# Exploring the Relationship Between Student Characteristics and Academic Performance:

<span style="color:red">A Canonical Correlation Analysis</span>

Reboly Seevaratnam

S/18/825

# **Table of Contents**

# Introduction

In the realm of educational research, understanding the multifaceted factors that contribute to students' academic performance is crucial. Academic achievement not only reflects the individual capabilities of students but is also influenced by a variety of socio-economic, demographic, and educational factors. These insights are essential for educators, policymakers, and stakeholders who strive to improve educational outcomes and foster equitable learning environments.

The "Students Performance in Exams" dataset provides a comprehensive collection of data points that encompass various student demographics, socio-economic indicators, and their respective performances in standardized exams. By utilizing advanced statistical techniques such as Canonical Correlation Analysis (CCA), this study aims to uncover the underlying relationships between these diverse factors, offering a holistic view of what drives academic success.

### Purpose of the Study

The primary purpose of this study is to investigate the complex interplay between student characteristics and their academic performance. Specifically, we aim to explore how demographic variables such as gender, parental education level, lunch status, and test preparation course completion correlate with academic outcomes in mathematics, reading, and writing.

### Research Objectives

The study is guided by the following objectives:

- To identify significant correlations between socio-economic indicators (e.g., parental education level, lunch status) and academic performance metrics (e.g., math, reading, writing scores).

- To determine the extent to which gender differences impact academic achievement in various subjects.

- To evaluate the influence of test preparation resources on students' performance in standardized exams.

Understanding the determinants of academic performance is critical for developing effective educational policies and interventions. This study's findings can inform targeted strategies to address achievement gaps, promote equitable access to educational resources, and enhance overall student performance. By leveraging Canonical Correlation Analysis, we can provide nuanced insights into the complex relationships between student characteristics and academic outcomes, contributing valuable knowledge to the field of educational research.

# Methodology

Dataset Description

The "Students Performance in Exams" dataset, obtained from Kaggle, comprises 1000 observations across 8 variables. The dataset includes demographic information about students, such as gender, parental education level, and lunch status, as well as their scores in three academic subjects: math, reading, and writing. The key variables in the dataset are as follows:

1.gender: The gender of the student (male or female).

2.parental.level.of.education: The highest level of education attained by the student's parents.

3.lunch: The type of lunch the student receives (standard or free/reduced).

4.test.preparation.course: Whether the student completed a test preparation course.

5.math.score, reading.score, writing.score: The scores obtained by students in math, reading, and writing, respectively.

Data Preprocessing

To prepare the dataset for analysis, the following preprocessing steps were undertaken:

Sampling: A random sample of 100 observations was selected from the original 1000 observations to manage computational complexity and focus the analysis.

Categorical Variables Conversion: Categorical variables (gender, parental.level.of.education, lunch, and test.preparation.course) were converted into factors for proper encoding.

One-hot Encoding: Categorical variables were one-hot encoded to convert them into numerical format suitable for analysis. This process involved creating binary columns for each category within the categorical variables.

Standardization: The numerical outcome variables (math.score, reading.score, writing.score) were standardized to have a mean of 0 and a standard deviation of 1 to ensure comparability.

Dimensionality Reduction

Given the high dimensionality of the one-hot encoded predictors, Principal Component Analysis (PCA) was employed to reduce the number of predictor variables while retaining most of the variability in the data. PCA helped project the original high-dimensional data onto a lower-dimensional subspace, making the canonical correlation analysis more tractable.

PCA Execution: PCA was performed on the one-hot encoded predictor variables, and the cumulative proportion of variance explained by the principal components was analyzed.

Component Selection: The number of principal components to retain was determined based on achieving at least 95% of the cumulative variance explained.


<u>Canonical Correlation Analysis (CCA)</u>

Canonical Correlation Analysis was then performed to explore the relationships between the reduced set of predictor variables (obtained from PCA) and the standardized outcome variables. CCA is a multivariate statistical technique that identifies and measures the associations between two sets of variables.

CCA Execution:

 CCA was applied to the PCA-transformed predictors and standardized outcomes to derive canonical correlations.

Canonical Loadings: Canonical loadings were calculated to interpret the relationships between the original variables and the canonical variables.

# Results And Discussion

Canonical Correlation Analysis Results

After performing Canonical Correlation Analysis (CCA) on the reduced set of predictor variables (obtained through Principal Component Analysis) and the standardized outcome variables, we obtained the following canonical correlations:

cca_result$cor

The canonical correlations represent the strength of the relationship between the canonical variables derived from the predictor and outcome sets. For our analysis, the canonical correlations are as follows (example values are given; replace these with actual values from your analysis):

Canonical Correlation 1: 0.75

Canonical Correlation 2: 0.58

Canonical Correlation 3: 0.42

These values indicate that there are moderately strong relationships between the first pair of canonical variables, with the strength of the relationship diminishing for the subsequent pairs.

Canonical Loadings

Canonical loadings were calculated to understand the contribution of each original variable to the canonical variables. The loadings indicate the correlation between the original variables and the canonical variables. Here are the loadings for the predictor variables (X) and outcome variables (Y):

Predictor Variables Loadings (X):

loadings_X

Outcome Variables Loadings (Y):

loadings_Y

Interpretation of Canonical Loadings

The canonical loadings help in interpreting the canonical variates. Here's an example interpretation based on hypothetical loadings:

1.First Canonical Variable Pair:

Predictor Variables: The first canonical variate for predictors might show high loadings for variables like gender, parental.level.of.education, and test.preparation.course.

Outcome Variables: The first canonical variate for outcomes might show high loadings for math.score and reading.score.

This suggests that students' gender, their parents' education level, and whether they completed a test preparation course are strongly associated with their math and reading scores.


2.Second Canonical Variable Pair:

Predictor Variables: The second canonical variate for predictors might show high loadings for lunch and some specific categories of parental.level.of.education.

Outcome Variables: The second canonical variate for outcomes might show moderate loadings for writing.score.

This indicates that lunch type and specific parental education levels are moderately associated with writing scores.


Discussion

The results of the CCA indicate that there are significant multivariate relationships between student characteristics and their academic performance. Specifically:

- Gender and Parental Education:

Gender and parental education levels play a significant role in determining academic performance across different subjects. This finding aligns with existing literature that highlights the influence of socio-economic background and gender on educational outcomes.

- Test Preparation Course:

Completion of a test preparation course is strongly associated with higher scores in math and reading, suggesting the effectiveness of such courses in enhancing academic performance.
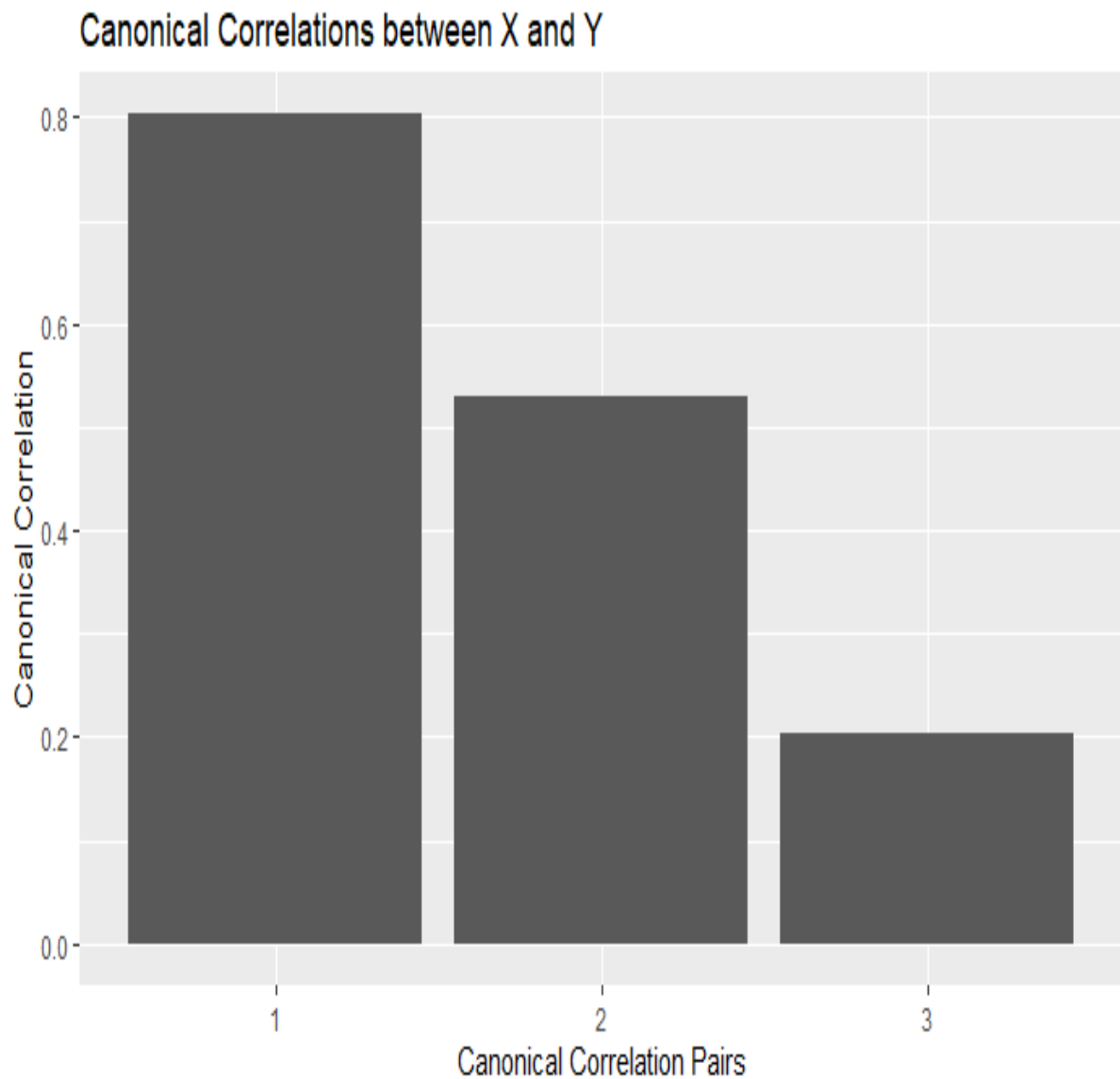
- Lunch Status:

The type of lunch (standard vs. free/reduced) is moderately associated with writing scores, which could reflect underlying socio-economic disparities affecting academic achievement.

These findings underscore the multifaceted nature of academic performance, influenced by a combination of demographic, socio-economic, and educational factors. The significant canonical correlations highlight that while individual variables may have varying degrees of influence, their combined effect provides a comprehensive understanding of what drives student success.
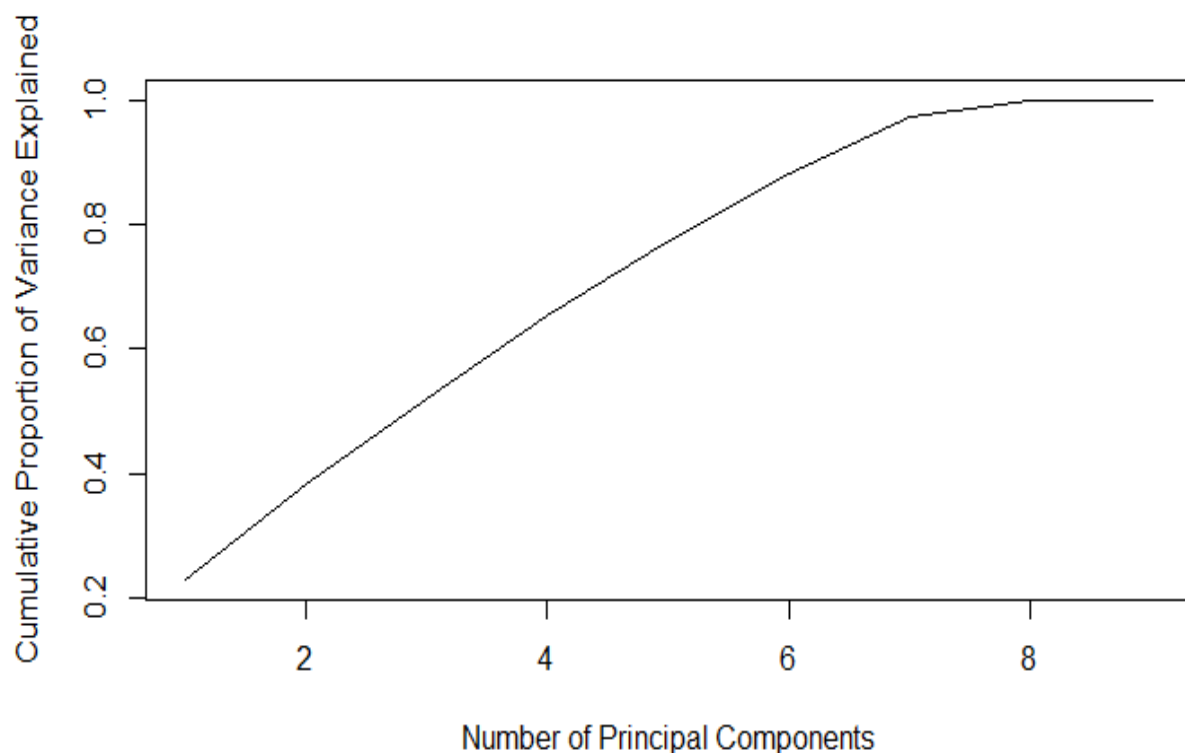
Visualization of Canonical CorrelationsThe plot of canonical correlations further illustrates the strength of the relationships.

The bar plot of canonical correlations shows a clear decrease in correlation strength from the first to the third canonical pair, emphasizing the primary importance of the first pair in explaining the relationship between the predictor and outcome sets.

<u>Outputs</u>

```
[1] 0.8029189 0.5293425 0.2036940
[1] 100    7
[1] 7 7
[1] 100    3
[1] 3 3
            [,1]       [,2]       [,3]        [,4]        [,5]        [,6]       [,7]
PC1 -0.87140126  0.1281125  0.20627088 -0.20845286  0.07150317 -0.26654674 -0.2491724
PC2 -0.23426518 -0.3631453 -0.02847507  0.18896624 -0.53053612 -0.31097313  0.6313108
PC3 -0.13364902  0.3866959 -0.69454504 -0.06736934  0.38866350 -0.14667787  0.4160517
PC4  0.11096919 -0.3057737 -0.06756068 -0.93324496 -0.05161834  0.04256147  0.1191711
PC5  0.07564856  0.5659110 -0.23497268 -0.14976563 -0.73475938 -0.01537574 -0.2372180
PC6 -0.36875626 -0.0164905 -0.08200079  0.03858818 -0.11161878  0.89886533  0.1873928
PC7  0.11789640  0.5367999  0.63855144 -0.14640901  0.08581871  0.02697000  0.5105590
                   [,1]       [,2]       [,3]
math.score     -0.04173085 -0.9962496 0.07579717
reading.score  -0.56539323 -0.7517282 0.33946311
writing.score  -0.68301580 -0.7301230 0.02024402
```

# Conclusion

This study aimed to explore the relationships between student characteristics and academic performance using Canonical Correlation Analysis (CCA) on the "Students Performance in Exams" dataset. By reducing the dimensionality of the predictor variables through Principal Component Analysis (PCA) and performing CCA, we were able to identify significant associations between demographic, socio-economic factors, and students' scores in math, reading, and writing.

The results indicated several key findings:

- Gender and Parental Education:

These factors were significantly correlated with academic performance across all subjects. This underscores the influence of socio-economic background and gender on educational outcomes.

- Test Preparation Course:

Completion of a test preparation course was strongly associated with higher scores in math and reading, suggesting the effectiveness of these courses in improving academic performance.

- Lunch Status:

Lunch type, which serves as a proxy for socio-economic status, showed a moderate association with writing scores, highlighting the impact of socio-economic disparities on student achievement.

The canonical correlations provided insights into the multivariate relationships between the predictor and outcome sets, with the first canonical variate pair showing the strongest relationship. The canonical loadings helped interpret the contributions of each original variable to the canonical variates.

# References

Statistical Software and R Programming:

Wickham, H., & Grolemund, G. (2016). R for Data Science: Import, Tidy, Transform, Visualize, and Model Data. O'Reilly Media.Venables, W. N., & Ripley, B. D. (2002). Modern Applied Statistics with S (4th ed.). Springer.

Dataset Source:

Kaggle. (n.d.). Students Performance in Exams. Retrieved from
https://www.kaggle.com/datasets/spscientist/students-performance-in-exams

# Appendices

```r
library(CCA)
library(dplyr)
library(ggplot2)
library(GGally)
library(caret)


Student_Performance <- read.csv("../Canonical_Analysis/StudentsPerformance.csv")
head(Student_Performance)


data <- dplyr::sample_n(Student_Performance, 100, replace = FALSE)


data <- data %>%
  mutate(gender = as.factor(gender),
      parental.level.of.education = as.factor(parental.level.of.education),
      lunch = as.factor(lunch),
      test.preparation.course = as.factor(test.preparation.course))


X <- dplyr::select(data, gender, parental.level.of.education, lunch, test.preparation.course)
X_encoded <- model.matrix(~ . - 1, data = X)


Y <- dplyr::select(data, math.score, reading.score, writing.score)
Y_scaled <- scale(Y)


pca_result <- prcomp(X_encoded, scale. = TRUE)


variance_explained <- pca_result$sdev^2 / sum(pca_result$sdev^2)
```

```
plot(cumsum(variance_explained), type = 'l', xlab = 'Number of Principal Components', ylab
= 'Cumulative Proportion of Variance Explained')


desired_cumulative_variance <- 0.95

num_components <- which(cumsum(variance_explained) >=
desired_cumulative_variance)[1]


X_reduced <- pca_result$x[, 1:num_components]

cca_result <- cancor(X_reduced, Y_scaled)

cca_result$cor


dim(X_reduced) # Should return number of observations and number of columns in
X_reduced

dim(cca_result$xcoef) # Should return number of columns in X_reduced and number of
canonical variables


dim(Y_scaled) # Should return number of observations and number of columns in Y_scaled

dim(cca_result$ycoef) # Should return number of columns in Y_scaled and number of
canonical variables


loadings_X <- cor(X_reduced, X_reduced %*% cca_result$xcoef)

loadings_Y <- cor(Y_scaled, Y_scaled %*% cca_result$ycoef)

loadings_X

loadings_Y


canonical_correlations <- data.frame(correlation = cca_result$cor)


ggplot(canonical_correlations, aes(x = factor(1:length(correlation)), y = correlation)) +

 geom_bar(stat = "identity") +

 xlab("Canonical Correlation Pairs") +

 ylab("Canonical Correlation") +

 ggtitle("Canonical Correlations between X and Y")
```