

A00573055 Rodrigo Martinez Vallejo Link del repositorio: <https://github.com/a00573055/Mastering-Analytics->

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

✓ Data Preparation

```
data = pd.read_csv('data.csv')

numeric_data = data.select_dtypes(include='number')
```

✓ Data Description

In this section I review the type of data each columns holds, as long as what it represents and the upper and lower limits of those columns.

```
print(data.dtypes)
```

```

student_id      object
age             int64
gender          object
study_hours_per_day float64
social_media_hours float64
netflix_hours   float64
part_time_job   object
attendance_percentage float64
sleep_hours     float64
diet_quality     object
exercise_frequency int64
parental_education_level object
internet_quality object
mental_health_rating int64
extracurricular_participation object
exam_score      float64
dtype: object
```

In total there are 15 variables and a total number of 1000 rows

- **student_id** *object* it is the unique identifier of each student
- **age** *int64* Age of student. Goes from 17 to 24
- **gender** *object* Male/Female/Other.
- **study_hours_per_day** *float64* Avg. daily study time. Goes from 0 to 8.3
- **social_media_hours** *float* Daily social media time. Goes from 0 to 7.2
- **netflix_hours** *float64* Avg. daily Netflix/binging time. Goes from 0 to 5.4
- **part_time-job** Yes/No.
- **attendance_percentage** *float64* Class attendace (0-100%).
- **sleep_hours** *float64* Avg. daily sleep

- **diet_quality** *object* Poor/Fair/Good. Goes from 3.2 to 10
- **exercise_frequency** *int64* Times per week. Goes from 0 to 7.
- **parental_education_level** *object* HighSchool/Bachelor/Other
- **internet_quality** *object* Good/Average/Other
- **mental_health_rating** *int64* Scale of 1 to 10
- **extracurricular_participation** *object* Yes/No
- **exam_scores** *float64* Final exam score (0-100)

Haz doble clic (o pulsa Intro) para editar

✓ Mean, median, and standard deviation

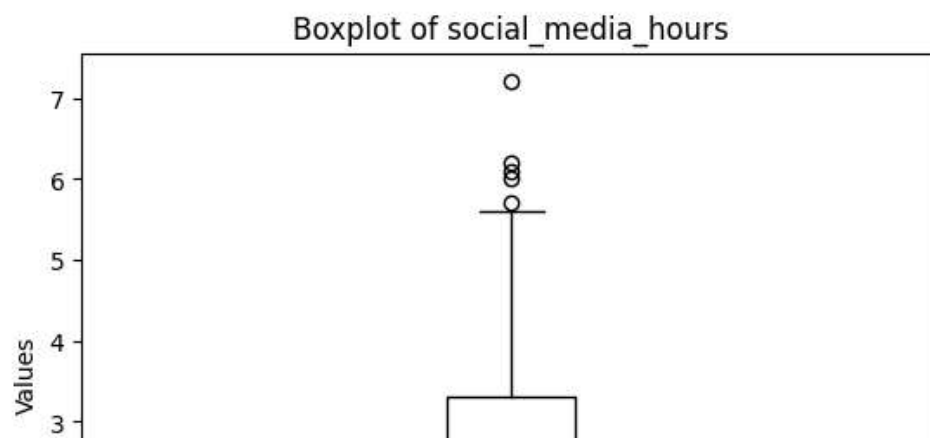
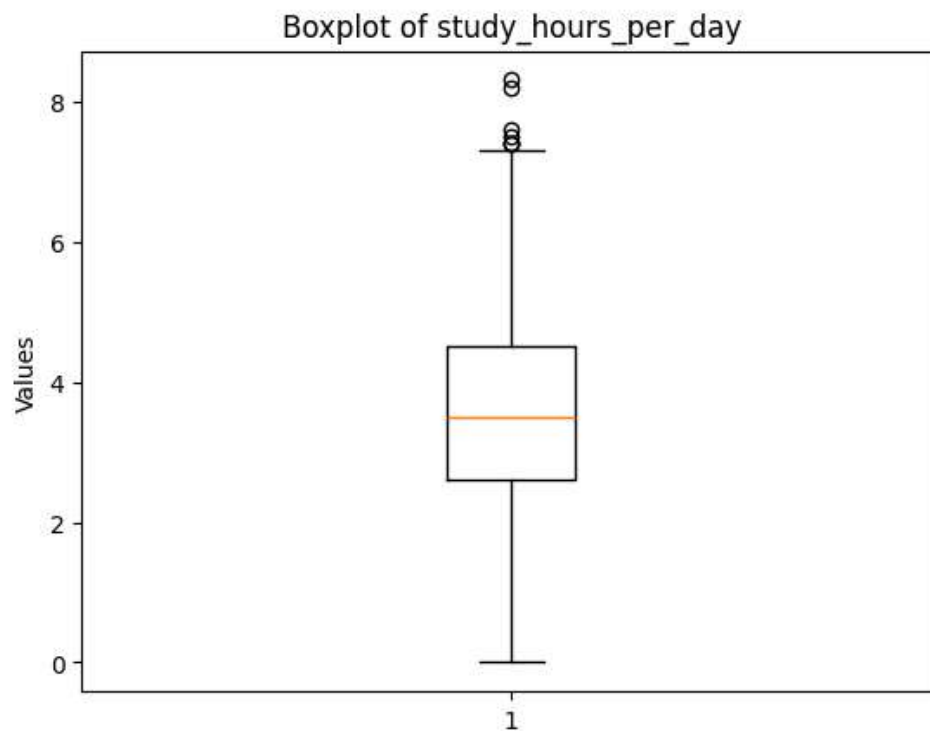
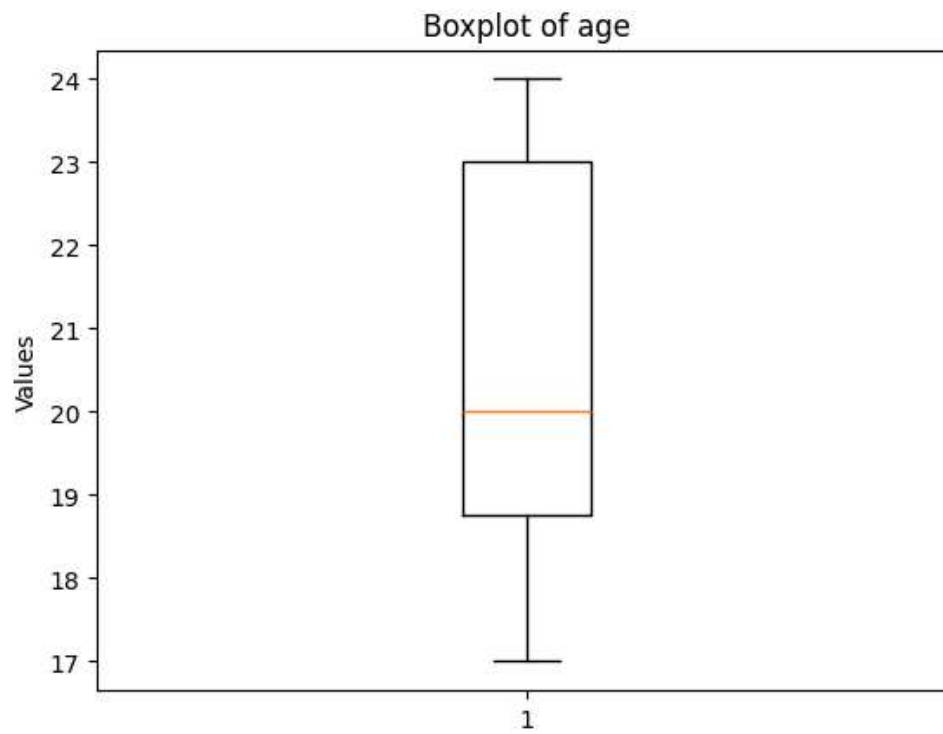
The first thing that I notice from this data is for example how the results given from Mean and Median are almost exactly the same which suggests a very symmetrical distribution of the data.

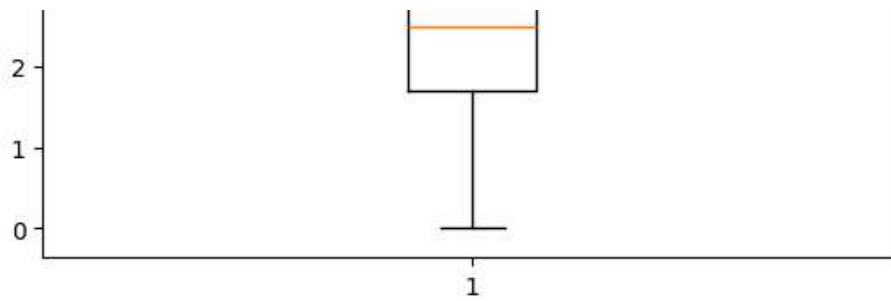
```
print("-->Mean")
print(numeric_data.mean())
print("-->Median")
print(numeric_data.median())
print("-->Standard deviation")
print(numeric_data.std())
```

```
→ -->Mean
age                20.4980
study_hours_per_day  3.5501
social_media_hours   2.5055
netflix_hours        1.8197
attendance_percentage 84.1317
sleep_hours          6.4701
exercise_frequency   3.0420
mental_health_rating  5.4380
exam_score           69.6015
dtype: float64
-->Median
age                20.0
study_hours_per_day  3.5
social_media_hours   2.5
netflix_hours        1.8
attendance_percentage 84.4
sleep_hours          6.5
exercise_frequency   3.0
mental_health_rating  5.0
exam_score           70.5
dtype: float64
-->Standard deviation
age                2.308100
study_hours_per_day  1.468890
social_media_hours   1.172422
netflix_hours        1.075118
attendance_percentage 9.399246
sleep_hours          1.226377
exercise_frequency   2.025423
mental_health_rating  2.847501
exam_score           16.888564
dtype: float64
```

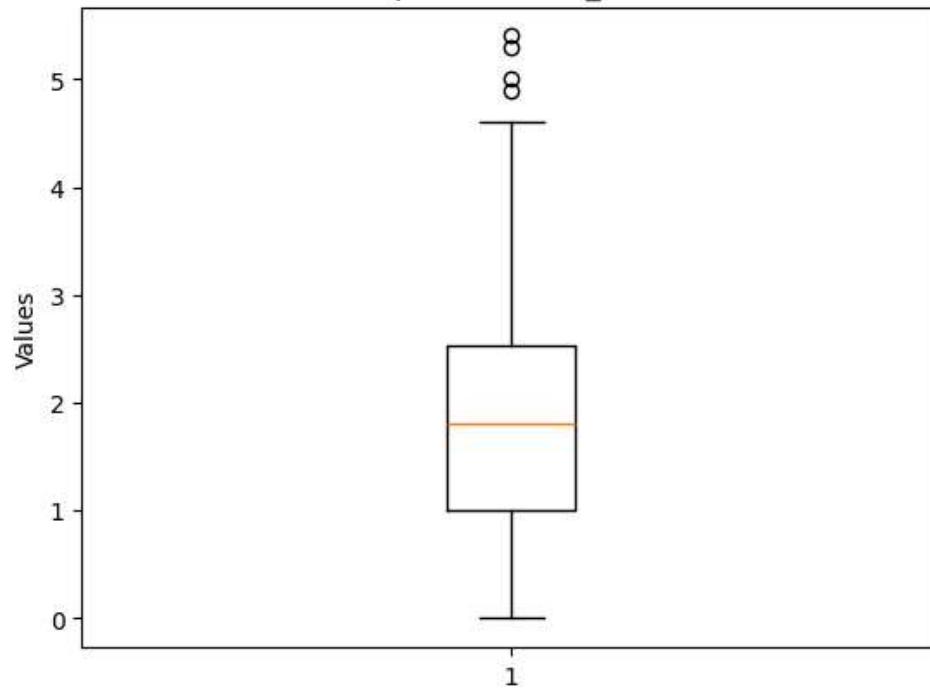
✓ Box Diagram

```
col_names = ['age', 'study_hours_per_day', 'social_media_hours', 'netflix_hours',  
             'attendance_percentage', 'sleep_hours', 'exercise_frequency',  
             'mental_health_rating', 'exam_score']  
  
for name in col_names:  
    plt.boxplot(numeric_data[name])  
    plt.title(f'Boxplot of {name}')  
    plt.ylabel('Values')  
    plt.show()
```

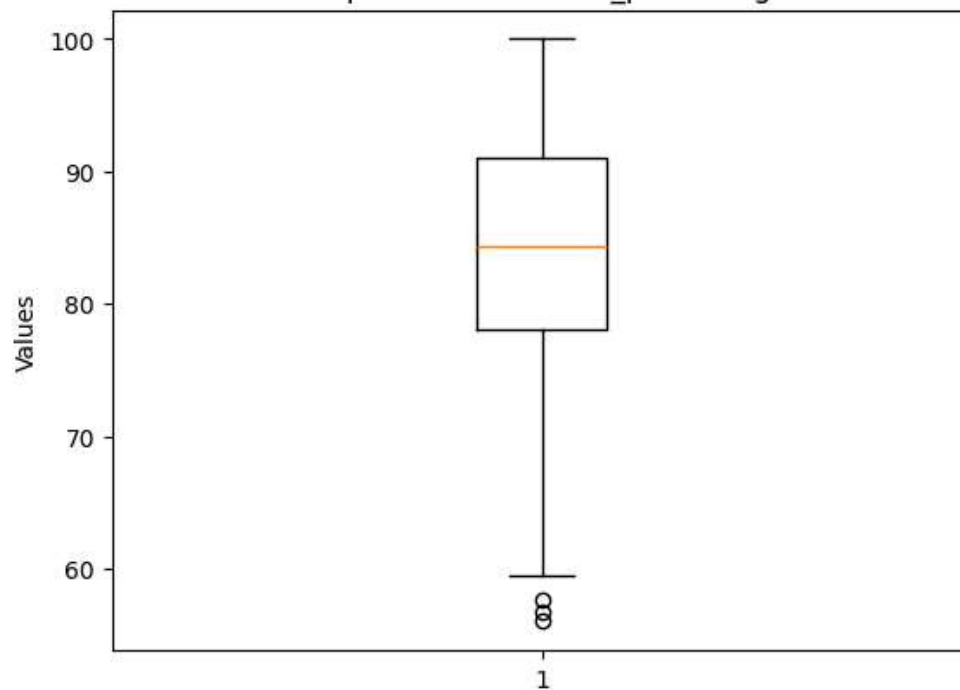




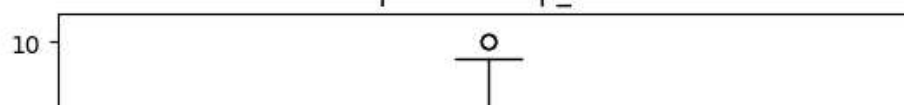
Boxplot of netflix_hours

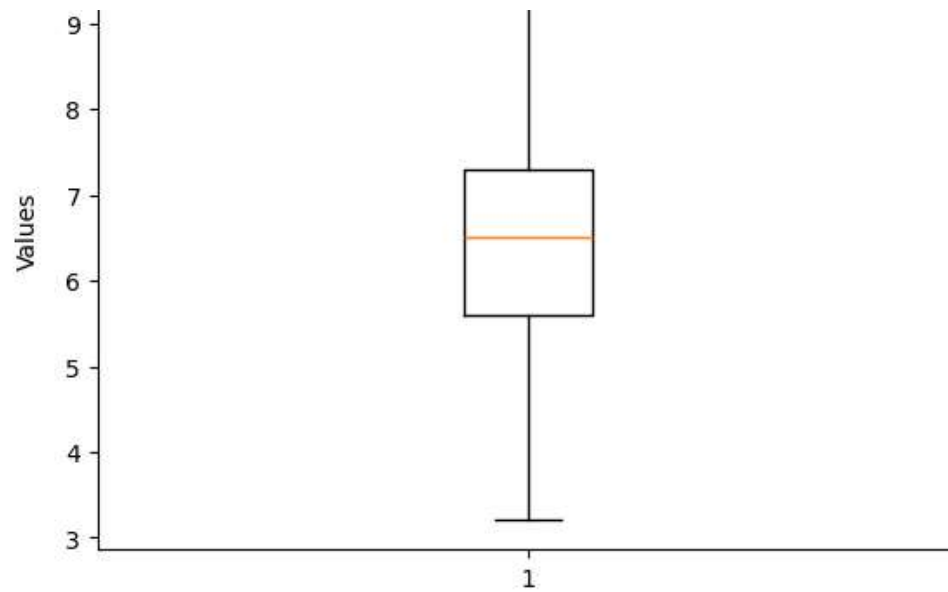


Boxplot of attendance_percentage

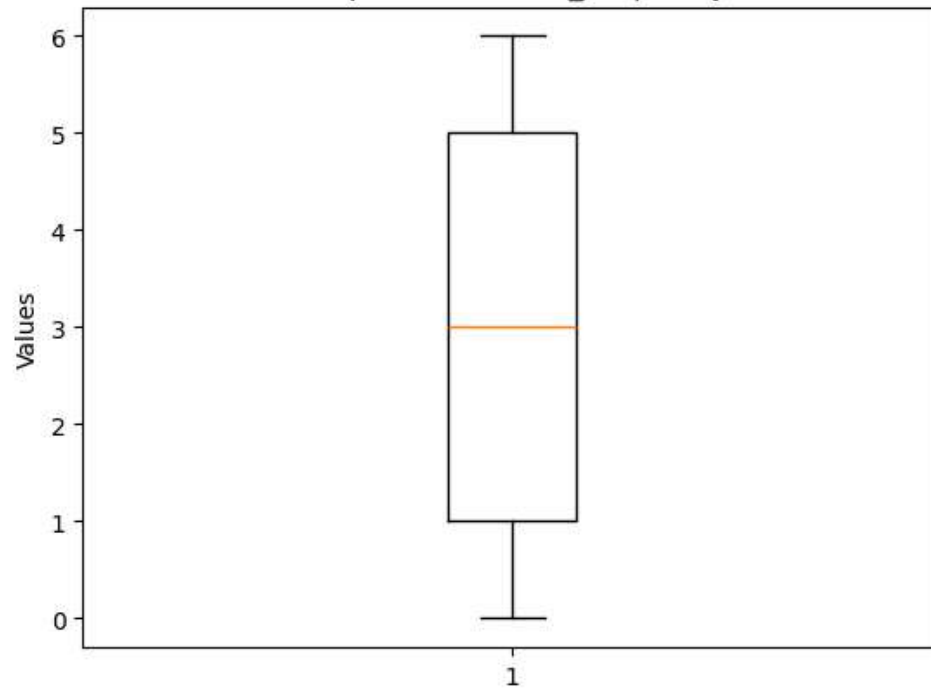


Boxplot of sleep_hours

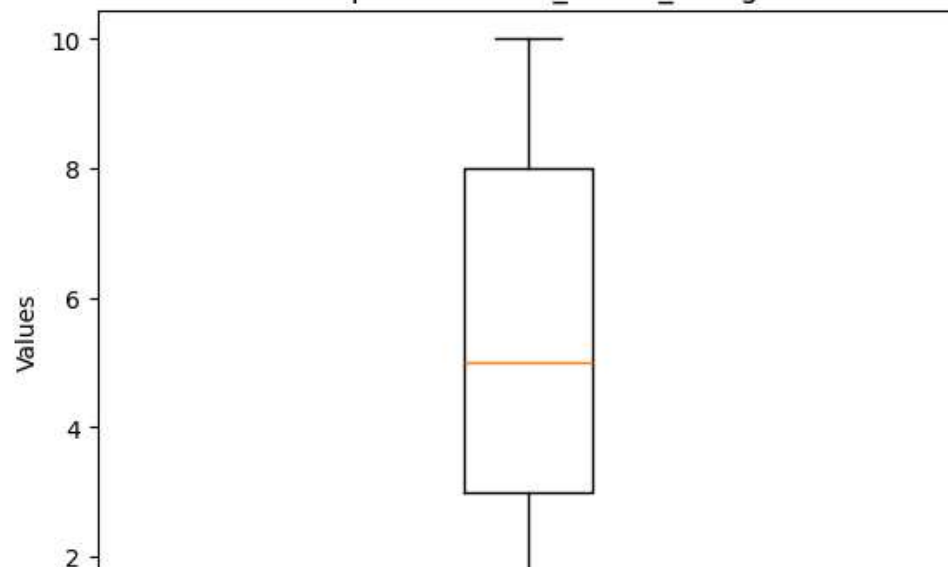


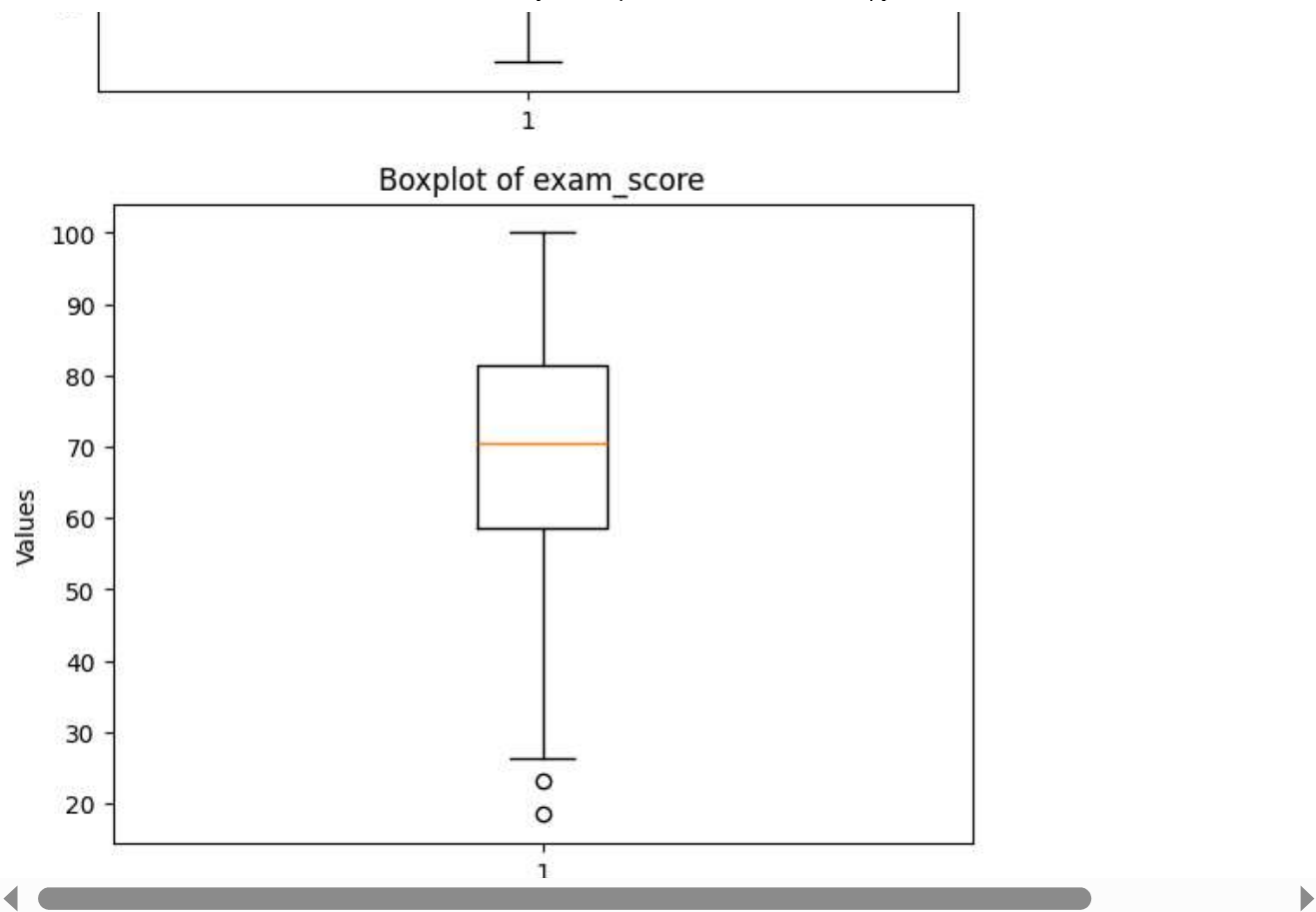


Boxplot of exercise_frequency



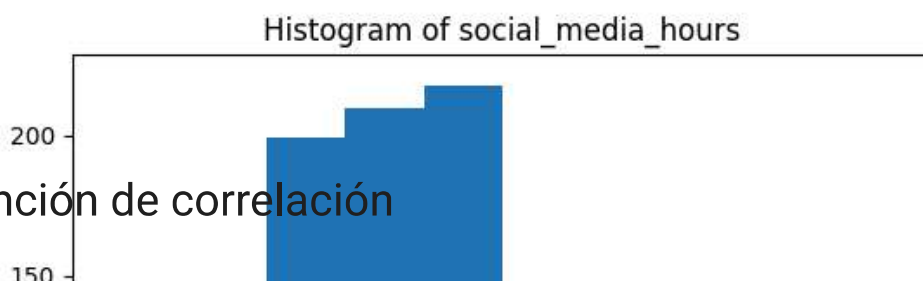
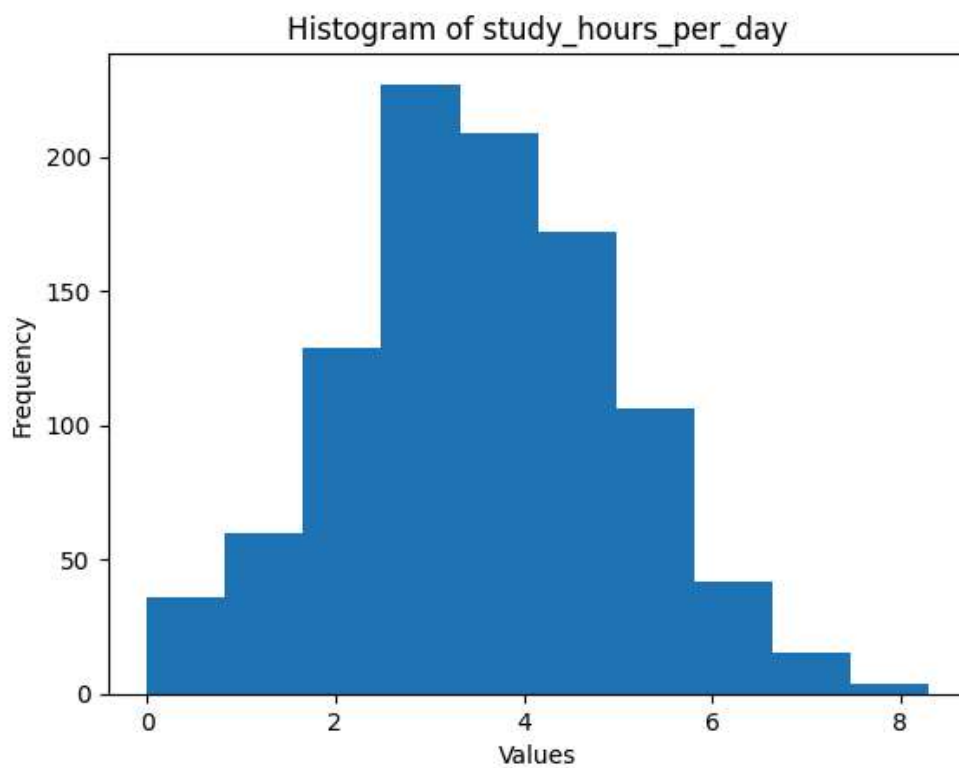
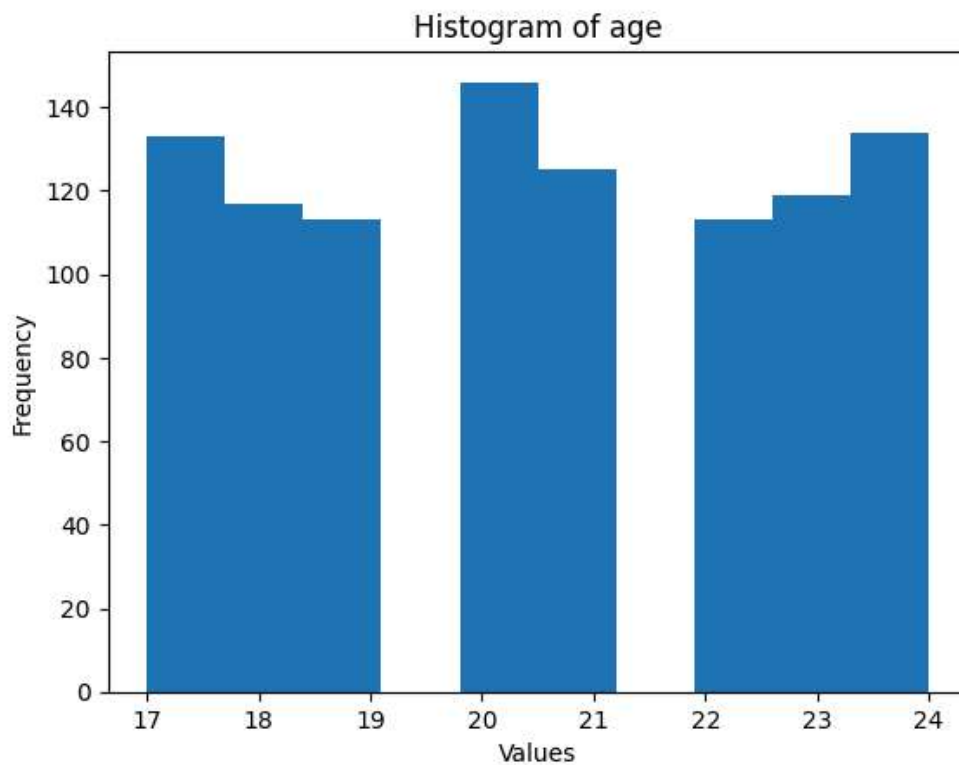
Boxplot of mental_health_rating





✓ .hits() function

```
for col in col_names:
    plt.hist(numeric_data[col])
    plt.title(f'Histogram of {col}')
    plt.xlabel('Values')
    plt.ylabel('Frequency')
    plt.show()
```

✓ Función de correlación

```
num_data_corr = numeric_data.corr()
print(num_data_corr)
```

