



# STUDENT PERFORMANCE FACTORS ANALYSIS

Dana Ortiz

01

## Background

The industry, company involved and the stakeholders role.

03

## Insights and Explanations

What did we do with the data and how it can be translated into tangible business value?

02

## The Dataset

The bare bones of the data file we are working with.

04

## Suggestions

What can your company do based on our professional insights



01

**BACKGROUND**

## Background - Company

**Pupil Progress Analytics** is a forward-thinking educational technology company focused on improving student outcomes through data-driven insights and innovative educational solutions. The company partners with schools, educational institutions, and tutoring organizations to integrate cutting-edge data analytics into their operations, helping them understand and address various academic challenges.



# Background - Stakeholder

As the Director of Data Analytics & Research, the Stakeholder (You!) leads a team that focuses on gathering and analyzing data to provide actionable insights that can improve student learning experiences. This role involves using data to help institutions better understand their students' needs, track academic performance trends, and suggest personalized interventions that can help improve student success rates and thus the school's competitiveness.



02

THE DATASET



**Target Variable:**

Final Exam Score (See all variables on next slide)



Our "Student Performance Factors" dataset provides comprehensive insights into various elements that affect student academic performance. It includes data on study habits, parental involvement, attendance, sleep patterns, access to resources, and more, which all contribute to the final exam score of students.



**Number of Records:**

6,607

**Number of Features:**

20





## Variables

- **Hours\_Studied:** Number of hours spent studying per week.
- **Attendance:** Percentage of classes attended.
- **Parental\_Involvement:** Level of parental involvement.
- **Access\_to\_Resources:** Availability of resources.
- **Extracurricular\_Activities:** Participation in activities.
- **Sleep\_Hours:** Average number of hours of sleep per night.
- **Previous\_Scores:** Scores from previous exams.
- **Motivation\_Level:** Student's level of motivation.
- **Internet\_Access:** Availability of internet access.
- **Tutoring\_Sessions:** Number of sessions attended monthly.
- **Family\_Income:** Family income level .
- **Teacher\_Quality:** Quality of the teachers .
- **School\_Type:** Type of school attended (Public, Private).
- **Peer\_Influence:** Influence of peers on performance.
- **Physical\_Activity:** Average number of hours of physical activity per week.
- **Learning\_Disabilities:** Presence of learning disabilities.
- **Parental\_Education\_Level:** Highest education level of parents (High School, College, Postgraduate).
- **Distance\_from\_Home:** Distance from home to school.
- **Gender:** Gender of the student (Male, Female).
- **Exam\_Score:** Final exam score.



03

# INSIGHTS AND EXPLANATIONS



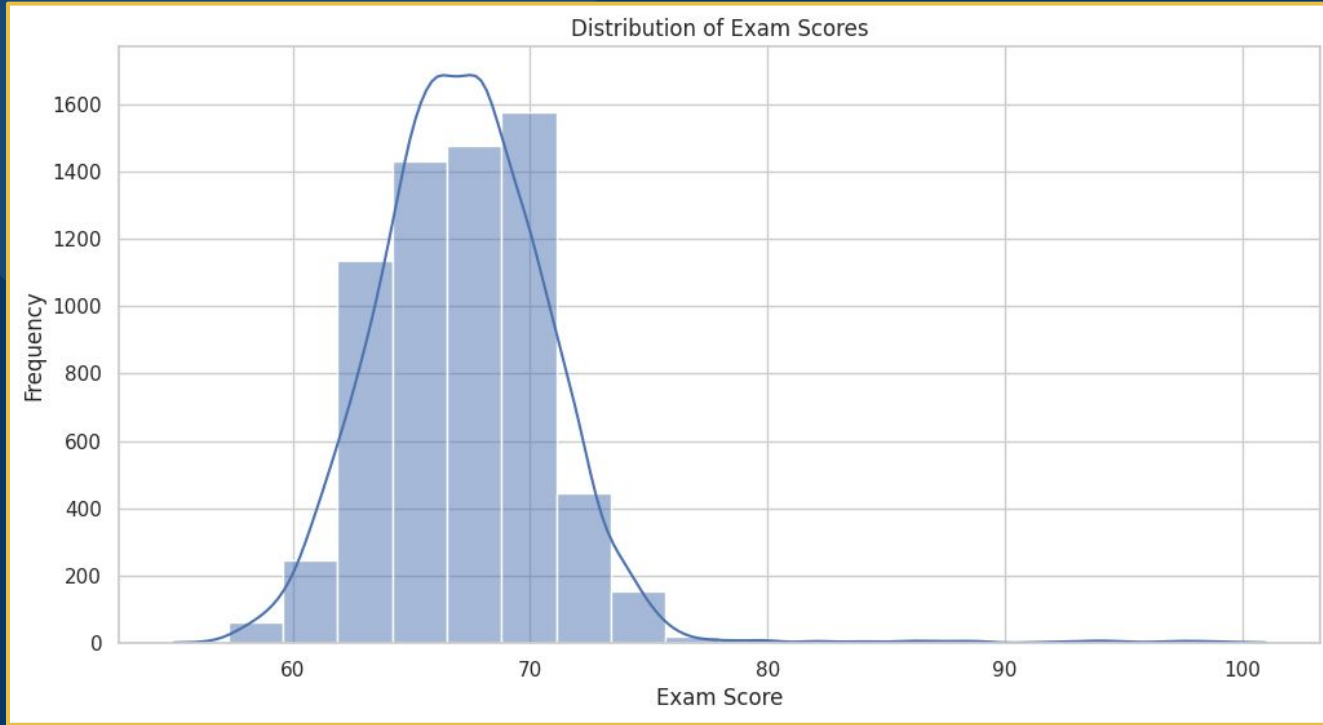
## How did we prep the data?

To prepare the data for analysis, all categorical variables were cleaned and converted into numerical form.

Ordinal variables such as "Parental Involvement" were mapped to numerical values from 1 to 3 based on levels of "Low," "Medium," and "High."

Binary variables like "Internet Access" were encoded as 1 for "Yes" and 0 for "No."





This histogram shows the overall distribution of student exam scores. The shape of the curve helps us understand how scores are spread out across the population. We observe that the majority of students fall within a middle range, with fewer students at the very low or very high ends. This normal-like distribution suggests a relatively balanced dataset with a few outliers, and it gives us a baseline for interpreting how various factors might shift a student's score.

# Key Insights

## Strong Positive Influences

**Access\_to\_Resources (1.02):** Students with better access to educational resources tend to score higher.

**Parental\_Involvement (1.00):** Active parental engagement is associated with improved student performance.

**Internet\_Access (0.90):** Reliable internet access facilitates learning and positively impacts scores.

**Motivation\_Level (0.55):** Highly motivated students achieve better exam results.

**Teacher\_Quality (0.55):** Quality teaching contributes significantly to student success.



## Negative Influences

**Learning\_Disabilities (-0.86):** Students with learning disabilities may require additional support to achieve comparable scores.

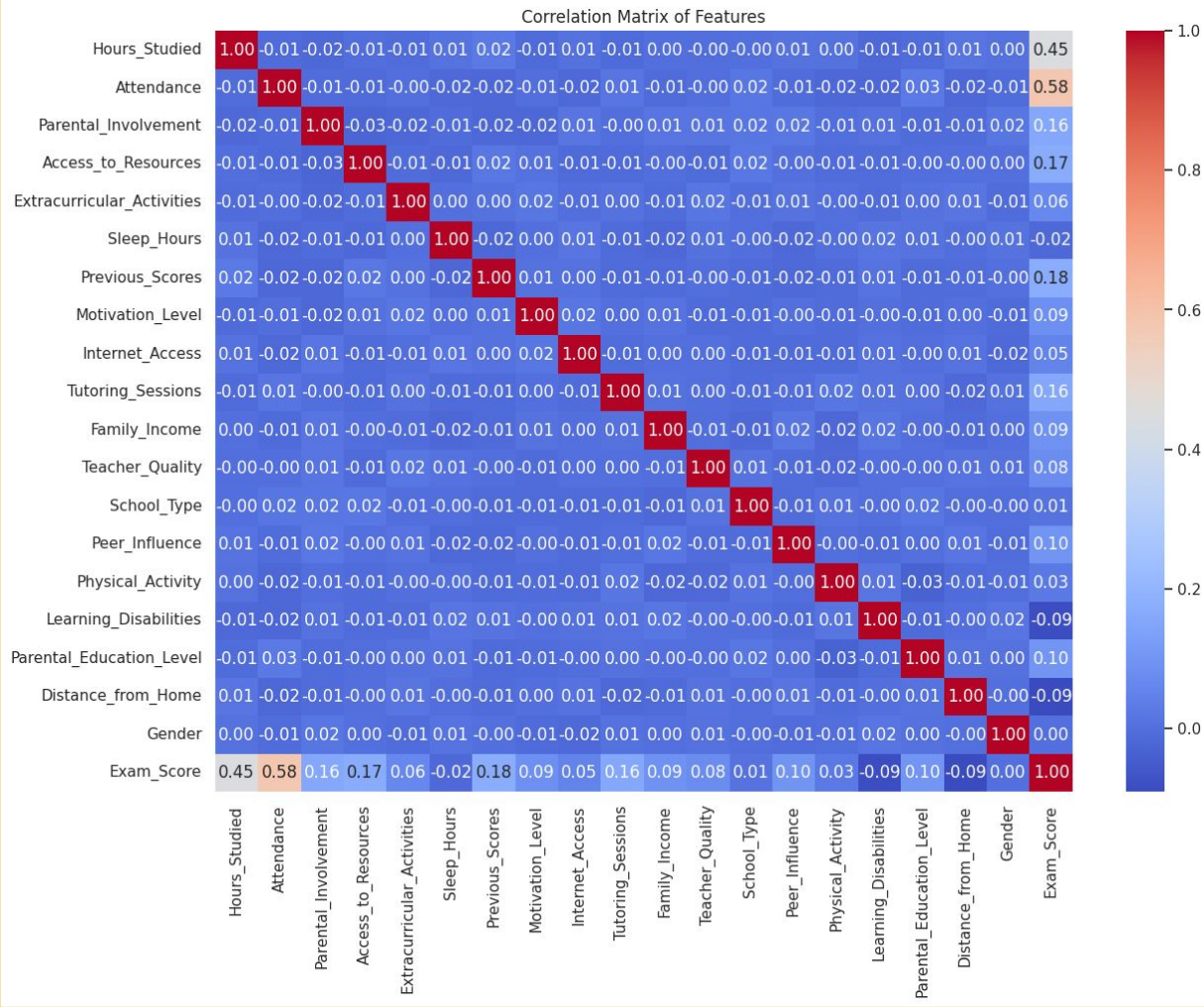
**Distance\_from\_Home (-0.46):** Longer commutes can negatively affect student performance, possibly due to fatigue or reduced study time.

**School\_Type (-0.05):** A slight negative coefficient suggests minimal impact, but further investigation may be warranted.



One of the most valuable visual tools used in this analysis was the correlation matrix. This heatmap allowed us to see which numerical variables were most strongly associated with final exam scores.

Notably, we observed strong positive correlations between Exam Score and variables such as Hours Studied and attendance. There were also many lower positive correlations with Parental Involvement and tutoring. Learning Disabilities and Distance from Home showed slightly negative relationships with exam performance.







## Linear Regression

Linear Regression was applied to predict student exam scores based on a combination of behavioral, environmental, and demographic factors. This model showed that certain factors had stronger positive coefficients, confirming their importance in academic outcomes. Although the model assumes linear relationships and doesn't capture complex interactions as well as a decision tree, it provides clear and actionable weightings for each variable, helping prioritize interventions and explain more variance for this specific dataset.



# Linear Regression Equation:

$$\begin{aligned} \text{Exam\_Score} = & (0.29 * \text{Hours\_Studied}) + (0.20 * \text{Attendance}) + (1.00 * \\ & \text{Parental\_Involvement}) + (1.02 * \text{Access\_to\_Resources}) + (0.56 * \\ & \text{Extracurricular\_Activities}) + (-0.01 * \text{Sleep\_Hours}) + (0.05 * \\ & \text{Previous\_Scores}) + (0.55 * \text{Motivation\_Level}) + (0.90 * \text{Internet\_Access}) \\ & + (0.48 * \text{Tutoring\_Sessions}) + (0.56 * \text{Family\_Income}) + (0.55 * \\ & \text{Teacher\_Quality}) + (-0.05 * \text{School\_Type}) + (0.50 * \text{Peer\_Influence}) + \\ & (0.20 * \text{Physical\_Activity}) + (-0.86 * \text{Learning\_Disabilities}) + (0.49 * \\ & \text{Parental\_Education\_Level}) + (-0.46 * \text{Distance\_from\_Home}) + (0.01 * \\ & \text{Gender}) + (0.34 * \text{Passed}) + 30.49 \end{aligned}$$



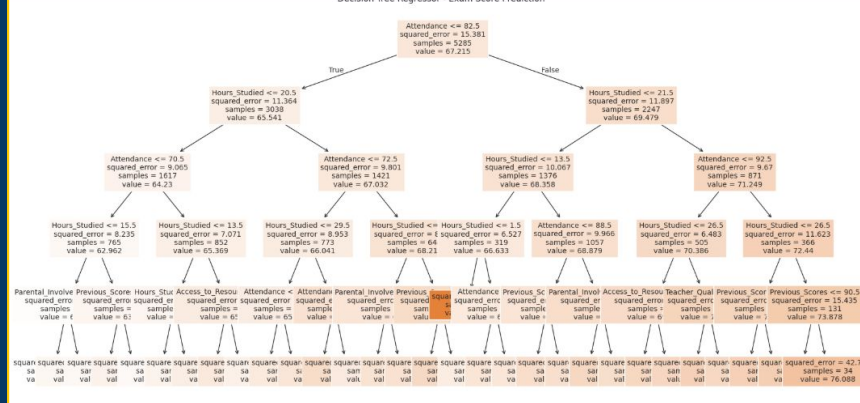
## Decision Tree

We also tried a Decision Tree Regressor to model the relationship between student factors and their final exam scores. This model was almost chosen because it captures complex, non-linear interactions between variables and provides a highly interpretable structure, perfect for identifying actionable drivers of student success.

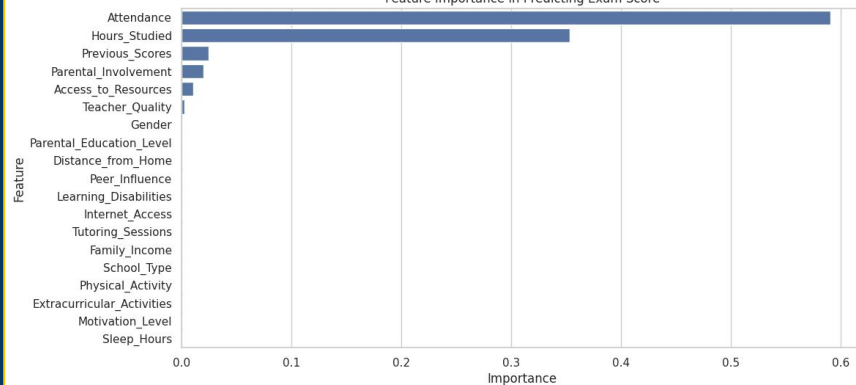
However, our decision tree has a Mean Squared Error of 6.44 and a  $R^2$  Score of 0.54 and our linear regression had a MSE of 4.15 and  $R^2$  of 0.73. This means that the linear regression is a better fit for our dataset and our goals.

This model reveals which factors most strongly affect performance and can guide interventions. This insight can help the company prioritize resources toward improving these specific areas to boost student success.

Decision Tree Regressor - Exam Score Prediction



Feature Importance in Predicting Exam Score



04

## SUGGESTIONS



# Suggestions at a Glance

Based on our model's findings, the following recommendations can be made:

- **Enhance Resource Availability:** Invest in educational materials and ensure students have access to necessary learning tools.
- **Promote Parental Engagement:** Implement programs to encourage parents to participate in their children's education.
- **Improve Internet Infrastructure:** Provide reliable internet access to support digital learning initiatives.
- **Support Motivational Programs:** Develop strategies to boost student motivation, such as goal-setting workshops or mentorship programs.
- **Invest in Teacher Development:** Offer professional development opportunities to enhance teaching quality.
- **Address Learning Disabilities:** Provide specialized support and resources for students with learning challenges.
- **Consider Transportation Solutions:** Explore options to reduce commute times, such as school transportation services or remote learning alternatives.

# Providing Access to Open Educational Resources

## Implementation Steps:

- **Curate Quality OER Materials:** Identify and select high-quality, curriculum-aligned open resources.
- **Train Educators:** Offer professional development to help teachers effectively integrate OER into their instruction.
- **Develop an OER Repository:** Create a centralized digital library where students and teachers can easily access materials.

Studies have shown that implementing OER can save students significant amounts on textbook costs. For example, students at schools participating in OER initiatives paid at least \$65 less per course on average.

## Anticipated Benefits:

- **Cost Savings for Students:** Reduces the financial burden of purchasing textbooks and materials.
- **Increased Accessibility:** Ensures all students have equal access to required learning materials.
- **Enhanced Learning Outcomes:** OER has been shown to increase student learning while breaking down barriers of affordability and accessibility.

## Influence on Other Factors:

- **Motivation Level:** Easier access to materials can increase student motivation and preparedness.
- **Parental Involvement:** Parents can utilize OER to better support their children's learning at home.

All of our code was written and  
ran in Google Collab. Compiled  
here is an ipynb for the code at the  
date of this presentation and the  
full dataset csv.

[Download ipynb](#)

[Download csv](#)



A low-angle shot of four graduates in black gowns and mortarboards, holding rolled-up diplomas tied with red ribbons. They are looking up and throwing their mortarboards into the air. In the background is a large, ornate university building with many windows and a clock face. The sky is overcast.

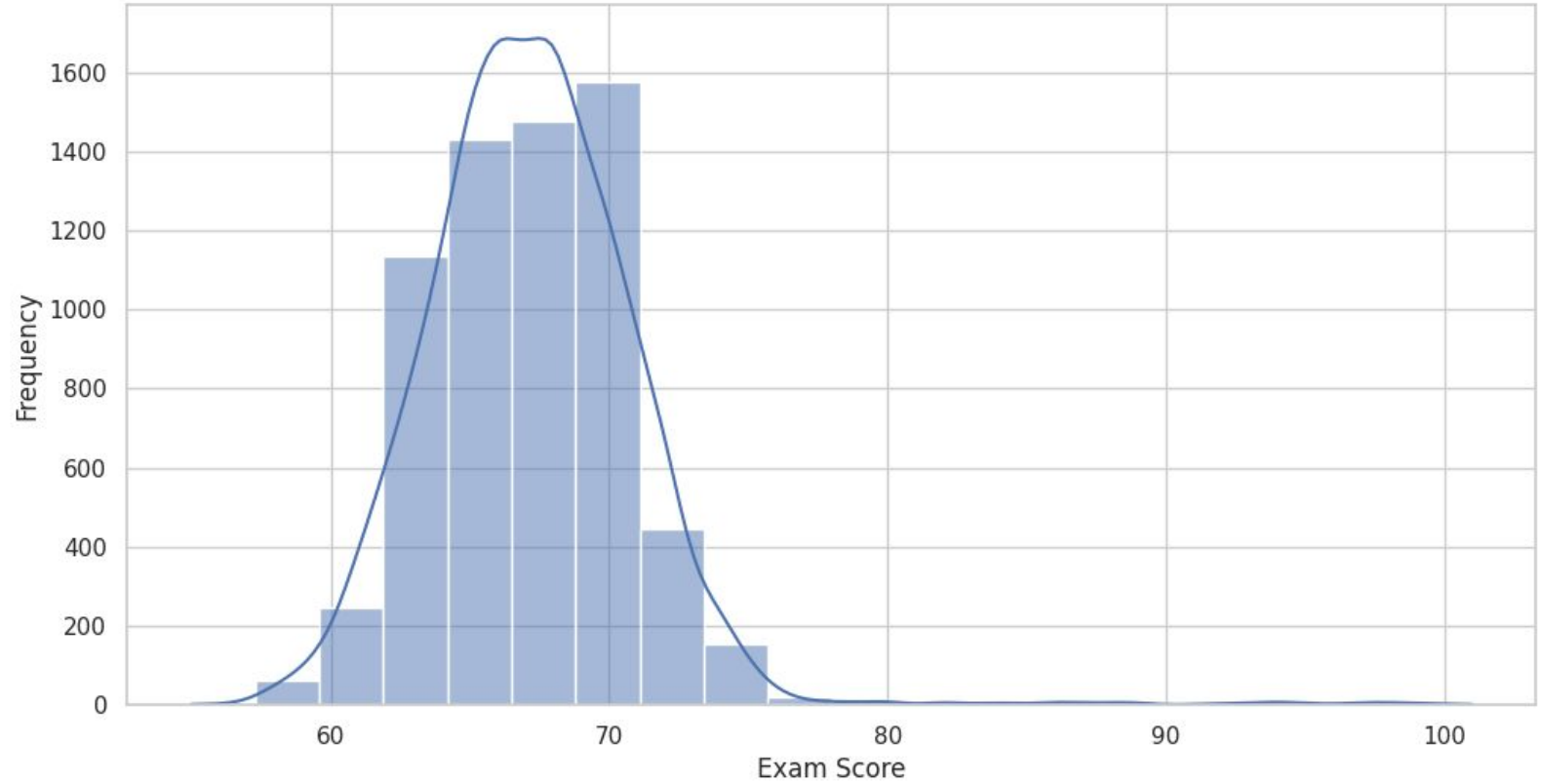
Thank you for your time!

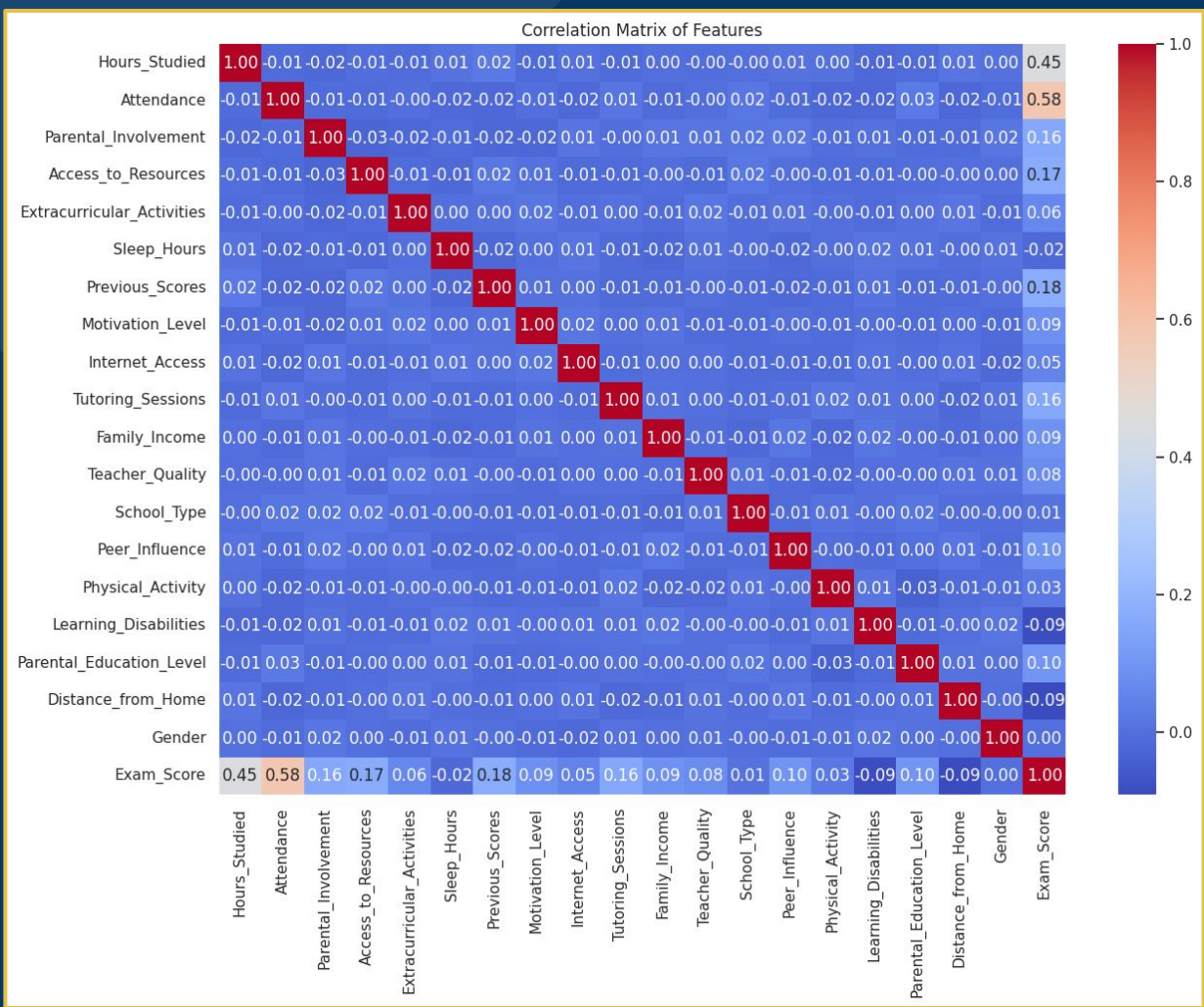
**A**

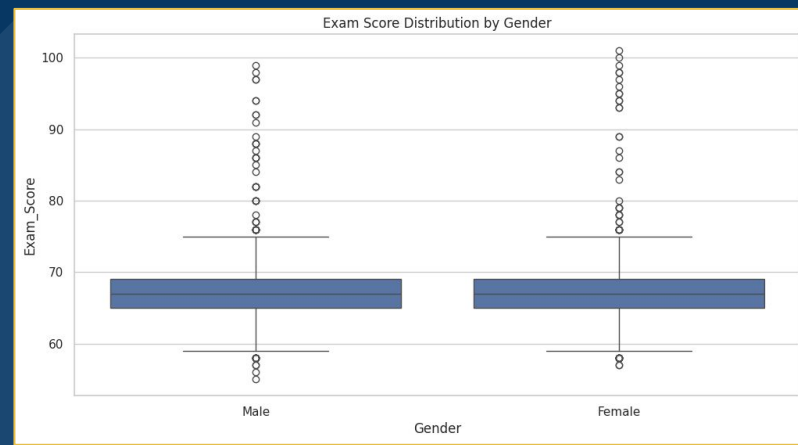
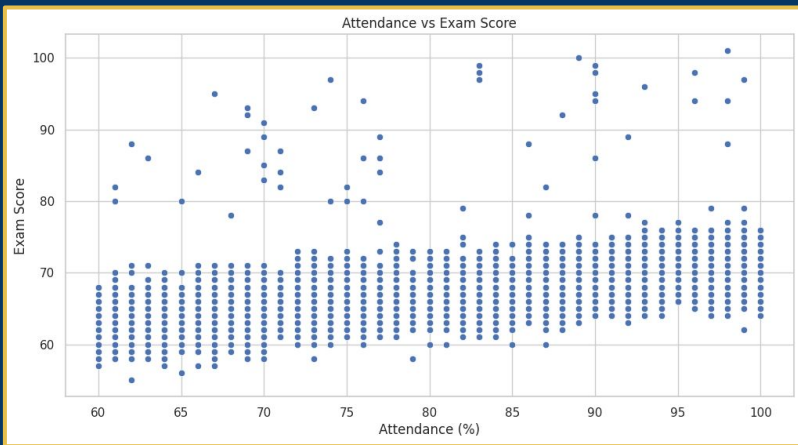
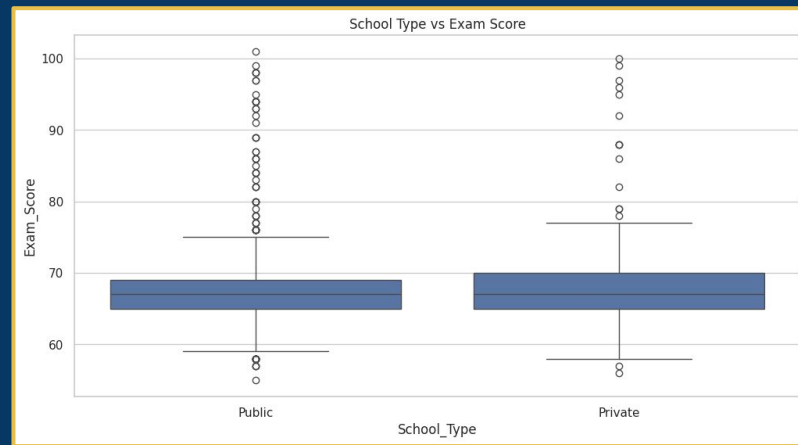
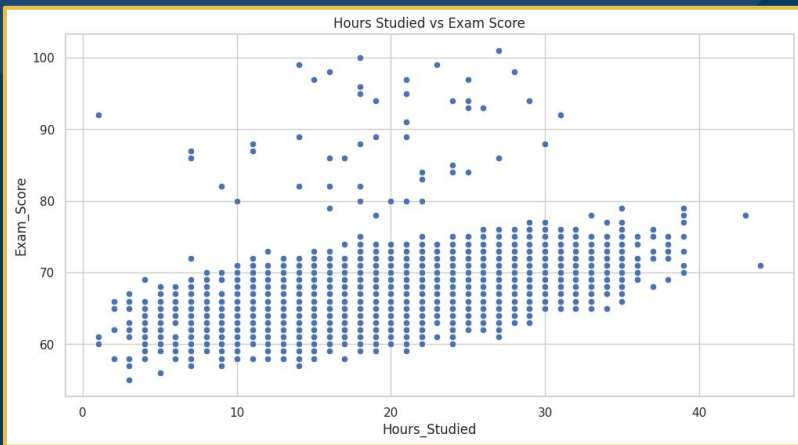
# APPENDIX



Distribution of Exam Scores







# Linear Regression Equation:

$$\begin{aligned} \text{Exam\_Score} = & (0.29 * \text{Hours\_Studied}) + (0.20 * \text{Attendance}) + (1.00 * \\ & \text{Parental\_Involvement}) + (1.02 * \text{Access\_to\_Resources}) + (0.56 * \\ & \text{Extracurricular\_Activities}) + (-0.01 * \text{Sleep\_Hours}) + (0.05 * \\ & \text{Previous\_Scores}) + (0.55 * \text{Motivation\_Level}) + (0.90 * \text{Internet\_Access}) \\ & + (0.48 * \text{Tutoring\_Sessions}) + (0.56 * \text{Family\_Income}) + (0.55 * \\ & \text{Teacher\_Quality}) + (-0.05 * \text{School\_Type}) + (0.50 * \text{Peer\_Influence}) + \\ & (0.20 * \text{Physical\_Activity}) + (-0.86 * \text{Learning\_Disabilities}) + (0.49 * \\ & \text{Parental\_Education\_Level}) + (-0.46 * \text{Distance\_from\_Home}) + (0.01 * \\ & \text{Gender}) + (0.34 * \text{Passed}) + 30.49 \end{aligned}$$



## Decision Tree Regressor - Exam Score Prediction

