Introduction
00

Data Visualization
00000000

Dimensionality Reduction
00

Regression
00000

Conclusion
000

References
00

# Mathematics Performance in Secondary Education

Yao-Chih Hsu, Xuan-Chun Wang, and Wen-Lee Sin

June 12, 2025

Introduction
oo

Data Visualization
oooooooo

Dimensionality Reduction
oo

Regression
ooooo

Conclusion
ooo

References
oo

# Outline

1 Introduction

2 Data Visualization

3 Dimensionality Reduction

4 Regression

5 Conclusion

6 References

## Introduction

- Dataset: *Math-Students Performance Data* from Kaggle (Shamim, 2025).

## Introduction

- Dataset: *Math-Students Performance Data* from Kaggle (Shamim, 2025).
- Variables G1, G2, G3, and absences were provided by the school.
- Remaining variables were collected via questionnaires and are mostly categorical.
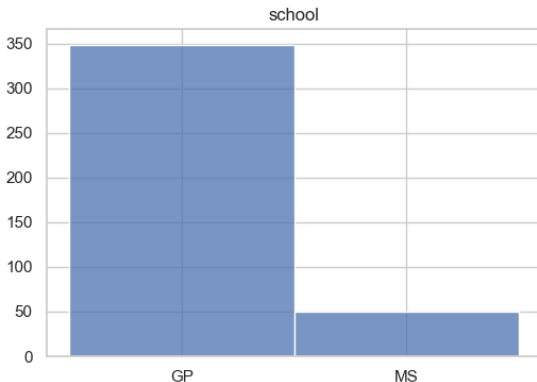
## Goal

Analyze G3 to identify influential variables.

Introduction
○○
Data Visualization
●○○○○○○○
Dimensionality Reduction
○○
Regression
○○○○○
Conclusion
○○○
References
○○

## Data Grouping

We grouped the data into the following:

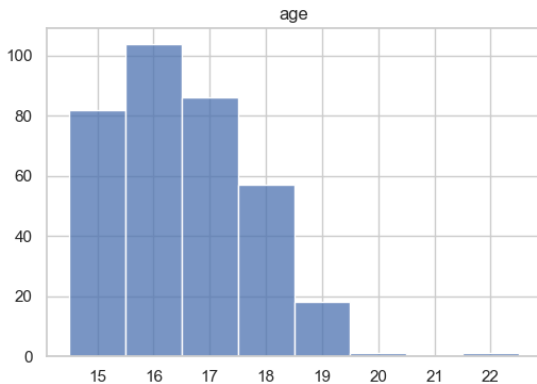| Groups | Variables |
|--------|-----------|
| support | schoolsup, famsup, paid |
| family | address, famsize, Pstatus, guardian, traveltime, famrel |
| parents | Medu, Fedu, Mjob, Fjob |
| performance | failures, studytime, absences |
| alcohol | Dalc, Walc, health |
| after_class | activities, freetime, goout |
| school_choice | reason, nursery, higher |
| score | G1, G2, G3 |

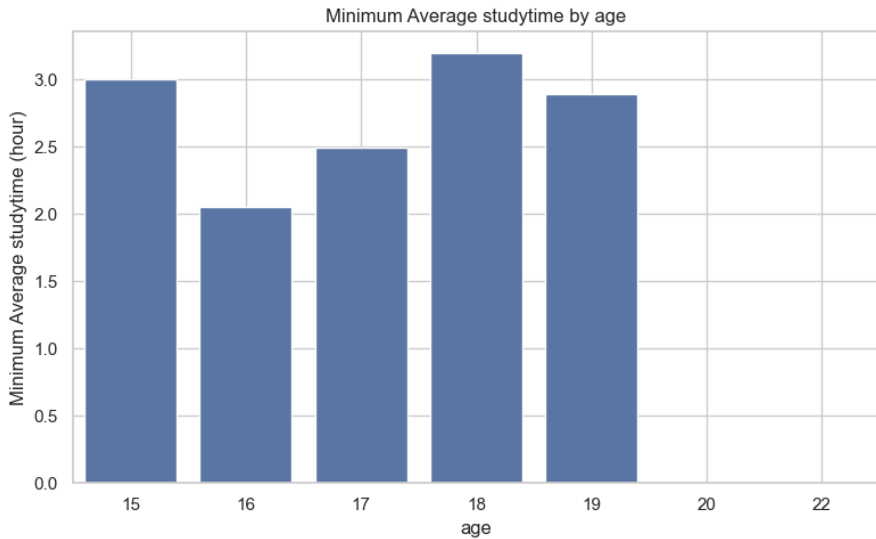The variables not yet assigned to any group are:
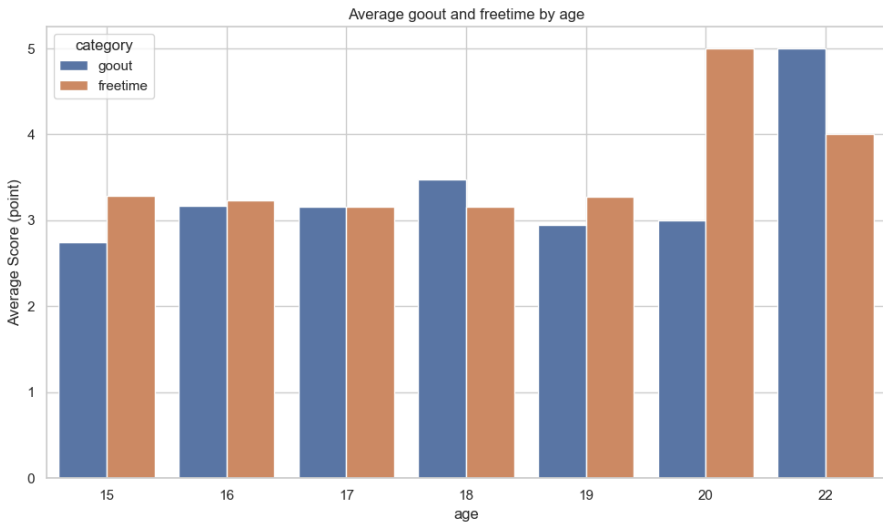
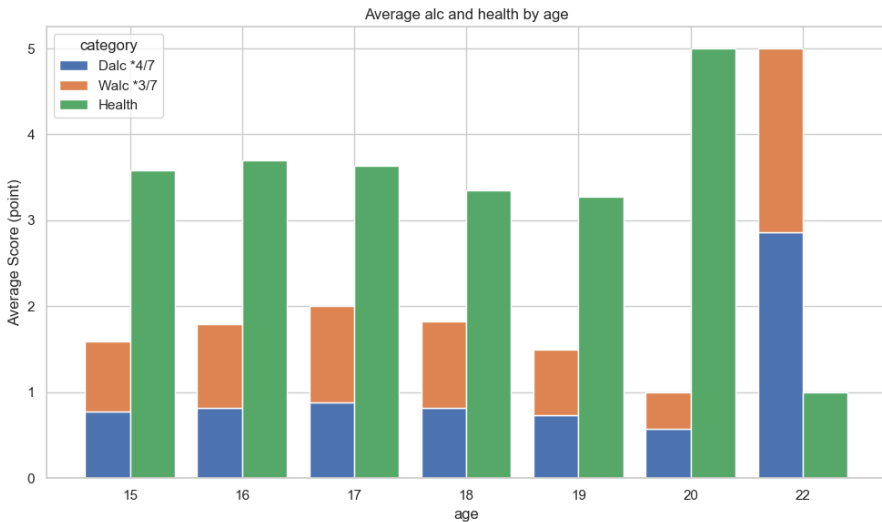sex, age, internet, romantic.
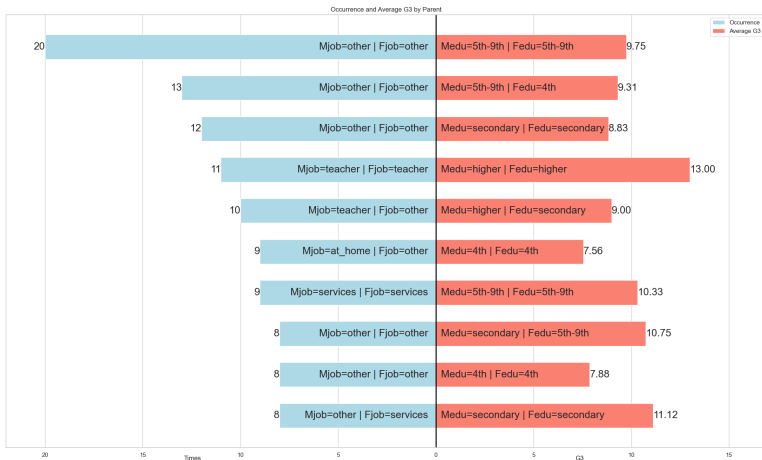
## Data Visualization



**GP (Gabriel Pereira)    MS (Mousinho da Silveira)**

Introduction
00

Data Visualization
00000●0000

Dimensionality Reduction
00

Regression
00000

Conclusion
000

References
00

Minimum Average studytime by age

Introduction
oo

Data Visualization
ooooo●ooo

Dimensionality Reduction
oo

Regression
ooooo

Conclusion
ooo

References
oo

Average goout and freetime by age

Introduction
oo

Data Visualization
ooooooo●oo

Dimensionality Reduction
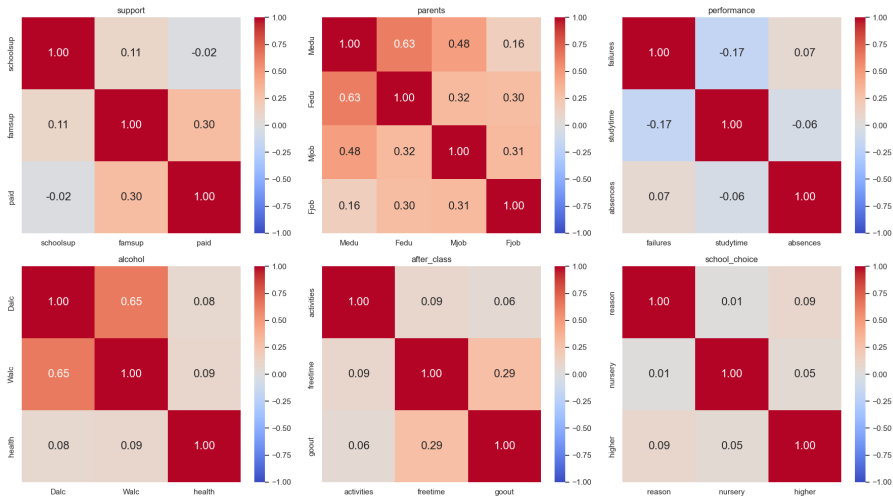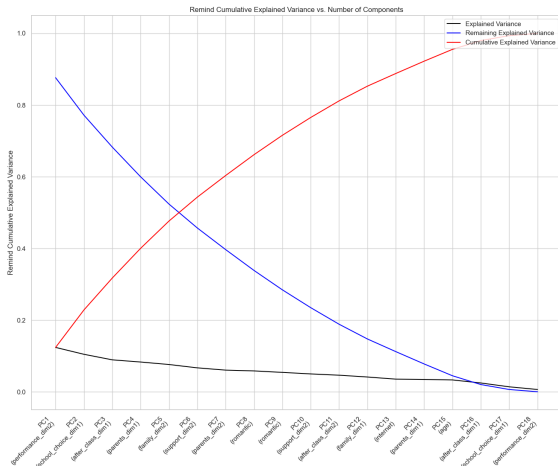oo

Regression
ooooo

Conclusion
ooo

References
oo

Top 10 Parent Combinations by Highest Average G3

# Multidimensional Scaling

Introduction
oo

Data Visualization
oooooooo

Dimensionality Reduction
o●

Regression
ooooo
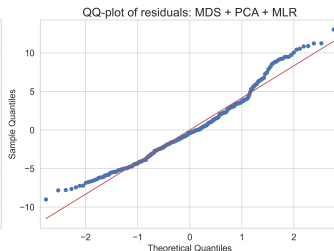
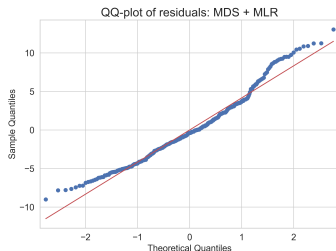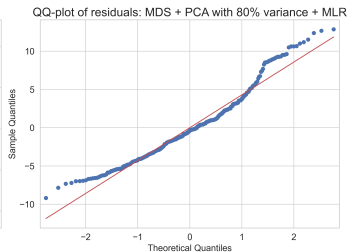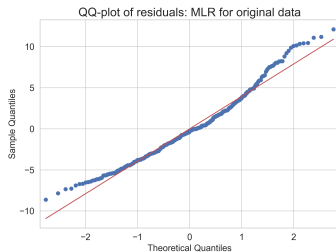Conclusion
ooo

References
oo

# PCA

## Variables of MLR

We compared four models:

- Original + MLR
- MDS + MLR
- MDS + PCA + MLR
- MDS + PCA 80% + MLR

The variables in the dataset transformed by MDS are as follows:

# Q-Q plots of MLR

## Variables of Reduced MLR

| Models | Selected Features |
|---|---|
| Original + MDS + Stepwise + MLR | `failures`, `goout`, `parents_dim1`, `after_class_dim1` `family_dim2`, `romantic` `support_dim2`, `after_class_dim2` `support_dim1`, `age` `performance_dim1` |
| MDS + PCA + Stepwise + MLR | `parents_dim1`, `after_class_dim1` `family_dim2`, `romantic` `support_dim2`, `after_class_dim2` `support_dim1`, `age` `performance_dim1` |

Table 1: Selected Features by Different Methods

Introduction
oo

Data Visualization
ooooooooo

Dimensionality Reduction
oo

Regression
oooeo

Conclusion
ooo

References
oo

## Original data + MDS + Stepwise + PCA

Introduction
oo

Data Visualization
oooooooo

Dimensionality Reduction
oo

Regression
ooooo●

Conclusion
ooo

References
oo

PCA: PC1 vs PC2

SW PCA: PC1 vs PC2

G3 [0, 5]
G3 [6, 14]
G3 [15, 20]

**PC1 (`school_choice_dim1`)**
**PC2 (`after_class_dim1`)**

**PC1 (`parent_dim2`)**
**PC2 (`romantic`)**

Introduction
00
Data Visualization
00000000
Dimensionality Reduction
00
Regression
00000
Conclusion
●00
References
00

## Conclusion

| Models | $R^2$ | Adj. $R^2$ | MSE | P-value | AIC |
|---|---|---|---|---|---|
| MLR for original data | **0.27** | **0.19** | **62.58** | 0.9157 | 2015.10 |
| MDS + MLR | 0.19 | 0.14 | 77.58 | 0.7494 | 2023.96 |
| MDS + PCA + MLR | 0.19 | 0.14 | 77.58 | 0.6143 | 2023.96 |
| MDS + PCA with 80% variance + MLR | 0.14 | 0.11 | 92.69 | 0.9442 | 2031.05 |
| Original data + MDS + Stepwise + MLR | 0.20 | **0.19** | 293.55 | **0.0367** | **1993.82** |
| MDS + PCA + Stepwise + MLR | 0.17 | 0.15 | 140.12 | 0.5113 | 2013.68 |
| Paper Proposed | 0.17 | 0.16 | 216.03 | 0.4624 | 2005.069 |

Table 2: Comparison of model performance metrics

The variables selected by the paper:
absences, schoolsup, higher, failures, Mjob

## Conclusion

| Model | TOP1 | TOP2 | TOP3 |
|---|---|---|---|
| **MLR for original data** | `failures` | `schoolsup` | `paid` |
| MLR + MLR | `romantic` | `support_dim2` | `support_dim1` |
| MDS + PCA + MLR | `romantic` | `support_dim2` | `age` |
| MDS + PCA with 80% variance + MLR | `romantic` | `support_dim2` | `parents_dim1` |
| **Original data + MDS + Stepwise + MLR** | `Mjob_at_home` | `failures` | `Mjob_other` |
| MDS + PCA + Stepwise + MLR | `age` | `support_dim1` | `family_dim2` |
| **Paper Proposed** | `higher` | `schoolsup` | `failures` |

Table 3: Most important feature to each model

## Conclusion

The variables we selected are:

$$Mjob, \texttt{failures}$$

# References

📄 Cortez, P. (2008). *Student Performance* [Dataset]. UCI Machine
   Learning Repository. `https://doi.org/10.24432/C5TG7T`

📄 Cortez, P., & Silva, A. M. (2008). Using data mining to predict
   secondary school student performance.

📄 Shamim, A. (2025). *Math students performance data*. Kaggle. `https://www.kaggle.com/datasets/adilshamim8/math-students`

Thanks for listening!