.xlabel('Average Review Score')
      plt.ylabel('Category')
      plt.show()
      # e) Correlation Matrix
      correlation matrix = df[['price', 'review score', 'review count', 'sales month 1']].corr()
      plt.figure(figsize=(8, 6))
      sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
      plt.title('Correlation Matrix')
      plt.show()
      # f) Boxplot of Sales in Month 1 by Category
      plt.figure(figsize=(10, 6))
      sns.boxplot(x='category', y='sales_month_1', data=df, color='cyan')
      plt.title('Boxplot of Sales in Month 1 by Category')
      plt.xlabel('Category')
      plt.ylabel('Sales in Month 1')
      plt.xticks(rotation-90)
      plt.show()
```

```
File Edit View Insert Runtime Tools Help All changes saved
+ Code + Text
 Table of cleaned data:
                               price review_score review_count sales_month_1 \
              product id
     count 1000.000000 1000.000000 1000.000000 1000.000000
                                                                 1000.000000
             500.500000 247.677130
                                        3.027600
                                                   526.586888
                                                                 498.814000
     mean
             288.819436
                         144.607983
                                        1.171243
                                                   282.269932
                                                                  289.512243
     std
              1.000000
                           7.290000
                                        1.000000
                                                     1.000000
                                                                    0.000000
     min
     25%
             250.750000
                         121.810000
                                        2.000000
                                                    283.750000
                                                                  246.000000
     50%
             500.500000
                         250.920000
                                        3.100000
                                                    543,000000
                                                                  508.000000
     75%
             750.250000 373.435000
                                        4,000000
                                                   772.000000
                                                                 740,750000
            1000.000000 499.860000
                                        5.000000
                                                   999.000000
                                                                 1000.000000
     max
            sales_month_2 sales_month_3 sales_month_4 sales_month_5 \
     count
             1000.000000 1000.000000
                                          999.000000
                                                         999.000000
               507.661000
                             506,739000
                                           504.324324
                                                          487.674675
     mean
     std
               285.992689
                             294.010873
                                           286.350194
                                                          287.586675
                2.000000
                              0.000000
                                            0.000000
                                                           0.000000
     min
     25%
               262.500000
                             243.750000
                                           263.000000
                                                          222,000000
     50%
               508.000000
                             493.000000
                                           502.000000
                                                          497.000000
     75%
               756.250000
                             777.250000
                                           750.000000
                                                         727,000000
              1000.000000
                             999.000000
                                         1000.000000
                                                       1000.000000
     max
            sales month 6 sales month 7 sales month 8 sales month 9 \
     count
              1000.000000
                            1000.000000
                                           999,000000
                                                        1000,000000
               491.653000
                             507.011000
                                           505.022022
                                                         491.934000
     mean
     std
               289.234018
                            291.047287
                                          289.736615
                                                         287.514731
     min
                 0.000000
                               0.000000
                                             5,000000
                                                           0.000000
     25%
               236.000000
                             254.000000
                                           241.500000
                                                          247.250000
     50%
               479.500000
                             522.500000
                                           500.000000
                                                          495.500000
     75%
               740.500000
                             757,250000
                                           762,500000
                                                         735.250000
              1000.000000
                                           1000.000000
                            1000.000000
                                                         1000.000000
     max
```

File	Edit Vi	ew Insert Runtin	ne Tools Help A	All changes saved	
Cod	le + Te	xt			
0		sales month 10	sales month 11	sales month 12	
_	count	1000.000000	1000.00000	1000.000000	
_	mean	514.798000	505.83800	500.386000	
	std	288.710119	288.82451	278.509459	
	min	1.000000	0.00000	4.000000	
	25%	267.000000	251.25000	259.000000	
	50%	532.000000	502.00000	500.500000	
	75%	770.250000	761.00000	730.000000	
	max	1000.000000	1000.00000	1000.000000	

