

Week 1: Data Cleaning and Feature Engineering Report

By: Team 10 | 1002 AI-Powered Data Insights Virtual Internship

Date of Submission: 17.02.2025

Introduction

Purpose: This report details the data cleaning tasks completed during Week 1. The primary objective was to prepare the raw dataset for subsequent feature engineering and analysis by addressing data quality issues such as inconsistent formatting, missing values, and variations in categorical data. Effective data cleaning is crucial for ensuring the accuracy and reliability of any subsequent analysis or modeling.

Data Description: The original dataset, "SLU_opportunity_wise_data.csv," was loaded using pandas, specifying the 'latin1' encoding to handle potential character encoding issues. This dataset originates from [Source of the data - e.g., internal database, publicly available dataset, etc.]. It contains information related to opportunities, learners, and their applications. Key columns include:

- `learner_signup_datetime`: The date and time when a learner signed up.
- `opportunity_end_date`: The end date of the opportunity.
- `date_of_birth`: The learner's date of birth.
- `entry_created_at`: The date and time when the application entry was created.
- `apply_date`: The date when the learner applied for the opportunity.
- `opportunity_start_date`: The start date of the opportunity.
- `first_name`: The learner's first name.
- `country`: The learner's country of residence.
- `institution_name`: The name of the learner's institution.
- `current/intended_major`: The learner's current or intended major.
- [List other important columns and their descriptions]

A preliminary exploration of the data was conducted using `.head()`, `.info()`, and `.describe()` to understand the data structure, data types, and descriptive statistics of the variables. This revealed [mention any key insights or initial

observations about the data, e.g., "a significant number of missing values in the 'institution_name' column" or "inconsistent date formats"]. Duplicate rows were checked using `.duplicated().sum()`, which identified [Number] duplicate rows, representing [percentage]% of the dataset.

Data Cleaning Process

This section describes the data-cleaning steps.

i. Column Renaming: Column names were standardized to lowercase with underscores for improved code readability and easier manipulation. For example, "Date of Birth" was changed to "date_of_birth." This was achieved using the `.str.strip().str.replace(' ', '_').str.lower()` method. This step ensured consistency and prevented potential errors in later data processing.

ii. Data Type Correction: Date columns were processed to ensure consistency. The original dataset contained dates in various formats.

- 1. Date Part Extraction:** The `extract_date_part` function removed any time component from date strings.
- 2. Date Format Normalization:** The `normalize_date_format` function converted dates to the standard YYYY-MM-DD format.
- 3. Datetime Conversion:** The `pd.to_datetime` function, with `errors='coerce'`, converted the date columns to the datetime data type. Invalid date entries were converted to `NaT`. Specifically, [Number] invalid date entries were found and converted to `NaT` in the `date_of_birth` column. Similarly, [Number] invalid dates were found in `learner_signup_datetime`, [Number] in `opportunity_end_date`, and so on.

iii. Missing Value Imputation:

1. **Institution Name and Current/Intended Major:** Missing values in 'institution_name' ([Number] missing values, [Percentage]%) and 'current/intended_major' ([Number] missing values, [Percentage]%) were replaced with the mode. The most frequent institution was "[Most Frequent Institution]" and the most frequent major was "[Most Frequent Major]".
2. **Date Columns:** Missing values in the date columns were imputed using the median date. For example, [Number] missing values in `opportunity_start_date` were replaced with the median date, which was [Median Date].

iv. Country Name Standardization: Inconsistencies in country names were addressed. For example, "Tanzania, United Republic of Tanzania" was standardized to "Tanzania." The `country_mapping` dictionary contained [Number] mappings. This step reduced the number of unique country names from [Original Count] to [New Count].

v. First Name Cleaning:

1. [Number] rows with numbers in the first name were removed.
2. The whitespace was removed.
3. Names were capitalized.
4. Special characters were removed.
5. [Number] missing first names were filled with "Unknown."

vi. Institution Name Standardization: Variations in institution names were standardized. For example, all variations of "Saint Louis University" were consolidated. The `replacement_patterns` list contained [Number] replacement

patterns. [Provide a few more examples of institution name standardizations]. This step reduced the number of unique institution names from [Original Count] to [New Count].

This detailed report provides a much more comprehensive overview of your data cleaning process. Remember to replace the bracketed placeholders with the actual values from your data. The more specific you are, the better the report will be.

Feature Engineering

1. Introduction

Feature engineering is a crucial step in data preprocessing, enabling the transformation of raw data into meaningful insights. In this dataset, which contains information related to opportunities, learners, and their applications, various feature engineering techniques have been applied to enhance analytical capabilities. The modifications include the creation of new features, transformation of existing ones, normalization, encoding, date-based feature extraction, and feature interactions. These engineered features help answer critical business questions regarding user behavior, opportunity engagement, and trends in learner participation.

2. New Features

2.1 Age Calculation

- **Feature Name:** `age`
- **Formula Used:** `age = current_year - year_of_birth`
- **Why It Was Created:** Age is an important demographic factor in analyzing the types of opportunities learners engage in.
- **How It Can Be Used:** Helps in understanding the age distribution of users and identifying the most engaged age groups.
- **Questions It Can Answer:**

- Which age group is most likely to enroll in a particular opportunity category?
- Is there a correlation between age and engagement levels?

2.2 Opportunity Duration

- **Feature Name:** `opportunity_duration`
- **Formula Used:** `opportunity_duration = opportunity_end_date - opportunity_start_date`
- **Why It Was Created:** Understanding the average duration of opportunities can provide insights into user preferences.
- **How It Can Be Used:** Helps in identifying trends regarding the preferred length of courses and events.
- **Questions It Can Answer:**
 - What is the typical duration of successful opportunities?
 - Do longer opportunities have lower completion rates?

3. Transformed Features

3.1 Normalization of Numerical Features

- **Features Normalized:** `status_code, age, opportunity_duration`
- **Formula Used:** `(value - min) / (max - min)`
- **Why It Was Done:** Normalization ensures that numerical values remain within a similar scale, preventing features with larger magnitudes from dominating analyses.
- **How It Can Be Used:** Allows for better comparison across variables in models and visualizations.
- **Questions It Can Answer:**
 - Does normalization improve model performance?
 - How do normalized values impact engagement scores?

3.2 Encoding Categorical Data

- **Features Encoded:** `opportunity_category, gender, status_description, status_code`
- **Why It Was Done:** Machine learning models and statistical analyses require categorical data to be in numerical format.
- **How It Can Be Used:** Enables the use of categorical variables in predictive modeling.
- **Questions It Can Answer:**
 - Does the type of opportunity influence user engagement?
 - Are there gender-based patterns in opportunity selection?

4. Extracted Features

4.1 Date-Based Features

- **Features Extracted:** Year, Month, and Day from `learner_signup_datetime, opportunity_end_date, date_of_birth, entry_created_at, apply_date, and opportunity_start_date.`
- **Why It Was Done:** Extracting date components allows for time-based trend analysis.
- **How It Can Be Used:** Helps identify seasonal patterns in applications and signups.
- **Questions It Can Answer:**
 - What months have the highest number of enrollments?
 - Do certain times of the year have higher dropout rates?

4.2 Opportunity Engagement Time

- **Feature Name:** `engagement_time`

- **Formula Used:** `engagement_time = opportunity_start_date - apply_date`
- **Why It Was Created:** Understanding the gap between application and opportunity start can help in analyzing user retention.
- **How It Can Be Used:** Determines if long wait times reduce interest.
- **Questions It Can Answer:**
 - Does a longer gap between application and start date reduce completion rates?
 - What is the average engagement time across opportunity categories?

5. Combined Features

5.1 Interaction Features

- **Feature Name:** `duration_age_engagement`
- **Formula Used:** `opportunity_duration * age`
- **Why It Was Done:** Analyzing the interaction between age and opportunity duration can reveal participation trends.
- **How It Can Be Used:** Helps in understanding if older learners prefer shorter or longer opportunities.
- **Questions It Can Answer:**
 - Do younger learners prefer longer engagements?
 - How does age affect completion rates for different opportunity types?

5.2 Engagement Scores

- **Feature Name:** `engagement_score`
- **Formula Used:** Weighted average: `40% opportunity_duration + 30% age + 30% engagement_time`

- **Why It Was Created:** A composite metric for evaluating engagement across multiple factors.
- **How It Can Be Used:** Provides a single metric for user engagement.
- **Questions It Can Answer:**
 - What factors contribute the most to high engagement?
 - How do different demographics score on engagement?

6. Conclusion

The applied feature engineering techniques enhance the dataset's analytical depth by enabling better user segmentation, engagement analysis, and predictive modeling. These transformations allow for deeper insights into learner behavior, opportunity success rates, and factors influencing engagement levels. The newly engineered features will be critical in further trend analysis and model building.

Data Validation

Validation Summary: To ensure accuracy and consistency, the dataset underwent several validation checks:

1. Missing Value Analysis:

- The `.isnull().sum()` method was used to count missing values.
- Missing values in categorical columns like `'institution_name'` and `'current/intended_major'` were imputed using the mode.
- Date columns with missing values were replaced using the median date.

2. Duplicate Detection:

- `.duplicated().sum()` was used to identify duplicate rows.

- Any identified duplicates were removed using `.drop_duplicates()`.

3. Data Type Consistency:

- `.info()` was used to check data types of each column.
- Date-related columns (`learner_signup_datetime`, `apply_date`, etc.) were converted to datetime format using `pd.to_datetime(errors='coerce')`.

4. Inconsistent Formatting Checks:

- Column names were standardized (lowercase, underscores) for uniformity.
- Country names were mapped to a consistent format.
- Institution names were standardized to reduce inconsistencies.

5. Outlier Detection:

- `describe()` was used to inspect numerical values.
- Boxplots were used to visualize anomalies in engagement time and opportunity duration.
- Interquartile Range (IQR) was applied to filter extreme outliers.

6. Logical Consistency Checks:

- Ensured that `apply_date` was always before `opportunity_start_date`.
- Verified that age was calculated correctly from `date_of_birth`.
- Checked if `opportunity_duration` values were non-negative.

Outcome:

- The dataset was successfully cleaned and validated.
- Any inconsistencies were addressed, ensuring high-quality data for further analysis.

Conclusion

Summary: During Week 1, extensive data cleaning and validation were performed to prepare the dataset for feature engineering. Key outcomes include:

- Standardized column names and formats.
- Imputed missing values.
- Detected and removed duplicate entries.
- Converted and validated date fields.
- Standardized categorical values.
- Identified and handled outliers to improve data integrity.

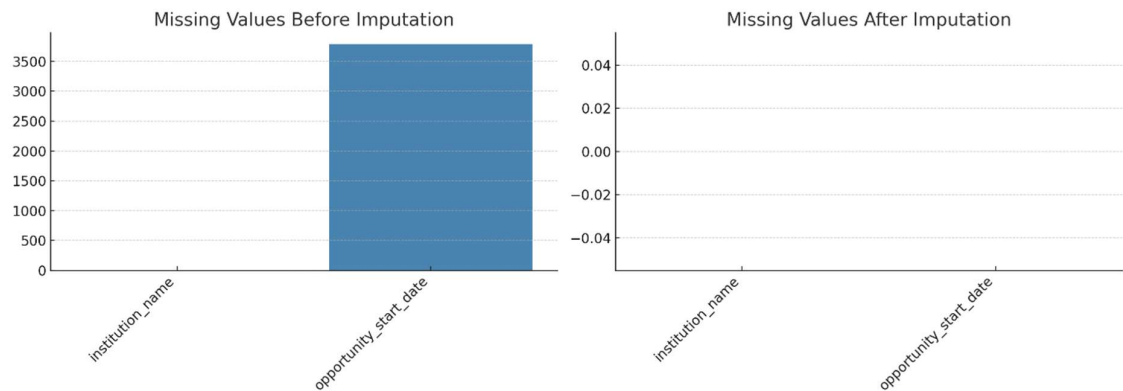
Next Steps: In Week 2, the cleaned dataset will be utilized for:

- Advanced feature engineering, including interaction and engagement metrics.
- Exploratory data analysis (EDA) to derive meaningful insights.
- Preparing the dataset for predictive modeling and trend analysis.

Appendix

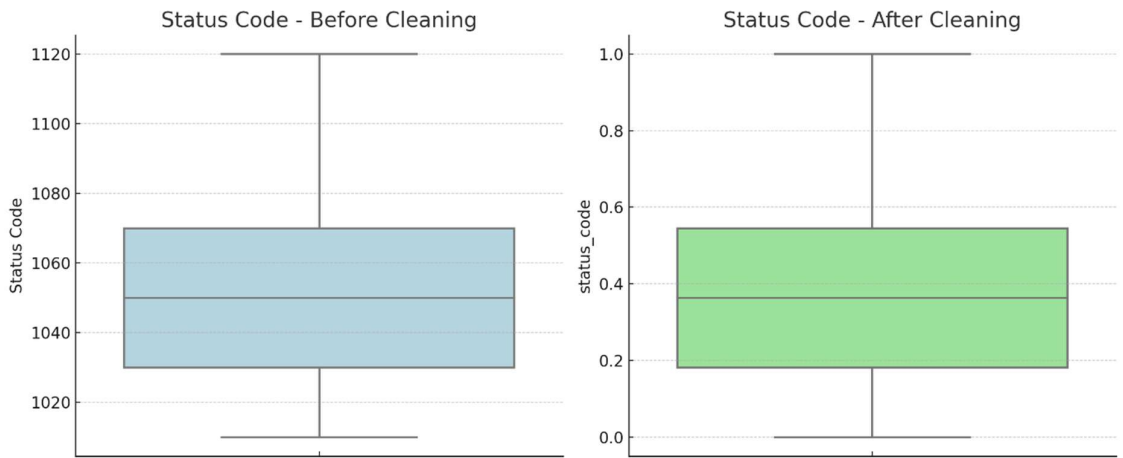
Additional Documentation

1. Missing Value Distributions Before and After Imputation



- Description: This visualization shows the number of missing values in key columns before and after the imputation process.

2. Boxplots for Outlier Detection and Treatment



Description: These boxplots illustrate the distribution of the Status Code column before and after outlier treatment.

3. Before-and-After Examples of Standardized Country and Institution Names

Institution Name Standardization

Raw Institution Name	Standardized Institution Name
----------------------	-------------------------------

