# Student Exam Scores Data Insights with PySpark

Sai Puneeth Pogula      Priyanka Kurimilla
Sai Teja Mallineni      Navya Ravipati

December 2023

## 1  Project Description

This dataset includes scores from three test scores of students at a (fictional) public school and a variety of personal and socio-economic factors that may have interaction effects upon them.

This dataset is fictional and should be used for educational purposes only.

There are a few attributes in the datasets:

1. **Gender:** Gender of the student (male/female).

2. **EthnicGroup:** Ethnic group of the student (group A to E).

3. **ParentEduc:** Parent(s) education background (from some high school to master's degree).

4. **LunchType:** School lunch type (standard or free/reduced).

5. **TestPrep:** Test preparation course followed (completed or none).

6. **ParentMaritalStatus:** Parent(s) marital status (married/single/widowed/divorced).

7. **PracticeSport:** How often the student practices sport (never/sometimes/regularly)).

8. **IsFirstChild:** If the child is the first child in the family or not (yes/no).

9. **NrSiblings:** Number of siblings the student has (0 to 7).

10. **TransportMeans:** Means of transport to school (school bus/private).

11. **WklyStudyHours:** Weekly self-study hours (less than 5 hrs; between 5 and 10 hrs; more than 10 hrs).

12. **MathScore:** Math test score (0-100).

13. **ReadingScore:** Reading test score (0-100).

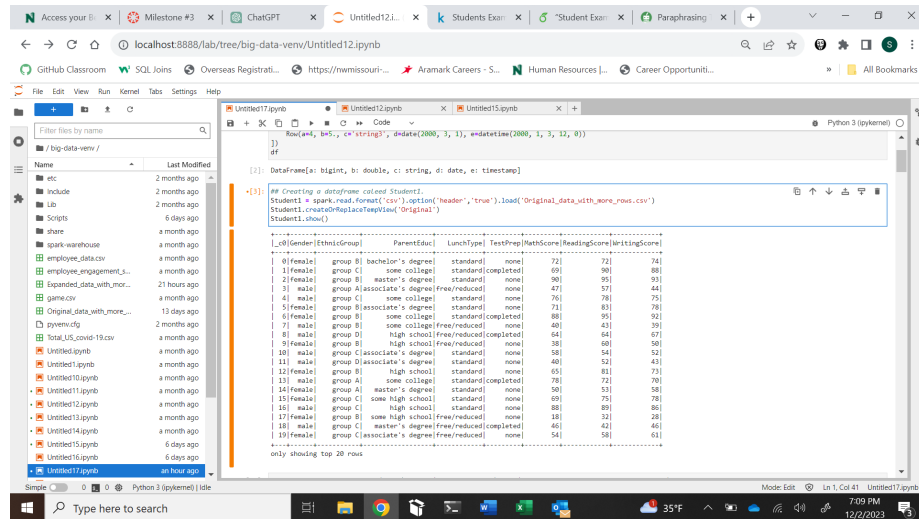14. **WritingScore:** Writing test score (0-100).

Here, all those attributes that I am using to find out data offer insights into the academic performance of the students and the different elements that may affect it.

I am using PySpark SQL query technologies to find out each student's performance in different sectors like TestPrep and Ethnic Group.

# 2 Results Summary

There are results of the goals:

**Goal 1:** Creating the dataframe

**Goal 2:**Determine the count of students based on their ethnic group and whether they underwent test preparation.



**Goal 3:**Assess the overall academic performance of students by examining average scores in Math, Reading, and Writing.

**Goal 4:** Determine the count of students based on their ethnic group and whether they underwent test preparation.



**Goal 5 :** Observing how different ethnic groups perform different acedemic subjects.

**Goal 6 :**Observing the impact of parental education on students' scores.



**Goal 7 :**Explore the impact of family-related factors like maritial status or number of siblings on academic performance.

**Goal 8:**Observing how the mode of transportation to school relates to weekly study hours drawing a graphs using matplotlib.

**Goal 9 :** Determine if students who completed test preparation performed better.

**Goal 10 :**Observing the relationship between weekly study hours and academic performance and also drawing a graphs between weekly study hours and academic performance using matplotlib.

# 3 Conclusion

We can conclude for this project that we can calculate the average scores in Math, Reading, and Writing to get an overall understanding of the students' performance where reading score average is highest while compare to other mathscore and writingscore and also Investigate whether the level of parental education has an impact on student scores and in various elements.

# 4 Citations

Provide all the necessary citations for the sources you utilized to finish this project successfully. Save all your work to your GitHub repo and provide the URL.

**GitHub Repo:** https://github.com/saipuneet/StudentExamScore.git