# Analysis and Prediction of GCSE English and Maths National Data 2019-20

**Executive Summary**

Maths results. The data-driven approach aimed to discern patterns and factors influencing student performance. The linear regression model developed provides a foundation for predicting educational outcomes and formulating interventions to enhance student achievement.

**Introduction**

The GCSE results are a pivotal indicator of educational attainment in the UK. Understanding the determinants of these outcomes can inform policy-making and teaching strategies. This project leverages data science techniques to analyze the national dataset for the academic year 2019-20, with the objective of identifying key predictors of student performance in English and Maths.

## Data Acquisition and Preprocessing

**Data Source**

The dataset was obtained in CSV format, containing various student attributes alongside English and Maths GCSE scores.

The data has been take from [www.data.gov.uk](www.data.gov.uk), which is a opensource database.

**Data Cleaning**

Data cleaning was crucial to ensure reliability. Numeric conversions and the removal of incomplete records resulted in a robust dataset for analysis.

## Exploratory Data Analysis (EDA)

**Descriptive Statistics**

Initial exploration provided insights into score distributions and central tendencies, highlighting the need for a closer examination of underlying factors.

### Visual Analysis

Histograms and bar charts revealed disparities in performance across different ethnic groups and other demographic factors.

### Correlation Analysis

A correlation matrix identified relationships between numerical variables, providing an empirical basis for feature selection.

## Feature Engineering and Selection

### Feature Identification

Features like ethnicity, gender, and educational support indicators were selected based on their potential impact on student performance.

### Data Preparation

Categorical variables were encoded, and the dataset was split into training (80%) and testing (20%) subsets to prepare for machine learning application.

Model Development

### Linear Regression

A linear regression model was chosen for its interpretability and relevance to continuous outcome prediction.

### Training

The model was trained using the processed training set, enabling it to learn the relationship between features and outcomes.

## Model Evaluation and Validation

### Performance Metrics

The model's performance was quantified using the Mean Squared Error (MSE) and R-squared metrics, yielding an MSE of 53.4 and an R-squared of 0.2 on the test set.

**Interpretation**

While the model could predict GCSE outcomes to a certain extent, the results suggest room for improvement. The moderate R-squared value indicates that additional variables and perhaps more complex modeling techniques may yield better prediction accuracy.

**Conclusion**

This analysis has highlighted the multifaceted nature of educational outcomes. The linear regression model serves as a baseline for future studies. However, the moderate predictive power suggests the need for a more nuanced approach, potentially incorporating additional data sources and employing more sophisticated models such as ensemble methods or neural networks.

# Recommendations

### Model Refinement

Further refinement of the predictive model is recommended, potentially through the inclusion of additional features, more advanced modeling techniques, and cross-validation.

### Educational Interventions

Insights from the analysis could be used to target educational interventions, such as additional support for groups identified as underperforming.

### Continued Research

Ongoing research is crucial to keep up with the evolving educational landscape, especially considering the potential impacts of recent global events on education.