



Danta
analytics

Caso de Análisis

Por Ronaldo Vindas Barboza

Agenda

Introducción

Metodología

Análisis Exploratorio Inicial de Características

Análisis Estadístico de Características

Análisis Visual de Características

Comportamiento de Datos

Modelo de Predicción

Hallazgos Importantes

Introducción

OBJETIVO:

Analizar el desempeño de los estudiantes en exámenes de matemáticas, lectura y escritura, utilizando el dataset "Students Performance in Exams" de Kaggle. A través de este análisis, se busca identificar aquellos patrones y factores que influyen en los resultados de los estudiantes.

kaggle

El Conjunto de Datos

Incluye las características:

- Género
- Grupo Étnico
- Nivel de Educación de Padres
- Tipo de Almuerzo
- Curso de Preparación
- Calificaciones

Importancia

Comprender qué factores afectan el rendimiento de las pruebas que puedan ayudar a desarrollar estrategias que mejoren los resultados de los estudiantes.



Metodología

LECTURA DE DATOS

PREPROCESAMIENTO DE DATOS

ANÁLISIS EXPLORATORIO
VISUAL Y ESTADÍSTICO

ANÁLISIS DE RESULTADOS

Lectura de Datos

Los datos fueron leídos por medio de la biblioteca Pandas de Python y se les almacenó en un Dataframe.



Label Encoding

| Food Name | Categorical # | Calories |
|-----------|---------------|----------|
| Apple | 1 | 95 |
| Chicken | 2 | 231 |
| Broccoli | 3 | 50 |



One Hot Encoding

| Apple | Chicken | Broccoli | Calories |
|-------|---------|----------|----------|
| 1 | 0 | 0 | 95 |
| 0 | 1 | 0 | 231 |
| 0 | 0 | 1 | 50 |

Preprocesamiento de Datos

Se aplicó One-Hot-Encoding para transformar features categóricos a numéricos y Normalización para establecer un intervalo de valores $[0, 1]$.

Análisis Exploratorio y Comportamiento de Datos

Análisis Exploratorio

VISUAL

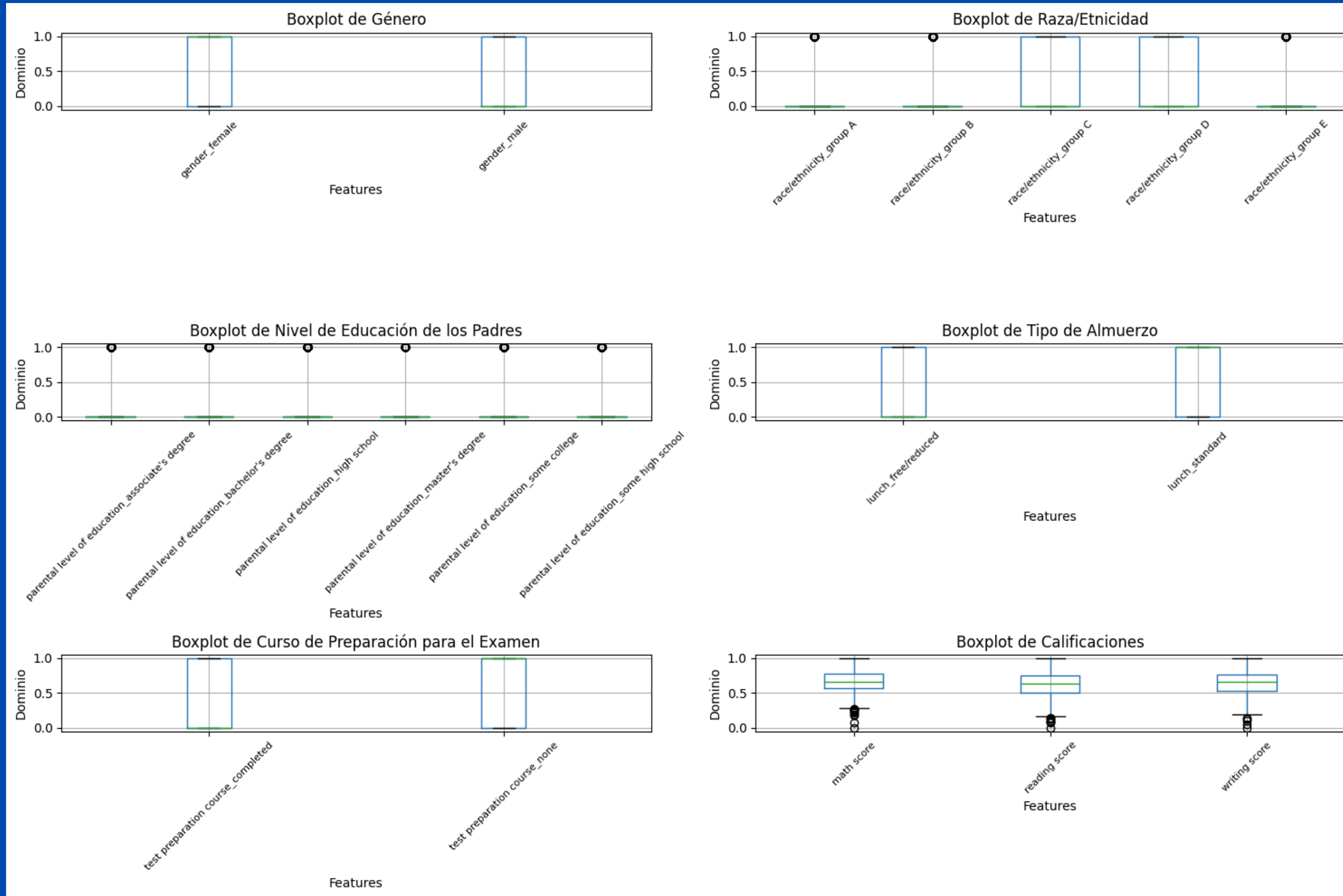
- Gráficos de Caja (BoxPlot)
- Gráfico de Dispersión
- Gráficos de Histogramas
- Matriz de Correlación

ESTADÍSTICOS

- Promedio
- Desv. Estándar
- Rangos de Cuartiles



Gráficos de Caja



Comportamiento

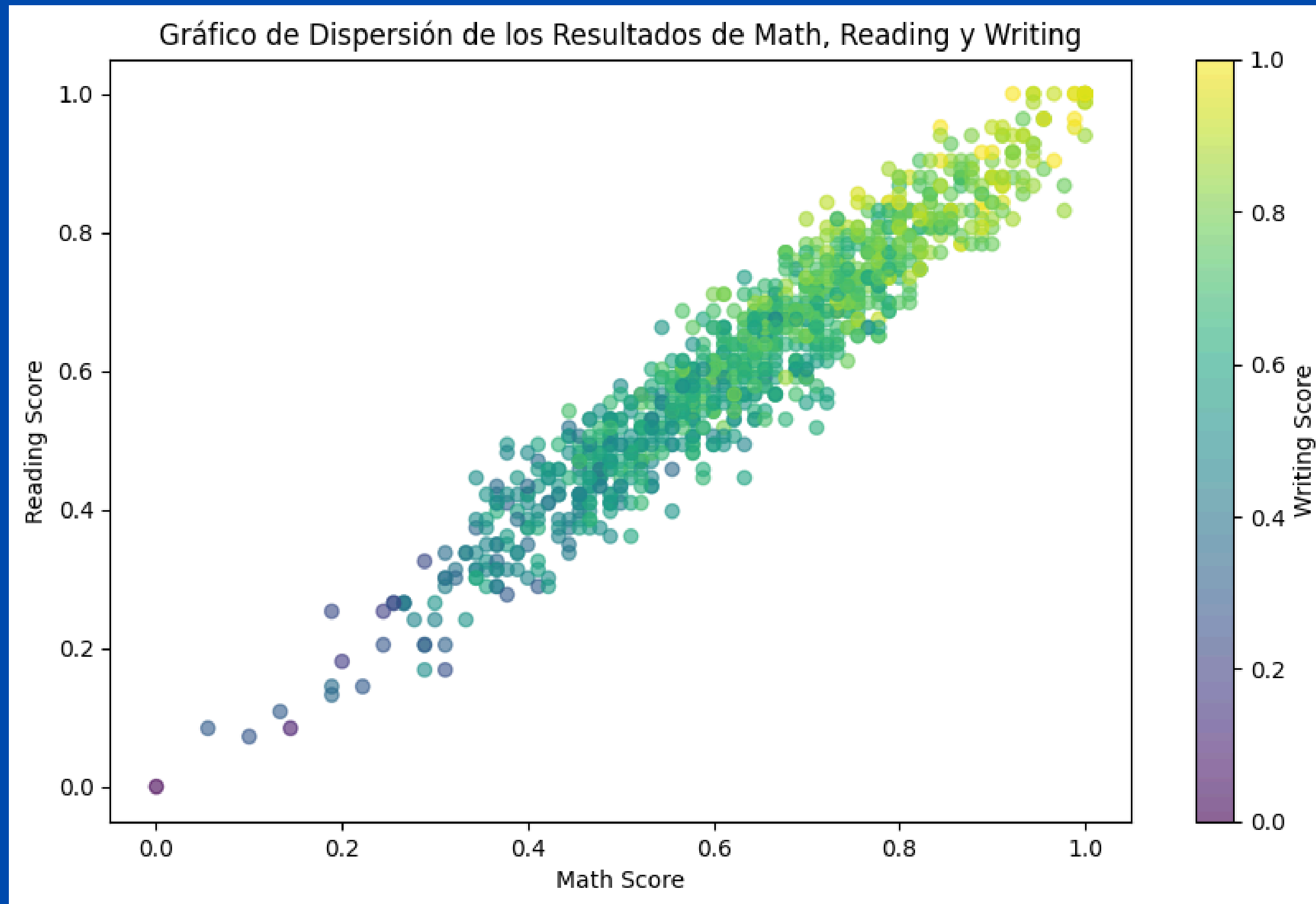
El Análisis Visual muestra:

Según los Gráficos de Cajas:

- Se presentan valores atípicos (**outliers**) principalmente en los grupos étnicos A, B y E así como en los niveles de educación de los padres. Se podrían trabajarán métodos de imputación de datos para reducir su impacto en el modelo de predicción.
- Las características de calificaciones presentan también valores atípicos, estos a causa de **variabilidad natural**, y por tanto se tomarán en cuenta en el modelo de predicción.



Gráfico de Dispersión



Comportamiento

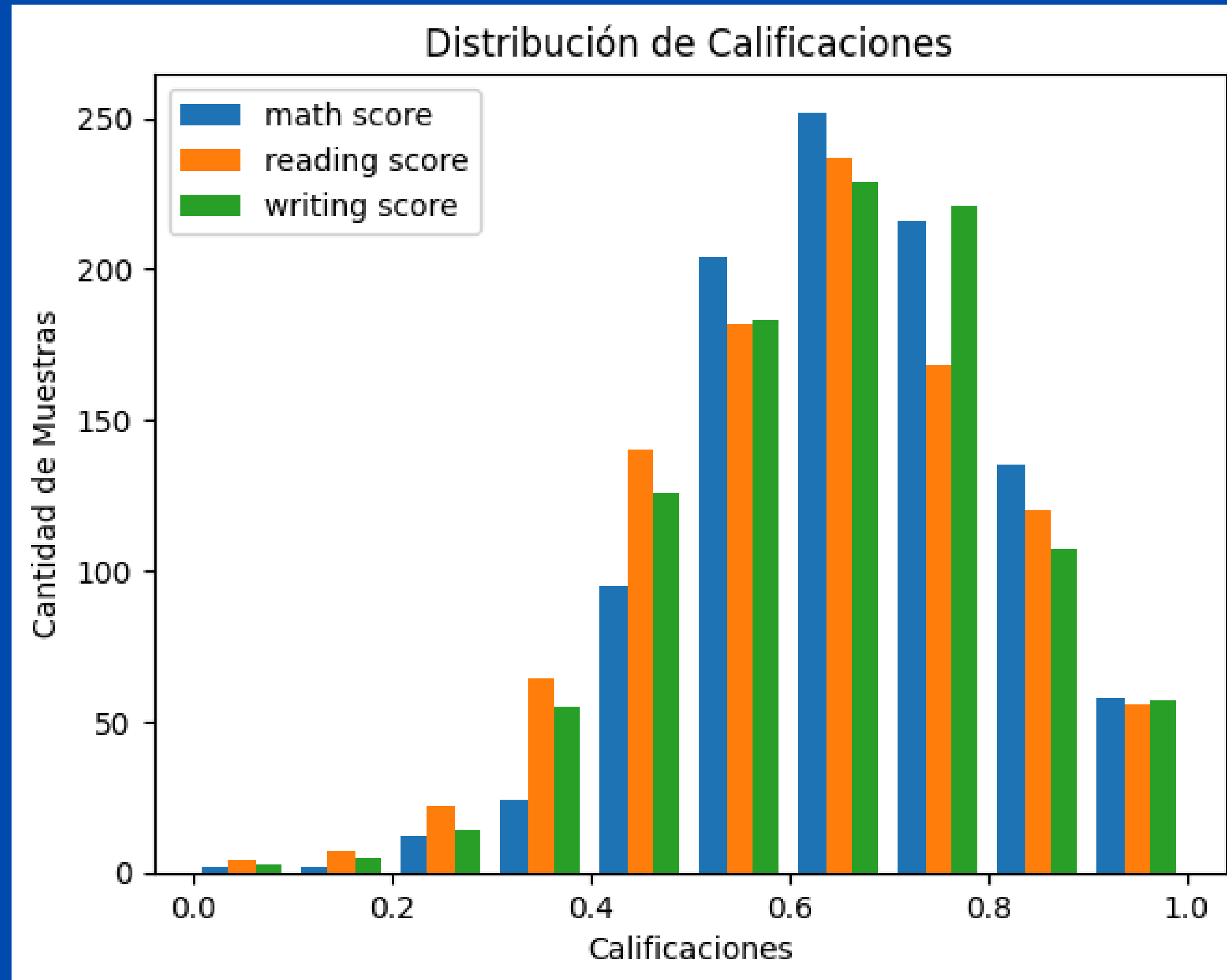
Análisis Visual:

De acuerdo al Gráfico de Dispersión:

- Los estudiantes que suelen **pasar** un examen de una materia suelen hacerlo también con las otras dos.
- Aquellos estudiantes que fallan un examen tienden a **fallar** los demás.
- La mayoría de estudiantes obtienen **notas regulares**.



Gráfico de Histogramas



Comportamiento

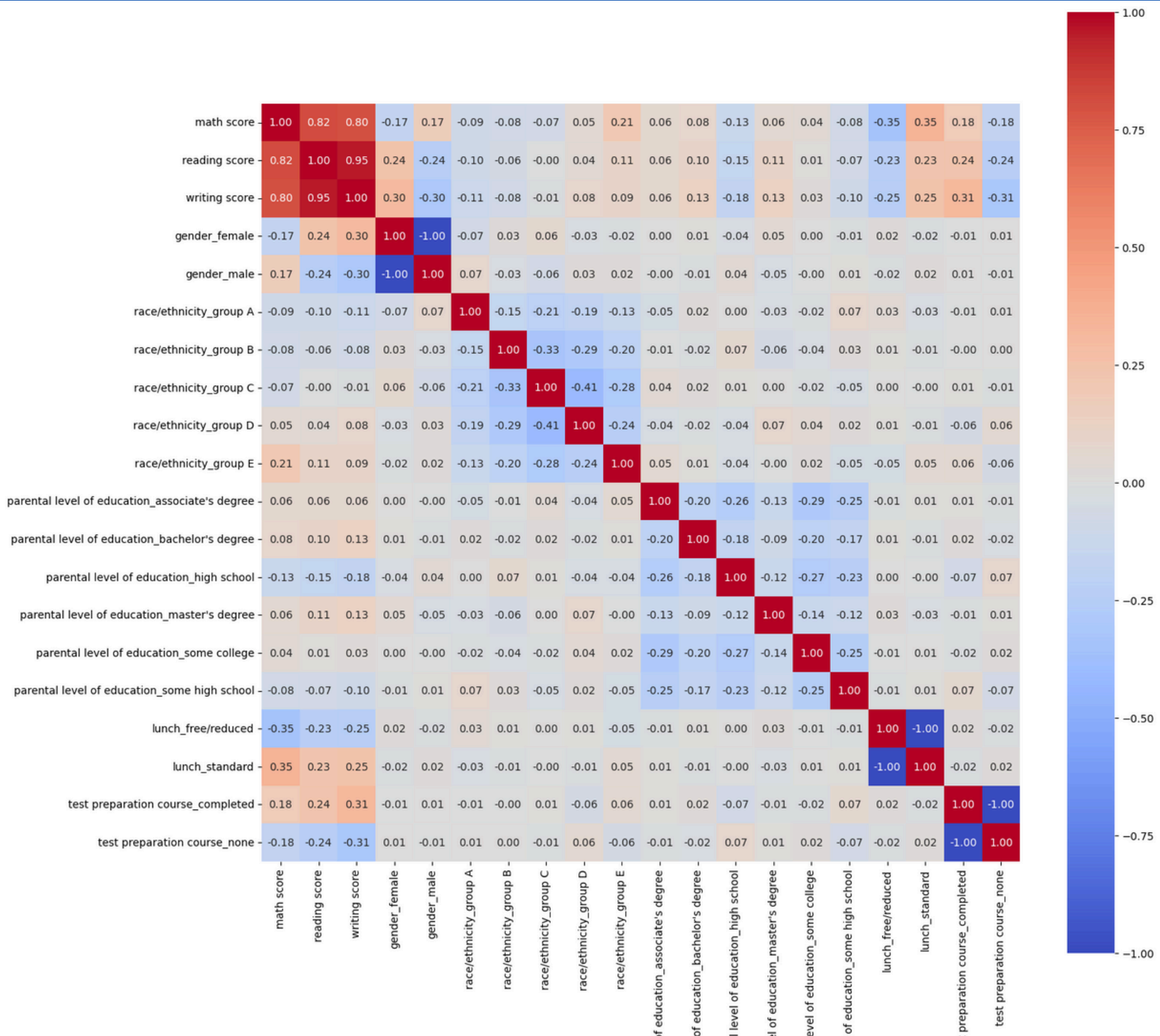
Análisis Visual:

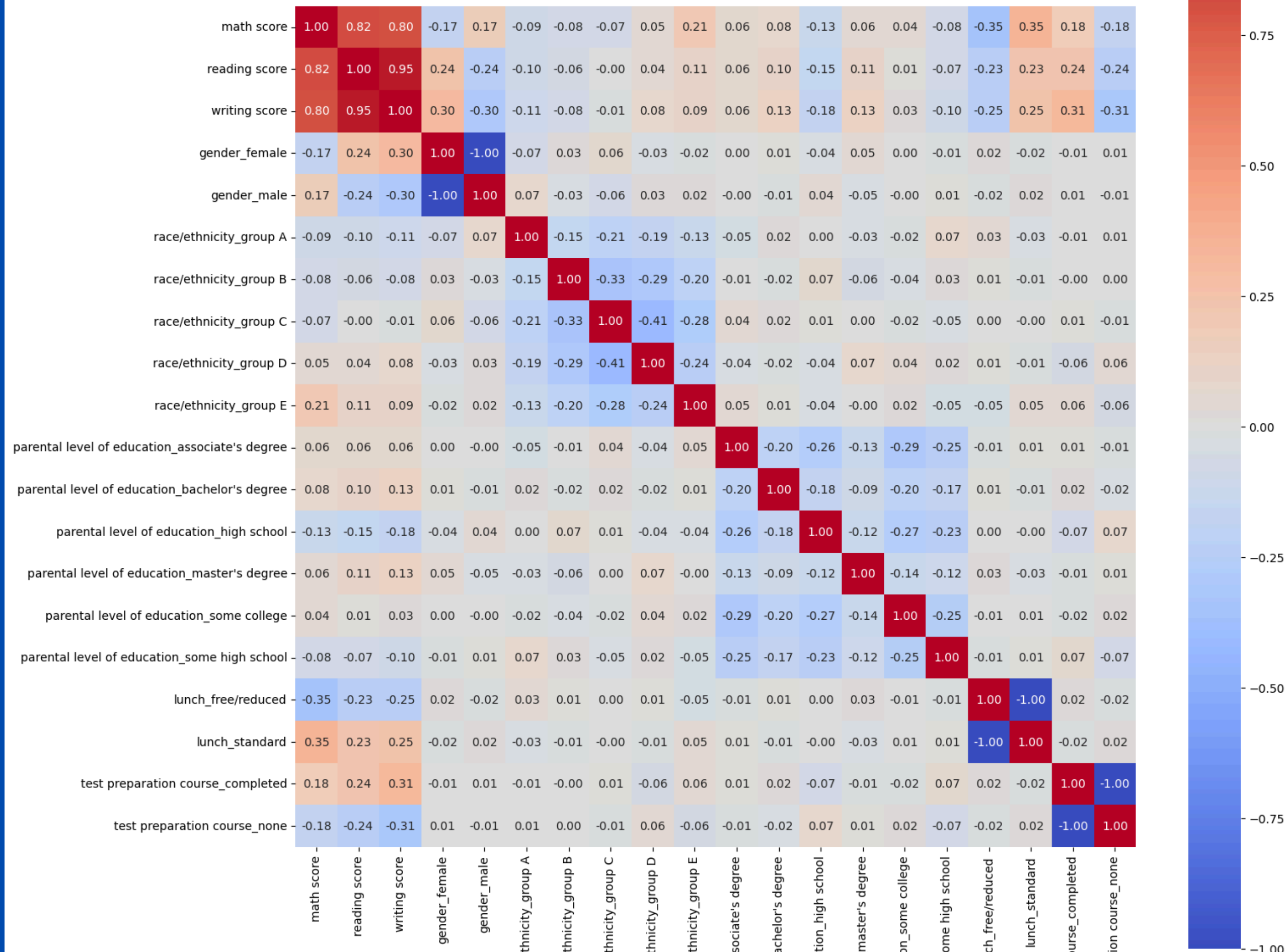


El gráfico de Histogramas nos dice que:

- Las muestras de puntuaciones siguen una **distribución normal**.
- Aproximadamente solo un 0.5% de los estudiantes sacó una **nota mayor** o igual a 90 en los tres exámenes.
- Menos de 25 estudiantes obtuvieron una **nota menor** a 20 en los tres exámenes.

Matriz de Correlación





Comportamiento

Análisis Visual:



El Gráfico de Correlación arroja que:

- Sobresale el hecho de que la **calificación** en una materia se relaciona mucho con las demás, esto podría traducirse en que al igual a como se salga en una prueba, saldrá similar en las demás.
- El prepararse para los exámenes parece influir positivamente en el resultado de estos de manera baja/moderada.

Comportamiento

Análisis Visual:



El Gráfico de Correlación arroja que:

- Características como el **tipo de almuerzo** parece tener una correlación positiva y negativa **moderada** con las puntuaciones de los exámenes. Aquellos estudiantes que tienen un almuerzo estándar puede que salgan un poco mejor que aquellos que no almuerzan o que almuerzan poco.
- El **estudio de los padres** no parece influir demasiado con las calificaciones de los exámenes.

Comportamiento

Análisis Visual:



El Gráfico de Correlación arroja que:

- Las **mujeres** tienden a salir un poco mejor en Escritura y Lectura que los hombres.
- Los **hombres** suelen salir un poco mejor en Matemáticas que las mujeres.
- La mayoría de grupos étnicos no influyen mucho en el resultado de los exámenes, aunque el **grupo étnico E** pareciera influir de manera baja/moderada sobre las calificaciones.

Análisis Estadístico

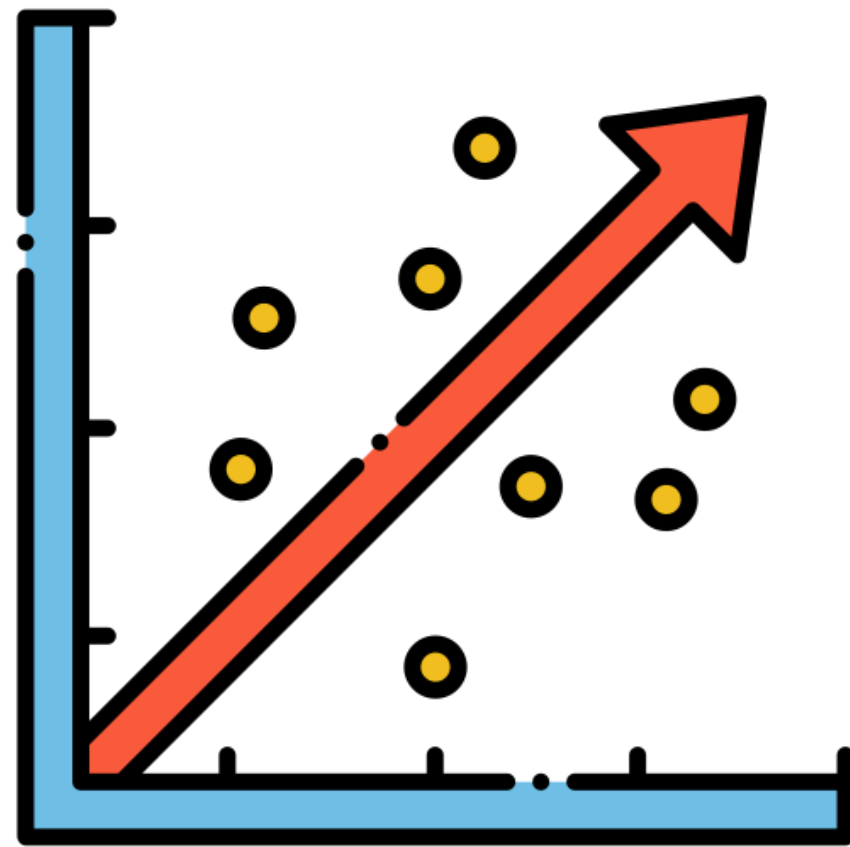
Comportamiento

Análisis Estadístico:



- En su mayoría, los estudiantes tienen calificaciones regulares (alrededor de los 60 pts).
- La media de estudiantes que no se preparan ante un examen es alta, aproximadamente un 64% de ellos.
- La mayoría de estudiantes pertenecen al grupo étnico C.
- Las características de "almuerzo" y "preparación" se encuentran algo desbalanceadas,. lo que podría llegar a sesgar el modelo predictivo.

Modelo de Predicción



Regresión Lineal

Uso de Biblioteca Scikit-Learn

Técnica 80/20.

Predecir "Math Score".

Evaluado mediante:

- Mean Squared Error
- Coeficiente de Determinación (R^2 Score)
- Coeficientes de Variables
- Gráficos
(Evaluación, Error y Coeficientes)



Experimento Adicional

Igual al modelo anterior pero entrenado con datos limpios de valores atípicos que puedan “entorpecer” el modelo.

Evaluación del Modelo



MEAN SQUARED ERROR (MSE)

0.0026945732490891833

R² SCORE

0.8788700517507204

Gráfico de Evaluación

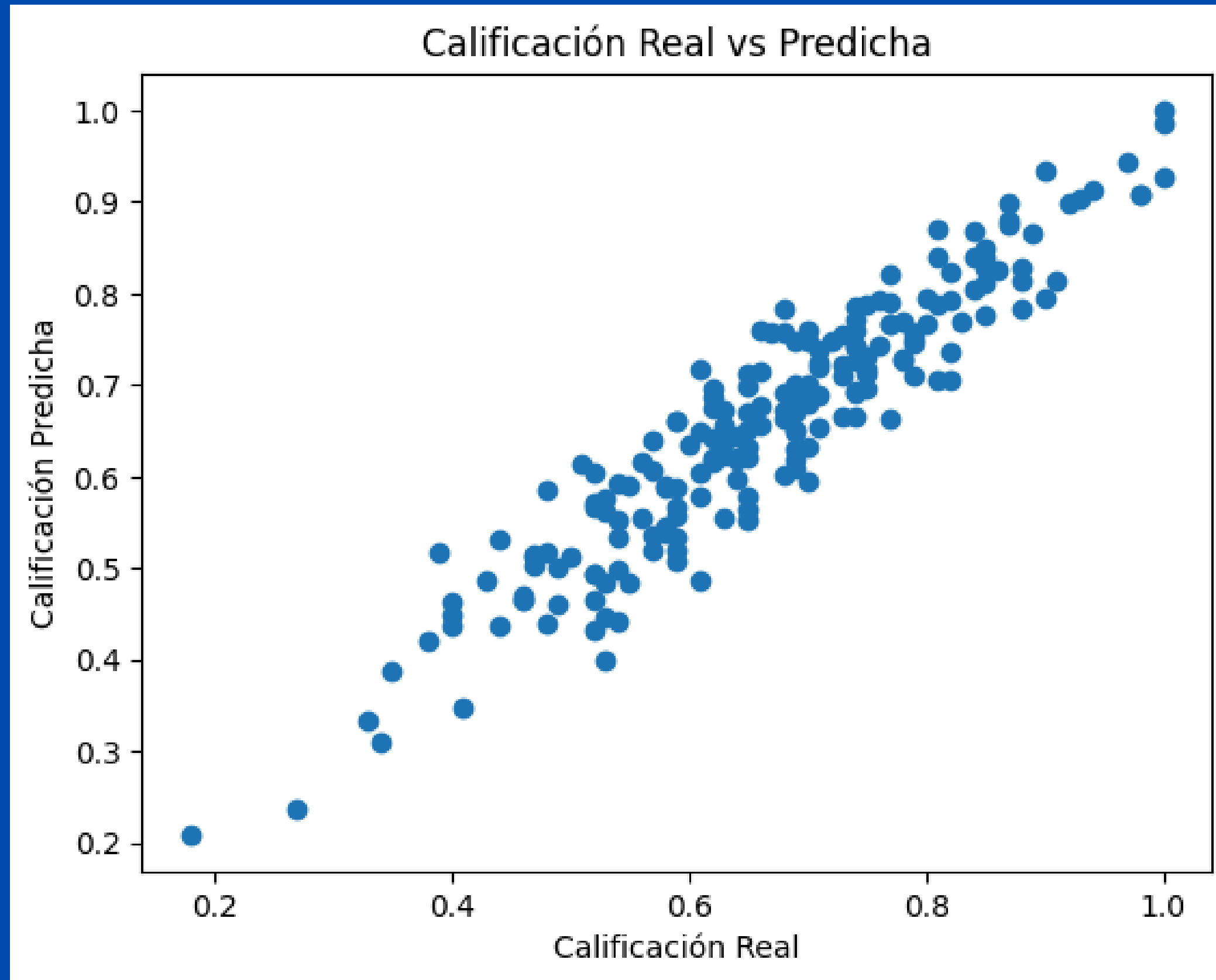


Gráfico de Error

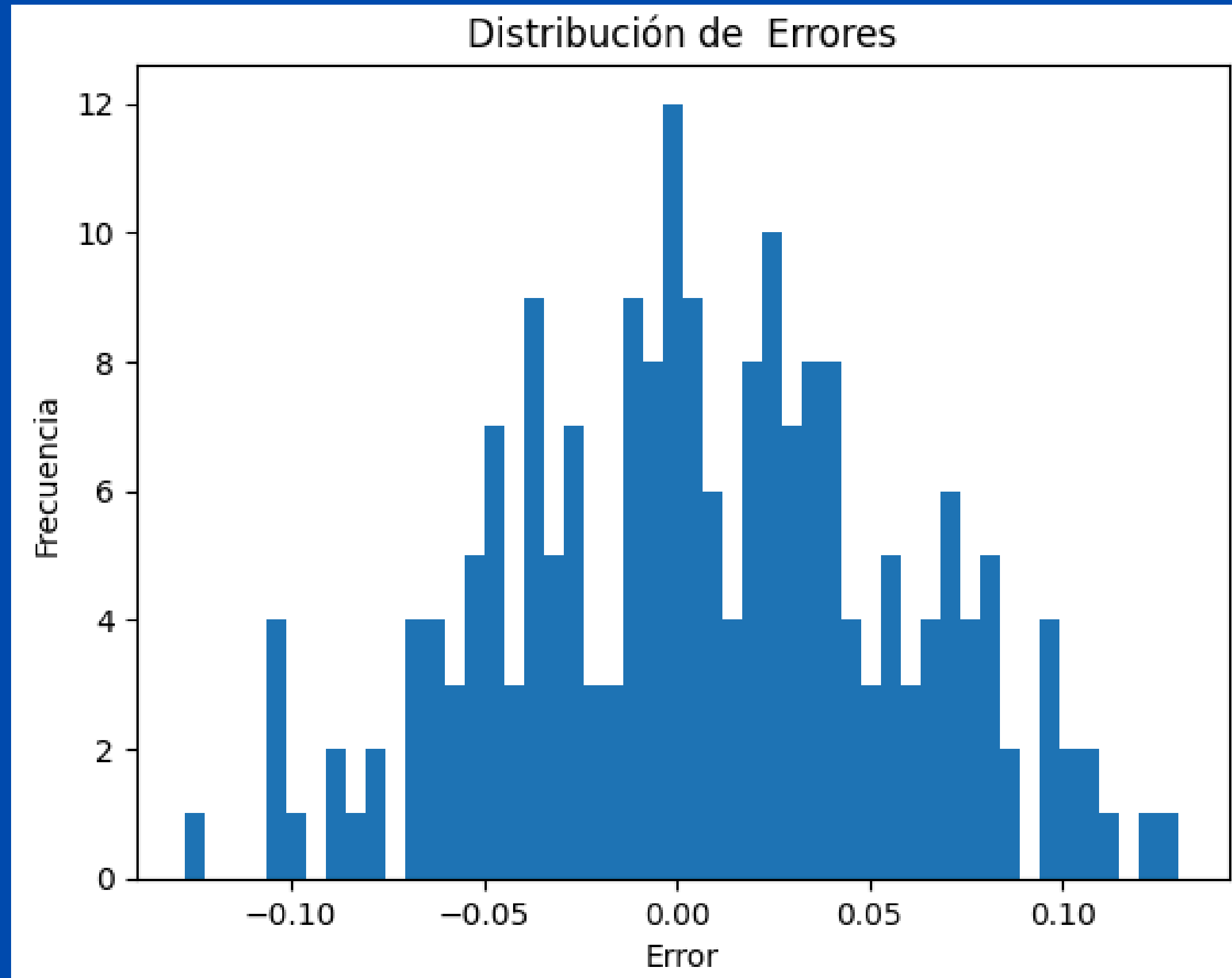
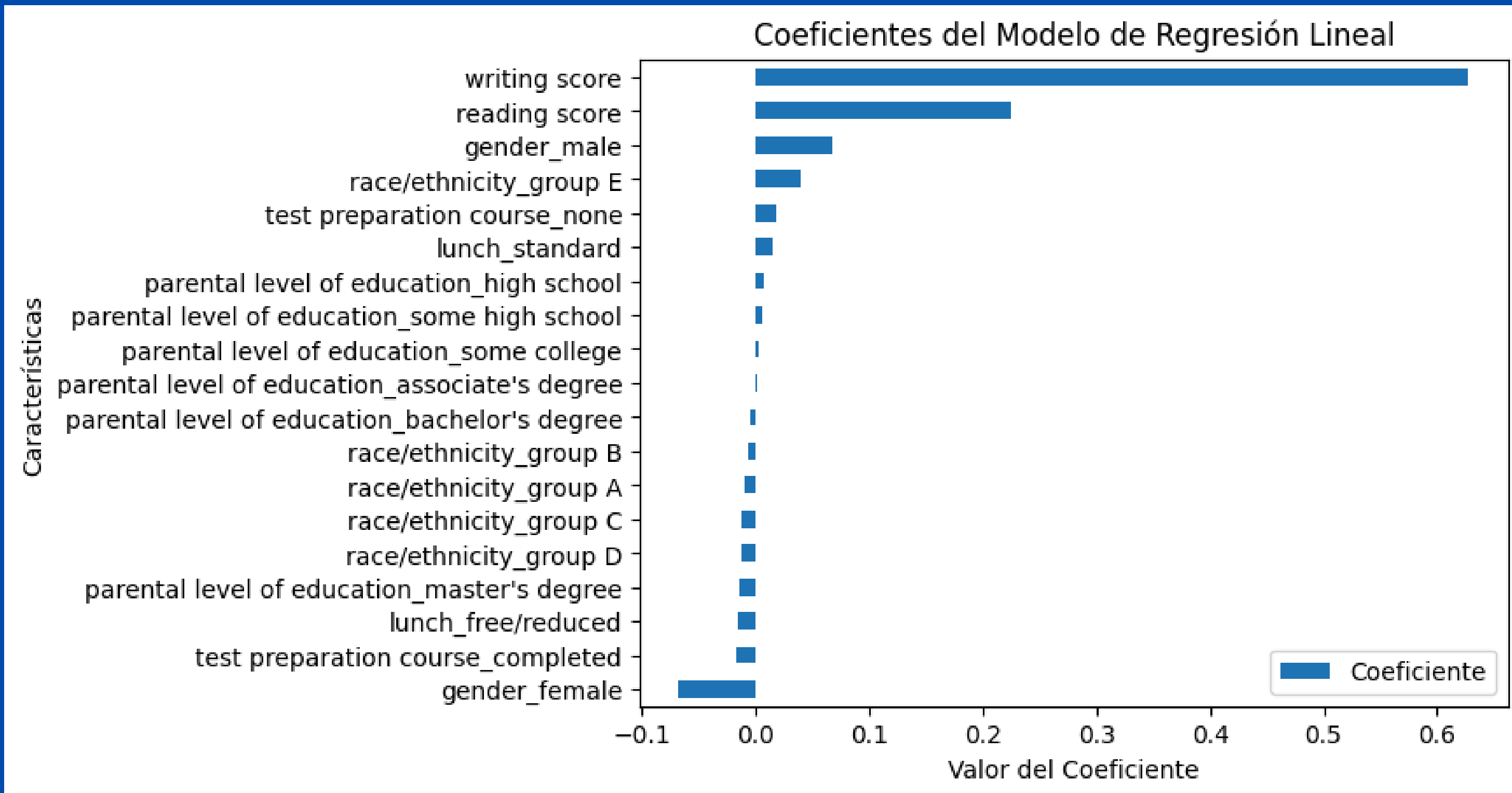


Gráfico de Coeficientes



Evaluación del Modelo Experimento



MEAN SQUARED ERROR (MSE)

0.0028874734116149433

R² SCORE

0.8439658801376765

Gráfico de Evaluación

Experimento Adicional

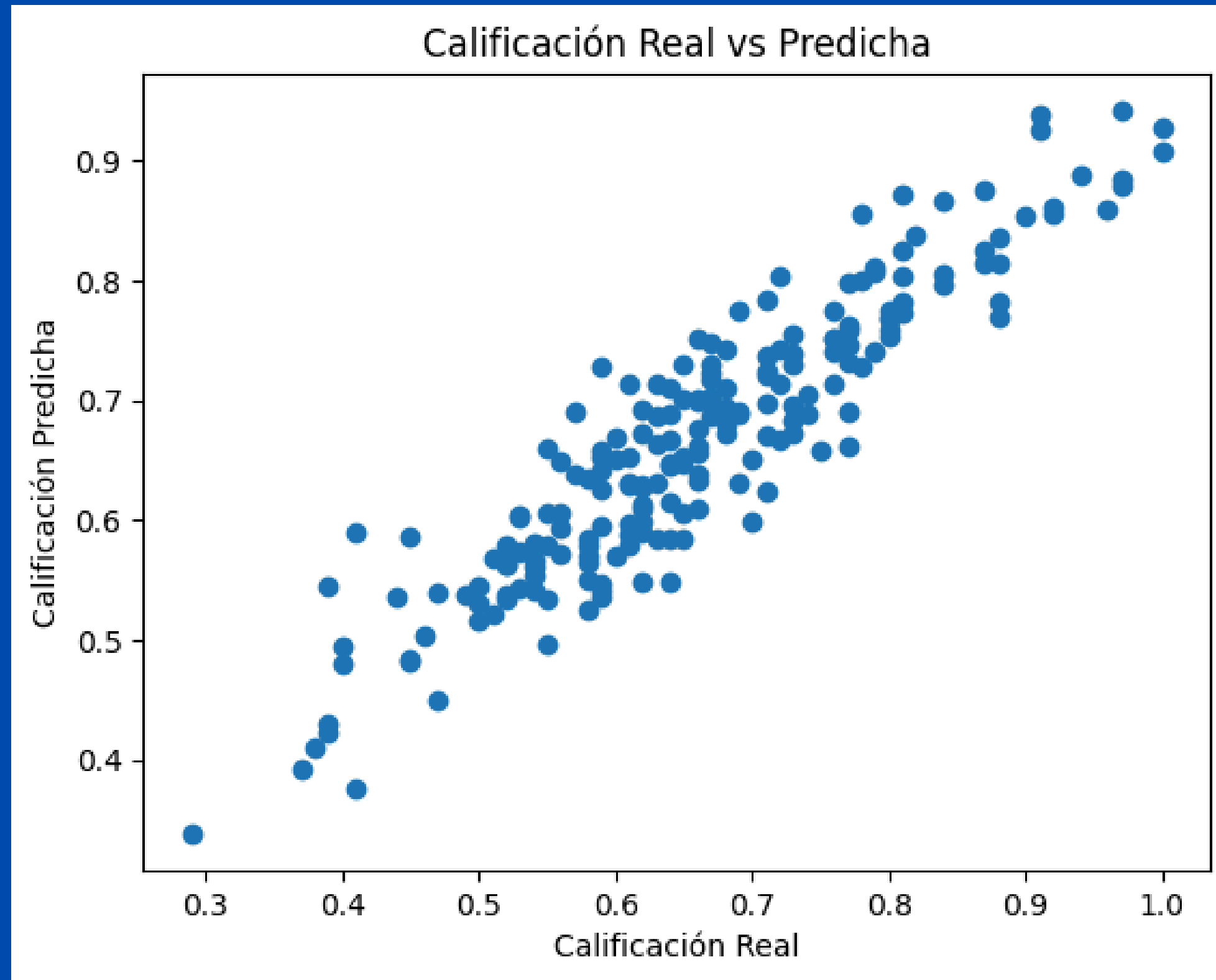


Gráfico de Error

Experimento Adicional

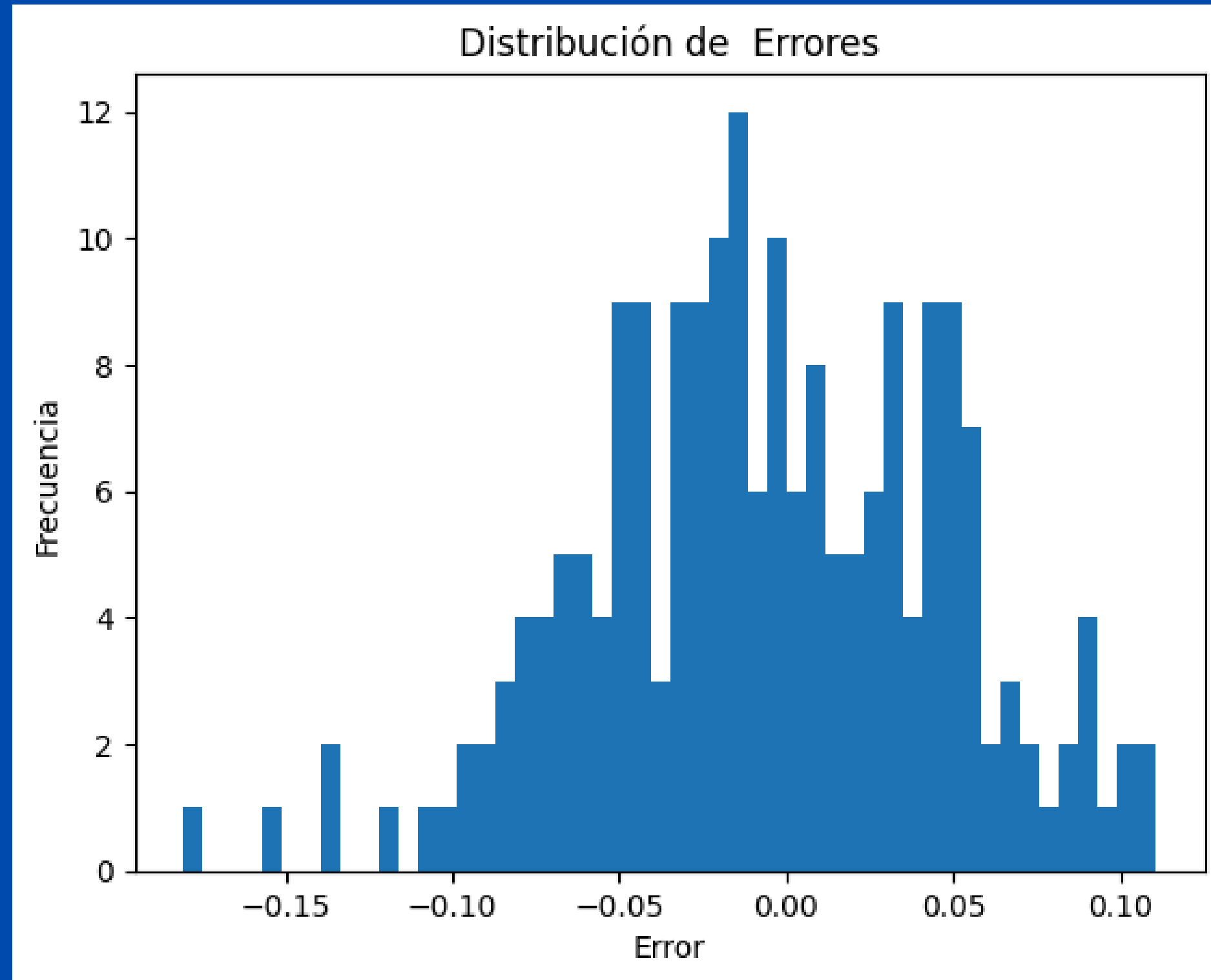
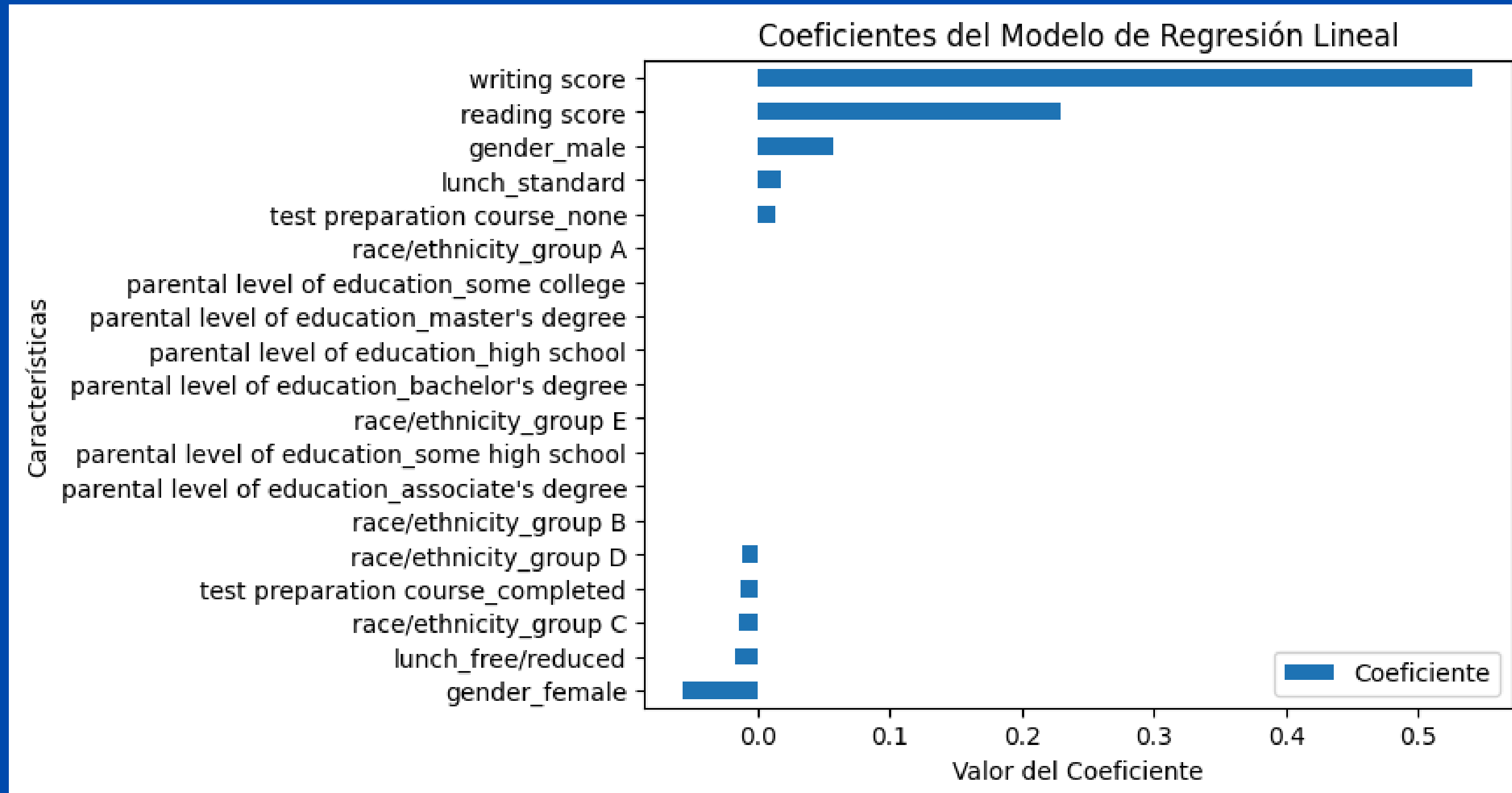


Gráfico de Coeficientes

Experimento Adicional



Hallazgos Clave

Resultados muestran:

De acuerdo a los resultados arrojados por la evaluación:

- La proporción de **error** del MSE en ambos modelos es **muy bajo**, lo que indica que estos modelos pueden llegar a predecir muy bien.
- El Coeficiente de Determinación (R^2 Score) tiene un valor de **ajuste muy bueno** sin llegar a un overfitting, además menciona que el 87% de la varianza en la calificación de matemáticas se puede explicar con las variables independientes.



Hallazgos Clave

Resultados muestran:

Según los gráficos de evaluación:

- El gráfico de dispersión muestra que los valores predichos siguen una distribución normal muy **similar** a la del dataframe **original** concentrando la mayoría de calificaciones alrededor de los 60 puntos. Además, este gráfico muestra la capacidad de aproximarse muy cercanamente al valor real de la calificación.
- El gráfico de distribución de errores muestra una **distribución** muy **simétrica** y con una varianza razonable aunque muestra aún algunos valores atípicos, esto podría deberse al haberse utilizado aún todos los features para el entrenamiento del modelo. Por lo que se pueden realizar más experimentos y hacer un nuevo entrenamiento sin los features menos significativos. Sin embargo, el modelo aún se comporta de muy buena manera.



Hallazgos Clave

Resultados muestran:

Según los gráficos de evaluación:

- Los índices de coeficientes de variables muestran que efectivamente, características como **"Reading Score"**, **"Writing Score"**, **"Gender Male"**, **"Race/Ethnicity Group E"** y **"Lunch Standard"** influyen de manera positiva en el resultado de la calificación de matemáticas.
- Sin embargo, resalta que la característica **"Test Preparation None"** influye aunque de poca manera en el resultado del examen. Y que, por otro lado, las características de **"Gender Female"** al igual que el análisis previo, parece influir de manera un poco negativa sobre el resultado del examen; Aunque curiosamente la característica **"Test Preparation Completed"** también influye negativamente.





¡Muchas Gracias!

¿Alguna consulta?



Enlace a Github

<https://github.com/RonaldoVindas/Danta-CasoAnalisis>