

```
1 from google.colab import files
2 uploaded = files.upload()
3
4 import pandas as pd
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7 df = pd.read_csv("student_habits_performance.csv")
8 df.head()
9 print(df)
10
11 sns.set(style="whitegrid")
12 #I use the code we work yesterday
13 # Crear múltiples subplots
14 fig, axs = plt.subplots(5, 2, figsize=(18, 24))
15 axs = axs.flatten()
16
17 # 1. Distribución del puntaje
18 sns.histplot(df['exam_score'], kde=True, ax=axs[0], color="skyblue")
19 axs[0].set_title("Distribución del Puntaje en el Examen")
20
21 # 2. Horas de estudio vs puntaje
22 sns.scatterplot(data=df, x='study_hours_per_day', y='exam_score', ax=axs[1])
23 axs[1].set_title("Horas de Estudio vs Puntaje en el Examen")
24
25 # 3. Salud mental vs puntaje
26 sns.boxplot(data=df, x='mental_health_rating', y='exam_score', ax=axs[2])
27 axs[2].set_title("Salud Mental vs Puntaje en el Examen")
28
29 # 4. Educación de los padres
30 sns.boxplot(data=df, x='parental_education_level', y='exam_score', ax=axs[3])
31 axs[3].set_title("Nivel Educativo de los Padres vs Puntaje")
32
33 # 5. Calidad del internet
34 sns.boxplot(data=df, x='internet_quality', y='exam_score', ax=axs[4])
35 axs[4].set_title("Calidad del Internet vs Puntaje")
36
37 # 6. Dieta y puntaje
38 sns.boxplot(data=df, x='diet_quality', y='exam_score', ax=axs[5])
39 axs[5].set_title("Calidad de la Dieta vs Puntaje")
40
41 # 7. Participación extracurricular
42 sns.boxplot(data=df, x='extracurricular_participation', y='exam_score', ax=axs[6])
43 axs[6].set_title("Actividades Extracurriculares vs Puntaje")
44
45 # 8. Redes sociales vs puntaje
46 sns.scatterplot(data=df, x='social_media_hours', y='exam_score', ax=axs[7])
47 axs[7].set_title("Horas en Redes Sociales vs Puntaje")
48
49 # 9. Trabajo medio tiempo
50 sns.boxplot(data=df, x='part_time_job', y='exam_score', ax=axs[8])
51 axs[8].set_title("Trabajo de Medio Tiempo vs Puntaje")
52
53 # 10. Sueño vs puntaje
54 sns.regplot(data=df, x='sleep_hours', y='exam_score', ax=axs[9], scatter_kws={'alpha':0.5})
55 axs[9].set_title("Horas de Sueño vs Puntaje")
56
57 plt.tight_layout()
58 plt.show()
59
60 # Heatmap de correlación entre variables numéricas
61 plt.figure(figsize=(12, 8))
62 corr = df.corr(numeric_only=True)
63 sns.heatmap(corr, annot=True, cmap="coolwarm", fmt=".2f")
64 plt.title("Heatmap de Correlaciones entre Variables Numéricas")
65 plt.show()
66
```

```
Elegir archivos student_hab...rmance.csv
• student_habits_performance.csv(text/csv) - 73663 bytes, last modified: 6/5/2025 - 100% done
Saving student_habits_performance.csv to student_habits_performance (2).csv

student_id age gender study_hours_per_day social_media_hours \
0 S1000 23 Female 0.0 1.2
1 S1001 20 Female 6.9 2.8
2 S1002 21 Male 1.4 3.1
3 S1003 23 Female 1.0 3.9
4 S1004 19 Female 5.0 4.4
.. ... ..
995 S1995 21 Female 2.6 0.5
996 S1996 17 Female 2.9 1.0
997 S1997 20 Male 3.0 2.6
998 S1998 24 Male 5.4 4.1
999 S1999 19 Female 4.3 2.9

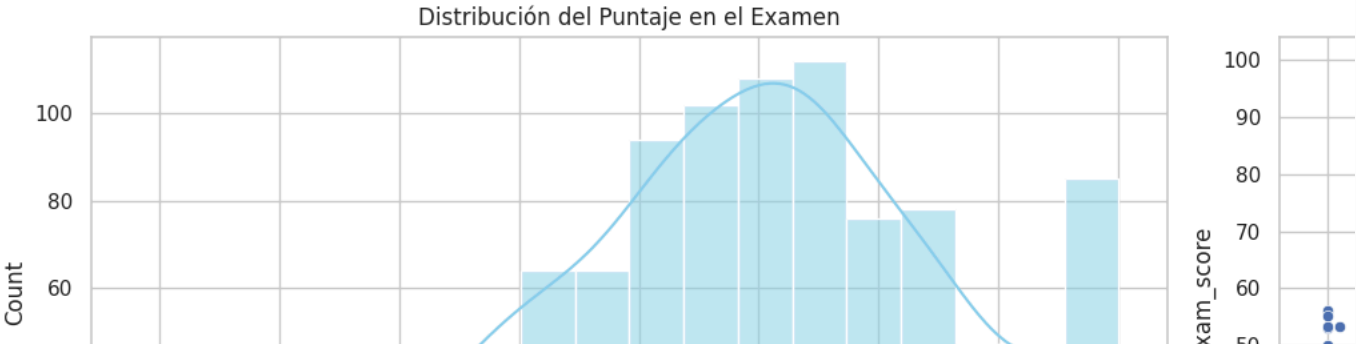
netflix_hours part_time_job attendance_percentage sleep_hours \
0 1.1 No 85.0 8.0
1 2.3 No 97.3 4.6
2 1.3 No 94.8 8.0
3 1.0 No 71.0 9.2
4 0.5 No 90.9 4.9
.. ... ..
995 1.6 No 77.0 7.5
996 2.4 Yes 86.0 6.8
997 1.3 No 61.9 6.5
998 1.1 Yes 100.0 7.6
999 1.9 No 89.4 7.1

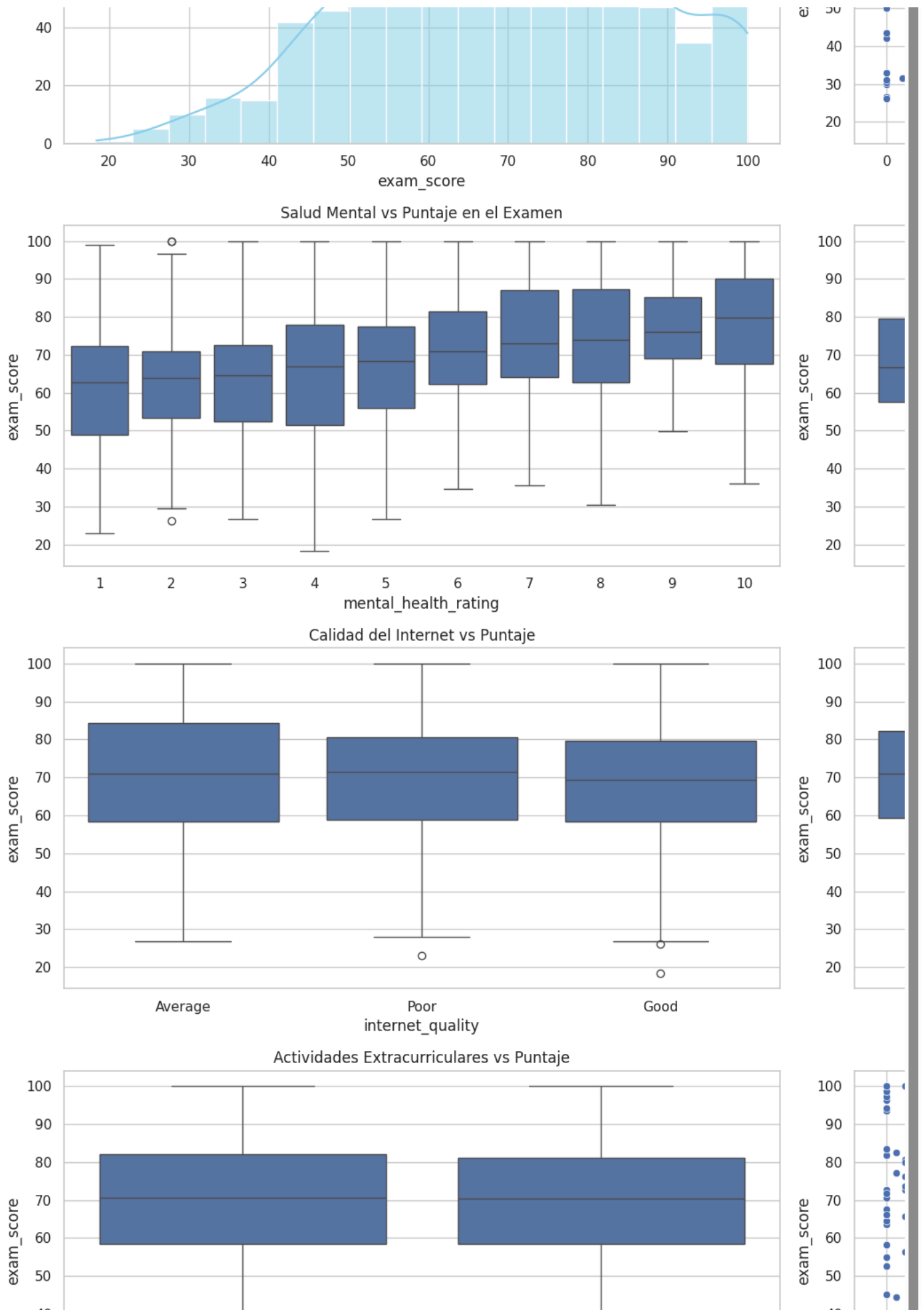
diet_quality exercise_frequency parental_education_level \
0 Fair 6 Master
1 Good 6 High School
2 Poor 1 High School
3 Poor 4 Master
4 Fair 3 Master
.. ... ..
995 Fair 2 High School
996 Poor 1 High School
997 Good 5 Bachelor
998 Fair 0 Bachelor
999 Good 2 Bachelor

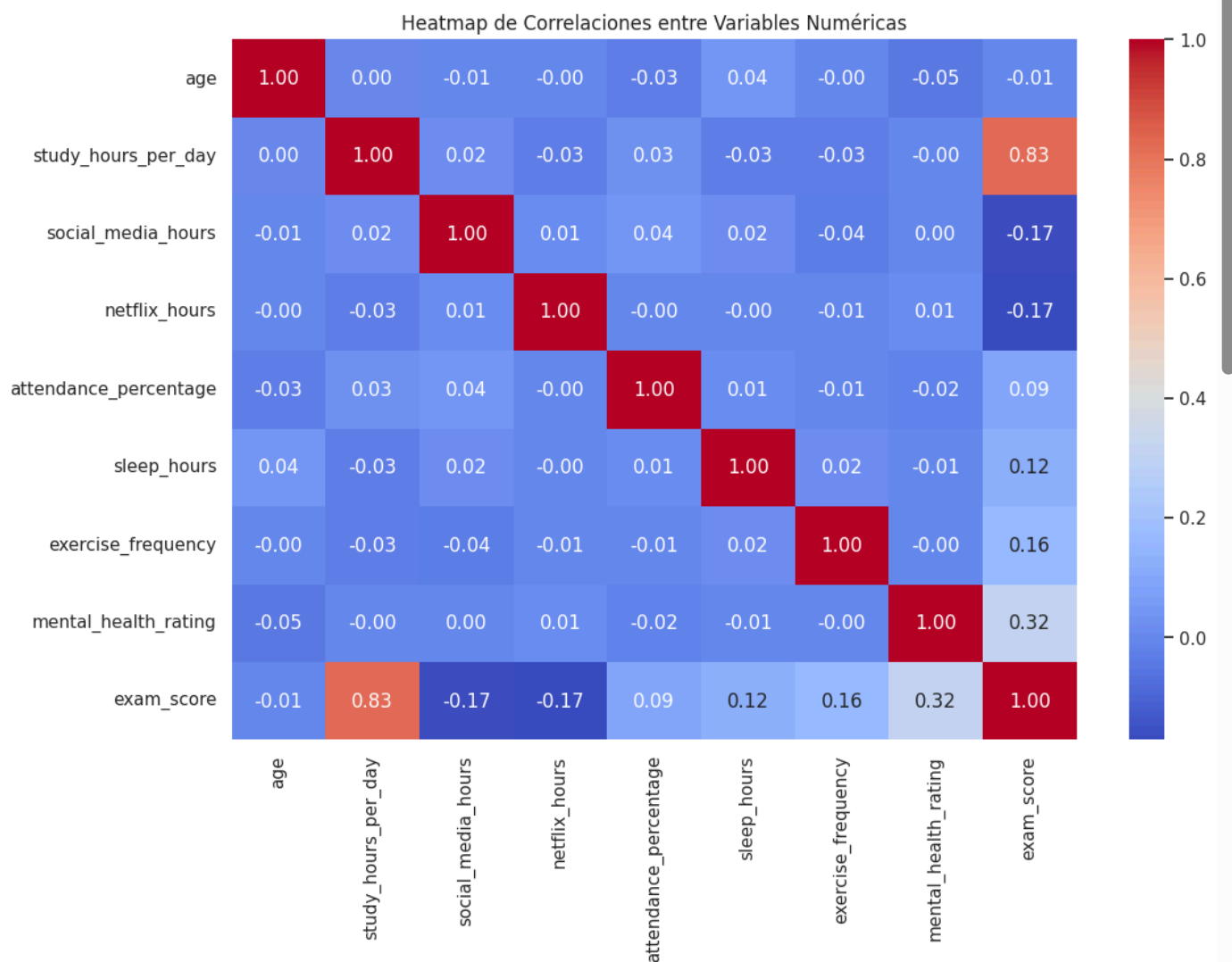
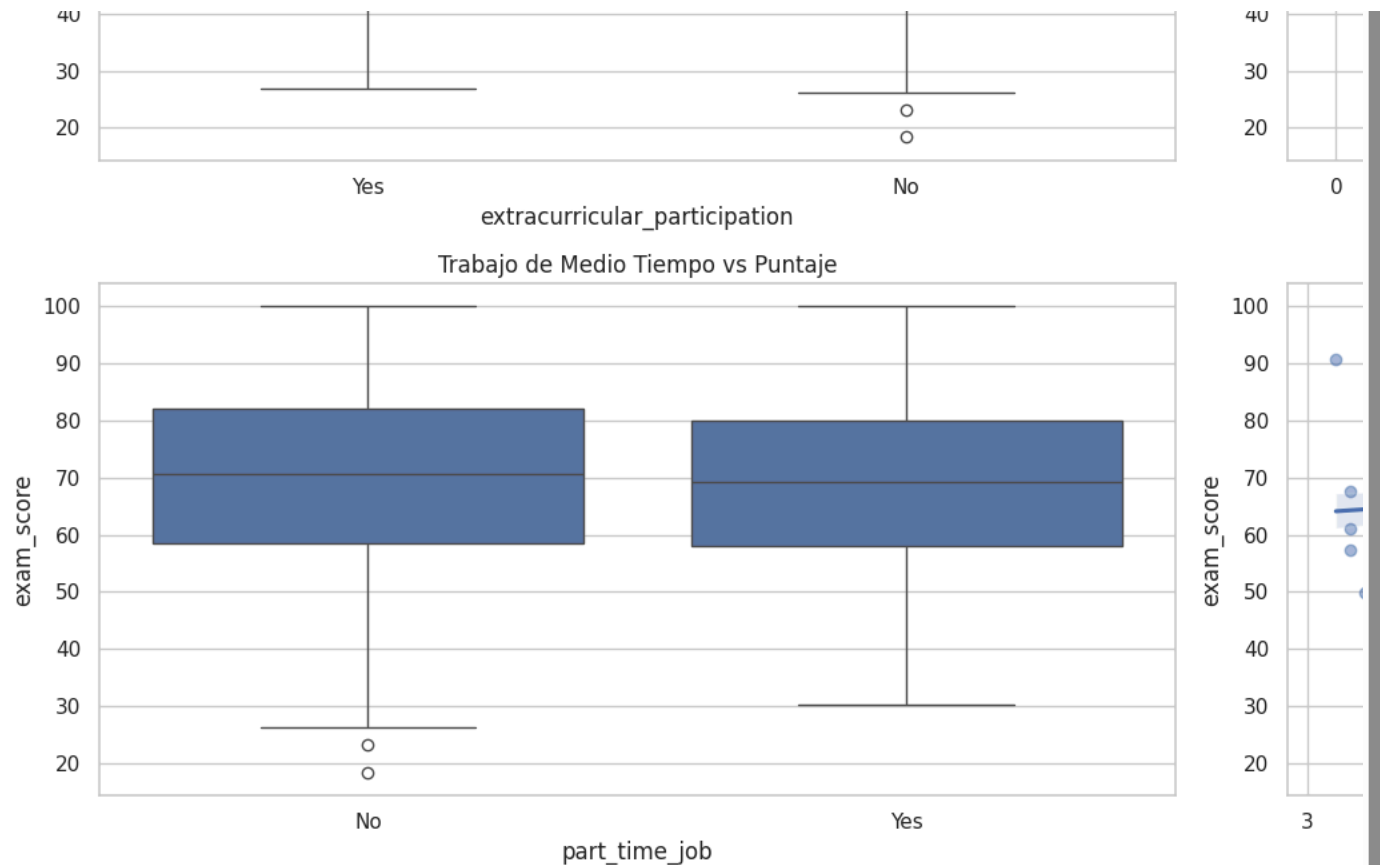
internet_quality mental_health_rating extracurricular_participation \
0 Average 8 Yes
1 Average 8 No
2 Poor 1 No
3 Good 1 Yes
4 Good 1 No
.. ... ..
995 Good 6 Yes
996 Average 6 Yes
997 Good 9 Yes
998 Average 1 No
999 Average 8 No

exam_score
0 56.2
1 100.0
2 34.3
3 26.8
4 66.4
.. ...
995 76.1
996 65.9
997 64.4
998 69.7
999 74.9
```

[1000 rows x 16 columns]









Are there any variables that do not provide information?

Variables like `student_id` are identifiers and do not provide useful information for analysis

If you had to eliminate variables, which ones would you remove and why?

I would remove `student_id` because it's an identifier

Are there any variables with unusual data?

Some outliers may be seen in the `exam_score` histogram or in the scatterplot of `study_hours_per_day`

If you compare the variables, are they all in similar ranges?

No, variables like `exam_score` (0–100) and `social_media_hours` or `sleep_hours` have different scales.

Do you think this affects the data analysis? Can you find any similar groups? What are these groups?

Yes, normalization or standardization is needed when comparing across different scales.

The scatterplots and boxplots show some trends, like students with more study hours or better mental health generally scoring higher. These may form clusters, such as:

- High performers with good habits
- Low performers with poor health or habits
- Average students with mixed patterns

<https://github.com/SypremeLemus/Act>