

### **Regression Dataset EDA and Preprocessing:**

The student math score dataset, incorporating gender and ethnic group data, underwent exploratory data analysis (EDA) and preprocessing. Key steps are summarized below:

- **Feature Selection:** Gender and ethnic\_group columns were excluded as they weren't directly relevant to math score prediction. Their omission streamlined the dataset.
- **Numerical Conversion:** Remaining string variables were transformed into numerical equivalents using the LabelEncoder class. This conversion was vital for LinearRegression to process the data effectively.
- **Data Splitting:** The dataset was divided into training and test sets to ensure unbiased model evaluation. Training on one set and evaluating on another mimics real-world predictions.
- **Model Selection:** The LinearRegression model was chosen for its simplicity and interpretability. Its suitability was reinforced by achieving a satisfactory mean squared error (MSE) score of 10.25 on the test set.

The rationale behind these steps:

- **Feature Omission:** Gender and ethnic\_group columns had no direct bearing on math score prediction. Although they might indirectly influence scores, they lacked direct predictive power.
- **Numerical Conversion:** LinearRegression mandates numerical data. The LabelEncoder maintained data semantics while enabling numerical processing.
- **Data Splitting:** Segregating training and test sets ensures fair model evaluation on new data. It simulates real-world predictive scenarios.
- **Model Choice:** LinearRegression's simplicity and interpretability make it suitable for this context. Its satisfactory test set performance supports its suitability.

### **Correlation Coefficient Insight:**

The correlation coefficient gauges variable relationships, ranging from -1 (perfect negative) to 1 (perfect positive). The correlation matrix revealed key correlations with MathScore, ReadingScore, and WritingScore:

- **ParentEduc:** Parental education level correlated with math, reading, and writing scores. Educated parents tend to provide a better educational environment.
- **LunchType:** Meal nature correlated with scores. Subsidized meals related to lower scores due to potential resource limitations in low-income families.
- **TestPrep:** Participation in test preparation linked to higher scores. Such courses enhance familiarity with test content and strategies.

Although weaker associations with scores were observed for other variables, potential evidence suggests their role in student achievement. For instance, sports engagement improved math and reading scores through skills like teamwork.

**Summary:**

The T-statistic and P-value analyses highlighted statistically significant differences in MathScores, WritingScores, and ReadingScores between males and females. A low MSE indicated the model's accurate prediction of scores. Overall, the model effectively captured gender-based score differences with statistical significance.

This comprehensive approach to EDA, preprocessing, and model evaluation equips us to gain insights from the dataset and make informed decisions based on meaningful statistical analysis.