

CCA-5: Ensure data integrity and consistency

[Overview](#)

[Objectives](#)

[Steps](#)

[Importance of Data Validity](#)

[Results Obtained](#)

[Next Steps](#)

[Conclusion](#)

Overview

It is very important to check data validity before analysis to ensure that the data which will be used in the analysis is proper and error-free. Data validation helps the analysts to know about errors, inconsistencies, or anomalies which could affect quality analysis. Please document what validation steps you took and the results of these validation steps in the process of finalizing the dataset for this project: it is also helpful if you give a brief note for each step on why you used a given method, or what you learned about the data from your result for each step.

Objectives

- To ensure that the data is free from any errors and is complete.
- To identify and handle the anomalies or inconsistencies that might affect the analysis.
- To ensure that the data transformations and feature engineering steps have been applied correctly.

Steps

1. Initial Data Integrity Checks

- **Objective:** Ensure that the data is complete and free from obvious errors.
- **Actions:**
 - **Missing Values:** Checked for missing values in the dataset.
 - **Data Types:** Ensured correct data types of all columns.
 - **Duplicated Records:** Identified 103 cases where records are duplicated within the dataset. After sufficient discussion and consideration with the team, we decided not to remove this set of duplicated records at this point. Given the timeline and possible impact on the integrity of the dataset, the removal of these records might add unnecessary complexity and risk to the analysis, especially at the stage we are at nearing task submission. Further, these duplicates are retained to allow for more in-depth analysis as all variations possible in customer behavior can be captured.
- **Code Snippet:**

```
1 # Checking for missing values
2 missing_values = df.isnull().sum()
3 print("Missing Values:\n", missing_values)
4
5 # Check data types
6 data_types = df.dtypes
7 print("Data Types:\n", data_types)
8
9 # Checking for duplicates
```

```

10 duplicate_count = df.duplicated().sum()
11 print(f"No of Duplicates: {duplicate_count}")
12

```

2. Consistency and Range Checks

- **Task:** To validate the consistency and range of the data values.
- **Steps Taken:**
 - **Outliers:** Checked for outliers that can significantly affect the analysis.
 - **Range Checks:** Checked that the numerical features fall between the expected ranges.
 - **Category Validation:** Checked that categorical features have valid categories only.
- **Code Snippet:**

```

1  # Find outliers using IQR
2  Q1 = df.quantile(0.25)
3  Q3 = df.quantile(0.75)
4  IQR = Q3 - Q1
5  outliers = ((df < (Q1 - 1.5 * IQR)) | (df > (Q3 + 1.5 * IQR))).sum()
6  print("Outliers found:\n", outliers)
7
8  # Inspect the range of numerical columns
9  for column in numerical_columns:
10     print(f"{column} range: {df[column].min()} to {df[column].max()}")
11
12 # Check for unexpected categories
13 for column in categorical_columns:
14     unique_categories = df[column].unique()
15     print(f"Categories in {column}:\n", unique_categories)
16

```

3. Transform Validation

- **Objective:** To make sure the data transformation has been applied correctly post feature engineering.
- **Steps Performed:**
 - **Scale Validation:** Checked if we had applied feature scaling (Min-Max, Standard) perfectly.
 - **Encoding Validation:** Checked if we had applied the categorical encoding perfectly.
- **Code Snippet:**

```

1  # Checking scaling on one column
2  original_min = df_original['MonthlyCharges'].min()
3  original_max = df_original['MonthlyCharges'].max()
4  scaled_min = df_scaled['MonthlyCharges'].min()
5  scaled_max = df_scaled['MonthlyCharges'].max()
6  print(f"Original range: {original_min} to {original_max}")
7  print(f"Scaled range: {scaled_min} to {scaled_max}")
8
9  # Checking if the encoding is consistent
10 encoded_columns = df_encoded.columns
11 print("Encoded Columns :\n", encoded_columns)
12

```

Importance of Data Validity

- **Data Accuracy:** It allows the analysis to be based on correct and reliable data, hence arriving at dependable insights.

- **Model Performance:** Accurate data makes the model more correct and reduces the chances of inaccuracy in predictions.
- **Compliance:** Maintains the standards of data integrity and keeps the project within the specifications provided for data quality.

Results Obtained

- **Data Integrity:** This has been confirmed, and it is noted that there is no missing value in the dataset, all entries in all features are valid, and duplication of 103 records is retained, assuming a conscious decision to maintain the complexity of the data.
- **Outlier Management:** Detected potential outliers in the data that will lead to distortion in analysis if not treated.
- **Transformation Verification:** Checks whether the data has undergone transformations such as scaling and encoding.

Next Steps

- **Cluster Analyzing:** Now that our data is validated, we can now cluster customer segmentation.
- **Feature Engineering Review:** If something is wrong during clustering, then revisit feature engineering to make sure all the features contribute positively to the analysis.

Conclusion

Data validity is one of the critical processes that ensure accuracy and dependability in our analysis. The decision to keep duplicates was made so that we do not lose richness in data and miss any kind of changing behavior. Now the validated data is ready for the following phase of clustering analysis.