# Data Collection and Preprocessing Phase

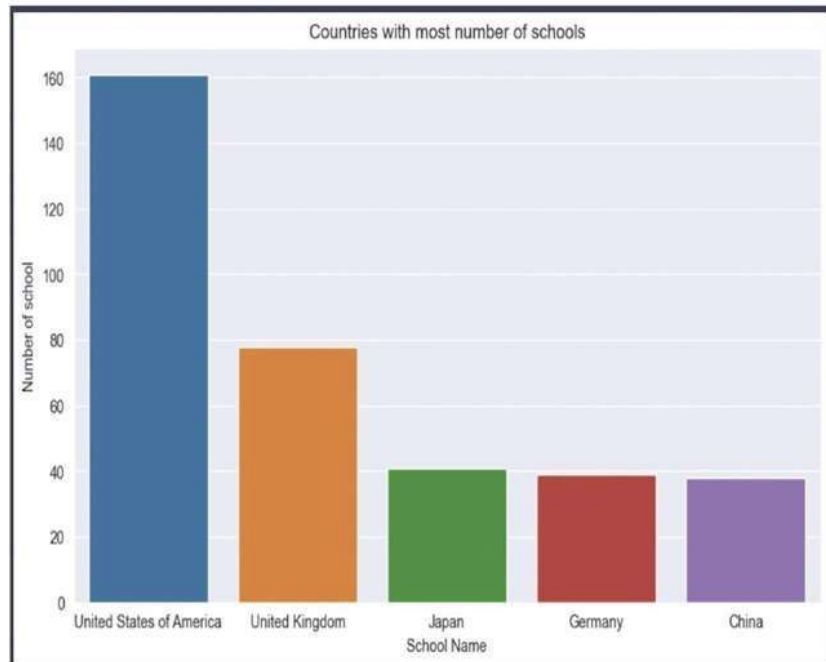| Date | 06 July 2024 |
|------|-------------|
| Team ID | 739824 |
| Project Title | SmartLender – Envisioning Success: Predicting University Scores With Machine Learning |
| Maximum Marks | 6 Marks |

## Data Exploration and Preprocessing Report

Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.
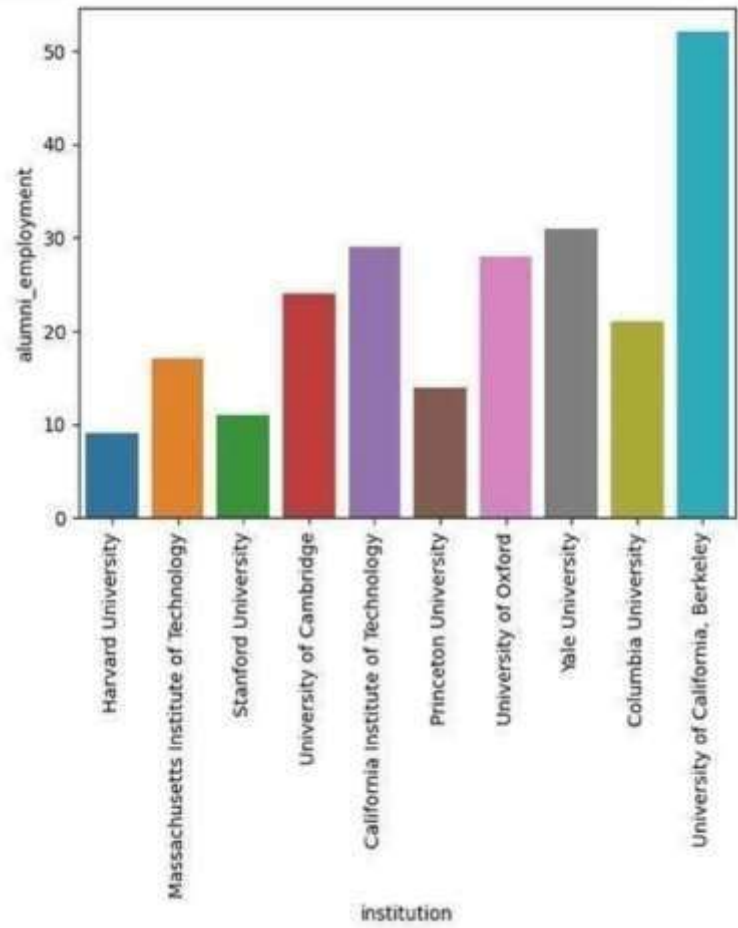
| Section | Description |
|---------|-------------|
| Data Overview |  |

| | |
|---|---|
| Univariate Analysis | *a.Univariate Analysis*<br><br>Countries with most number of schools<br><br> |

**B**ivariate Analysis

| Multivariate Analysis |  |
|---|---|

## Data Preprocessing Code Screenshots

| Loading Data |  |
|---|---|

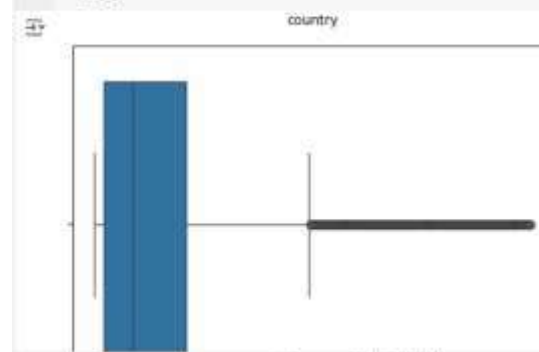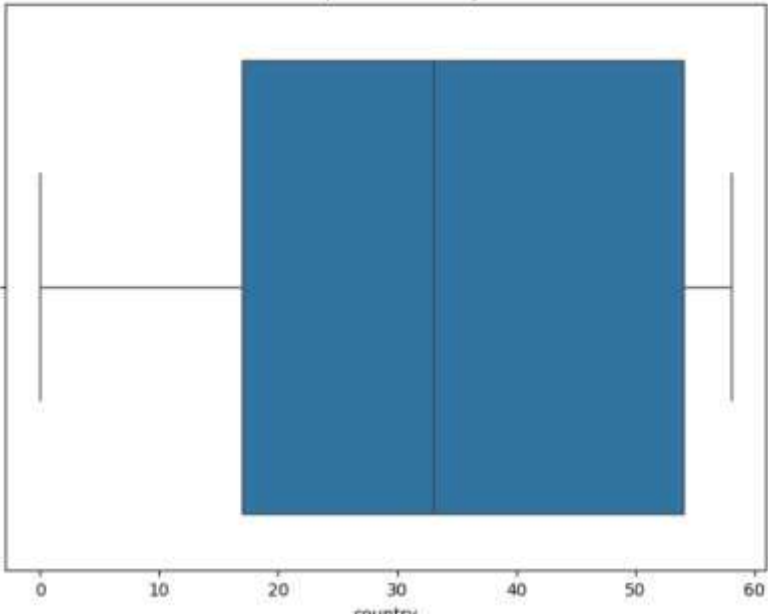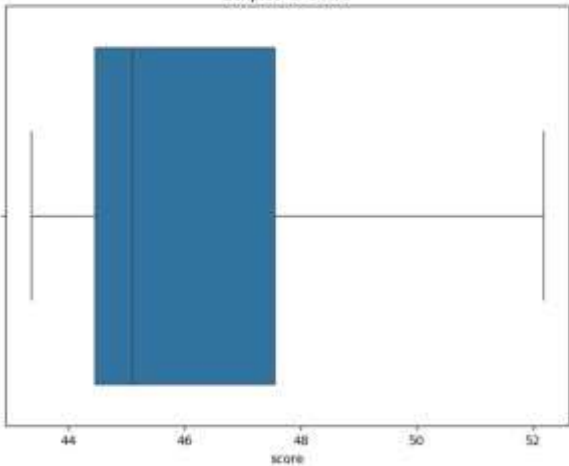| Handling Null Values | ## Handling null Values<br><br>```python<br>np.sum(cwur.isnull())<br>```<br><br>```<br>world_rank          0<br>institution         0<br>country             0<br>national_rank       0<br>quality_of_education 0<br>alumni_employment   0<br>quality_of_faculty  0<br>publications        0<br>influence           0<br>citations           0<br>broad_impact        0<br>patents             0<br>score               0<br>year                0<br>dtype: int64<br>``` |
|---|---|
| Viewing outliers | ### Handling outliers<br><br>```python<br>def fun(col):<br>    sns.boxplot(x=col,data=cwur)<br>    plt.show()<br><br>for i in cwur.columns:<br>    fun(i)<br>```<br><br>country<br><br>✓ 0s   completed at 9:50 PM |

| | |
|---|---|
| |  |
| Handling outliers | ```python
# Iterate over each column and plot boxplot
for column in cwur.columns:
    plt.figure(figsize=(8, 6))  # Adjust the figure size as needed
    sns.boxplot(x=cwur[column])
    plt.title(f'Boxplot for {column}')
    plt.xlabel(column)
    plt.show()
```<br><br>Boxplot for country |

| | |
|---|---|
| | **Boxplot for score**<br><br>![Boxplot for score showing box spanning from about 44.5 to 48 with median near 45.5, whiskers extending to about 43 and 52]<br>44    46    48    50    52<br>score |
| Saved Processed Data | ▶ cwur.shape<br><br>⇥ (2200, 14) |