

# Project Documentation

## Part1: College Distance Dataset Analysis

### Overview

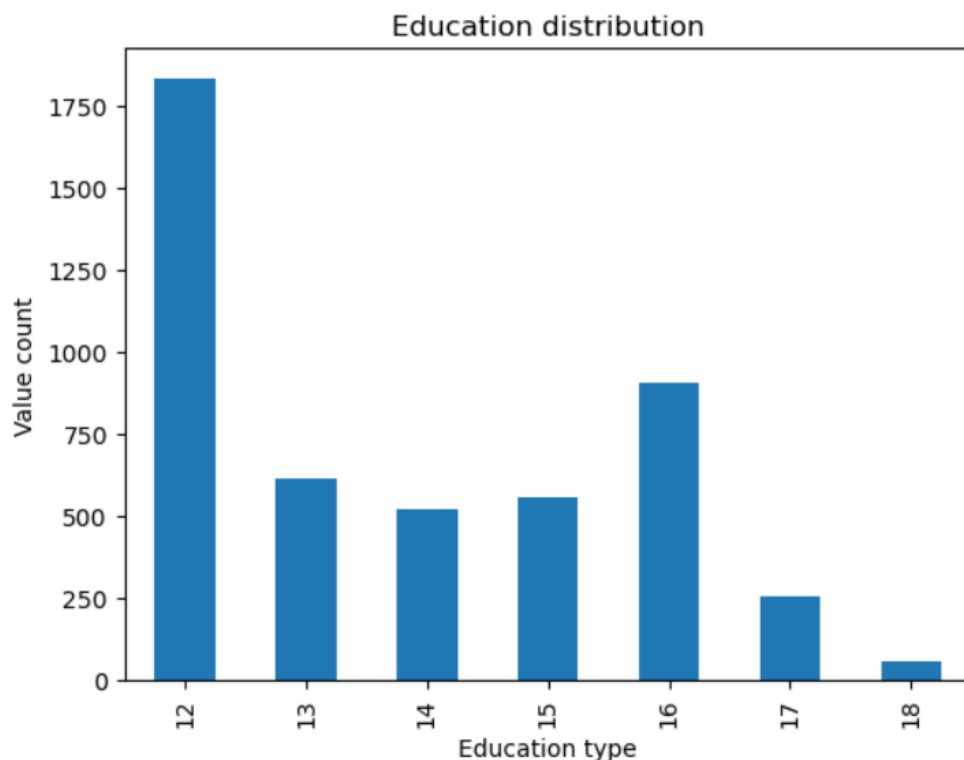
This script provides a data exploration, transformation, and encoding workflow for the College Distance dataset. The primary goals include visualizing data distributions, analyzing key statistics, handling categorical features, and saving a cleaned dataset and the encoder for future use.

### 1. Initial Exploration

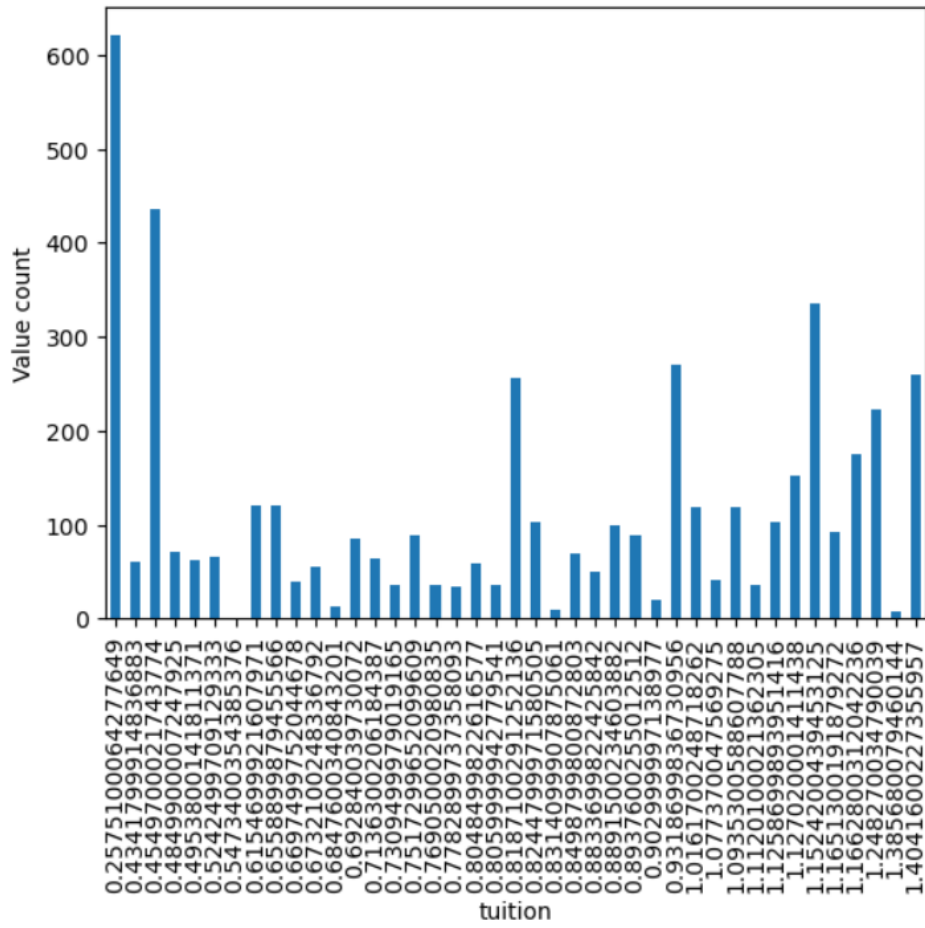
- \* Display the first few rows and statistical summary.
- \* Check for null values and data shape.
- \* Analyze unique values in each column.

### 2. Distribution Visualizations

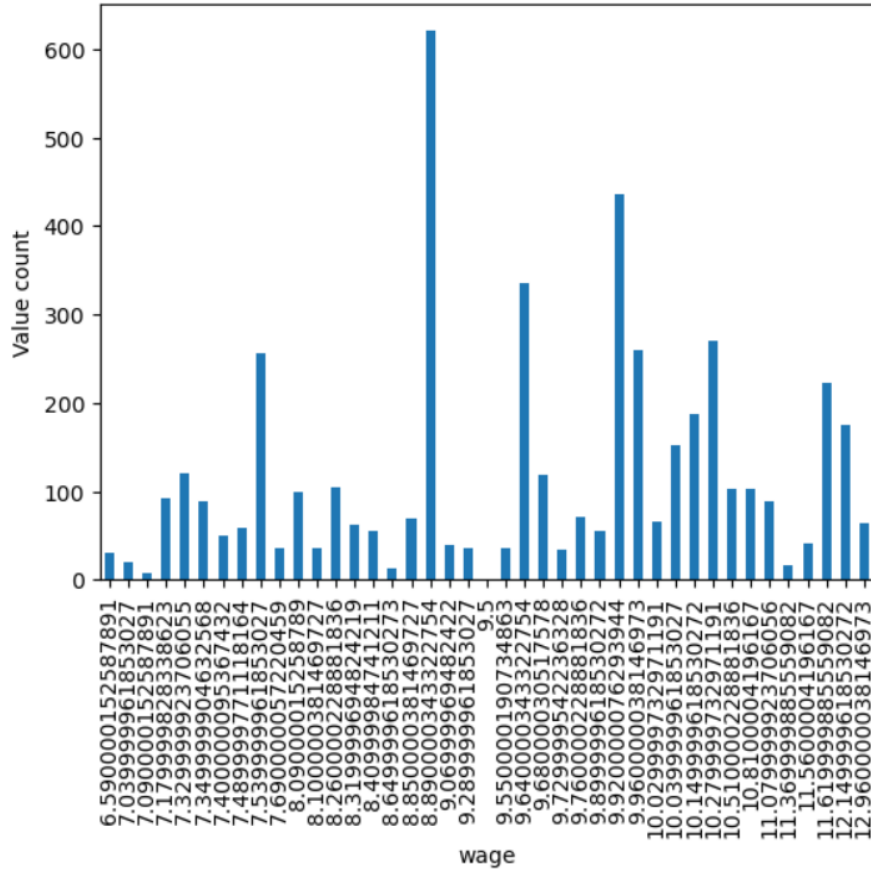
- \* Education Distribution: Bar plot of `education` levels.
- \* Tuition Distribution: Bar plot of `tuition` values.
- \* Tuition Distribution: Bar plot of `tuition` values.
- \* Wage Distribution: Bar plot of `wage` values.
- \* Score Distribution: Histogram of scores binned into ranges.

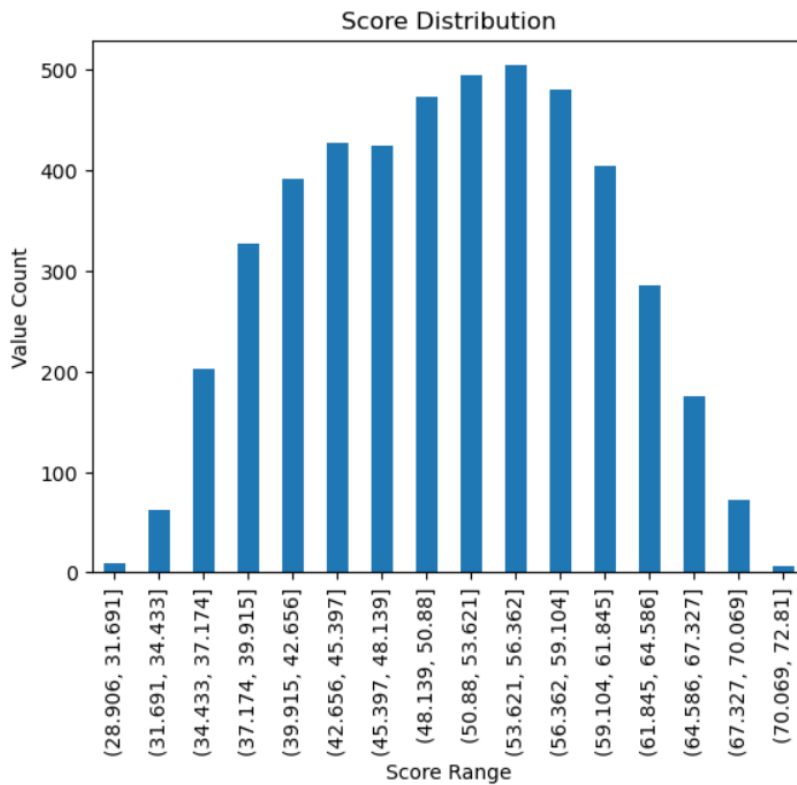


Tuition distribution



Wage distribution



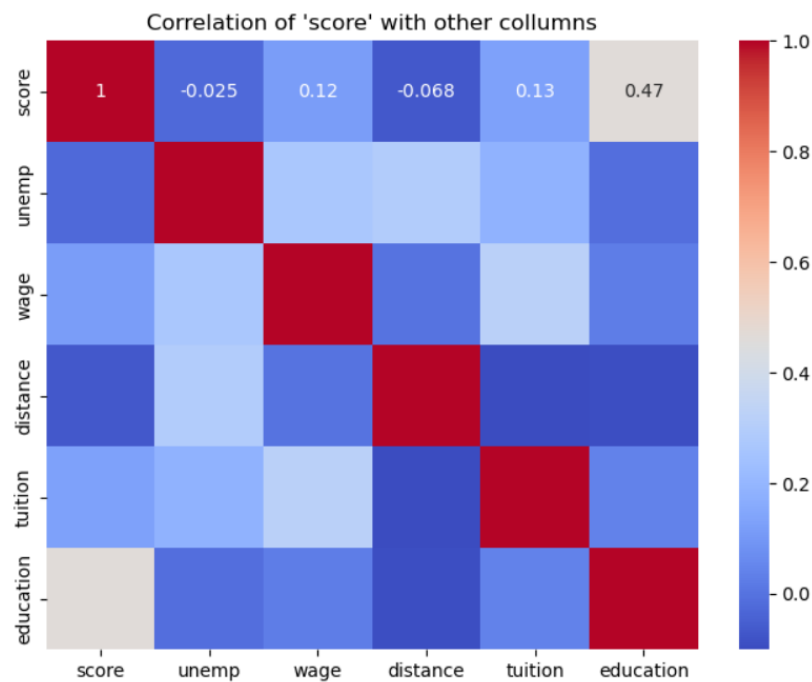


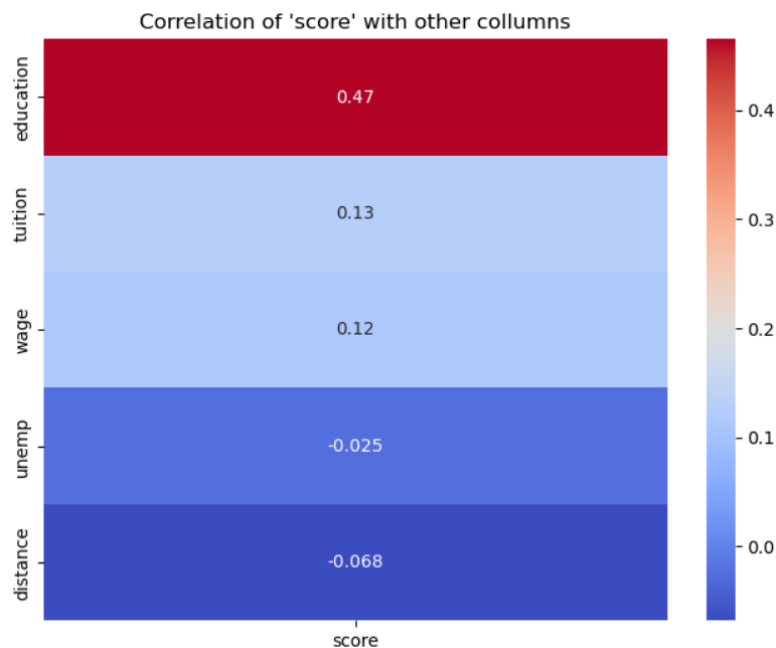
### 3. Statistical Analysis

\* Compute skewness and kurtosis for numeric columns:

```
score: Skewness = -0.03, Kurtosis = -0.88
unemp: Skewness = 1.56, Kurtosis = 5.41
wage: Skewness = 0.09, Kurtosis = -0.26
distance: Skewness = 3.00, Kurtosis = 13.04
tuition: Skewness = -0.15, Kurtosis = -1.05
education: Skewness = 0.44, Kurtosis = -1.23
```

\* Identify correlations between `score` and other columns and display correlation matrices:





#### 4. Categorical Encoding

- \* Identify and encode categorical columns using OrdinalEncoder.
- \* Convert encoded columns to integer type and create a dictionary for encoding mappings:
 

```
{
    'gender': {'female': 0, 'male': 1},
    'ethnicity': {'afam': 0, 'hispanic': 1, 'other': 2},
    'fcollege': {'no': 0, 'yes': 1},
    'mcollege': {'no': 0, 'yes': 1},
    'home': {'no': 0, 'yes': 1},
    'urban': {'no': 0, 'yes': 1},
    'income': {'high': 0, 'low': 1},
    'region': {'other': 0, 'west': 1}}
      
```

#### 5. Saving Results

- \* Save the cleaned dataset as 'CollegeDistanceCleaned.csv'.
- \* Save the encoder as 'ordinal\_encoder.pkl' for reuse.

### Part1: College Model Training and Evaluation

#### Overview

This script builds, trains, and evaluates several regression models to predict the `score` variable in the College Distance dataset. The models explored include Decision Tree, Random Forest, Linear Regression, and Gradient Boosting Regressor, with a focus on optimizing model performance using evaluation metrics and hyperparameter tuning.

#### 1. Data Import and

- \* Features (`X`) are defined by excluding the `score` column, while the target variable (`y`) is set to `score`
- \* Split data into training and testing sets with an 80-20 ratio.

#### 2. Model Training and Evaluation

#### Decision Tree Regressor:

- \* Model 1: Default settings with `squared\_error` criterion.
- \* Model 2 & 3: Modified models using `absolute\_error` criterion with `max\_depth` variations.
- \* Best scores:
  - MAE: 6.1114398879843925
  - MSE: 55.872398998729544
  - r2 score: 0.26978943190142146

#### Random Forest Regressor:

- \* Model 1 & 2: Tested with `n\_estimators` (100 and 150), different criterion values, and depth restrictions.
- \* Model 3: Lower `max\_depth` to 4 for testing overfitting vs. underfitting.
- \* Best scores:
  - MAE: 6.00735368290028
  - MSE: 54.13116539506546
  - r2 score: 0.29254605595389505

#### Linear Regression:

- \* Trained on all features as a baseline linear model.
- \* Best Scores:
  - MAE: 5.999723186927585
  - MSE: 53.029970954561485
  - r2 score: 0.306937845681873

#### Gradient Boosting Regressor:

- \* Initial Models: Tested with various `n\_estimators` and `learning\_rate` values to assess boosting effect.
- \* Modified Feature Set: Dropped columns `ethnicity`, `mcollege`, `fcollege`, `urban`, and `home` and trained the model to assess feature importance.
- \* Custom Function: Defined `within\_10\_percent\_accuracy` to evaluate accuracy within  $\pm 10\%$  tolerance of actual values.
- \* Best scores:
  - MAE: 5.642095694628343
  - MSE: 48.19469718805719
  - r2 score: 0.35878482763328656
  - Accuracy within  $\pm 10\%$ : 51.0548523206751%

### 3. Hyperparameter Tuning with RandomizedSearchCV

- \* Used a parameter grid on Gradient Boosting Regressor (`RandomizedSearchCV`) to identify the best model parameters based on negative MSE score.
- \* Output: Displayed best parameters, MAE, MSE,  $R^2$  score, and accuracy within  $\pm 10\%$  of the target variable.
- \* The search did not beat:
  - MAE: 5.642095694628343
  - MSE: 48.19469718805719
  - r2 score: 0.35878482763328656
  - Accuracy within  $\pm 10\%$ : 51.0548523206751%

#### Key Results

- \* Decision Tree: Performed well with minimal tuning but sensitive to depth settings.
- \* Random Forest: Higher accuracy with more trees, though computationally intensive with larger depths.

- \* Linear Regression: Baseline model provided reference results with limited feature interactions.
- \* Gradient Boosting: Effective in improving accuracy, especially after tuning and dropping less influential features.

Best Overall scores:

MAE: 5.642095694628343

MSE: 48.19469718805719

r2 score: 0.35878482763328656

Accuracy within +-10%: 51.0548523206751%