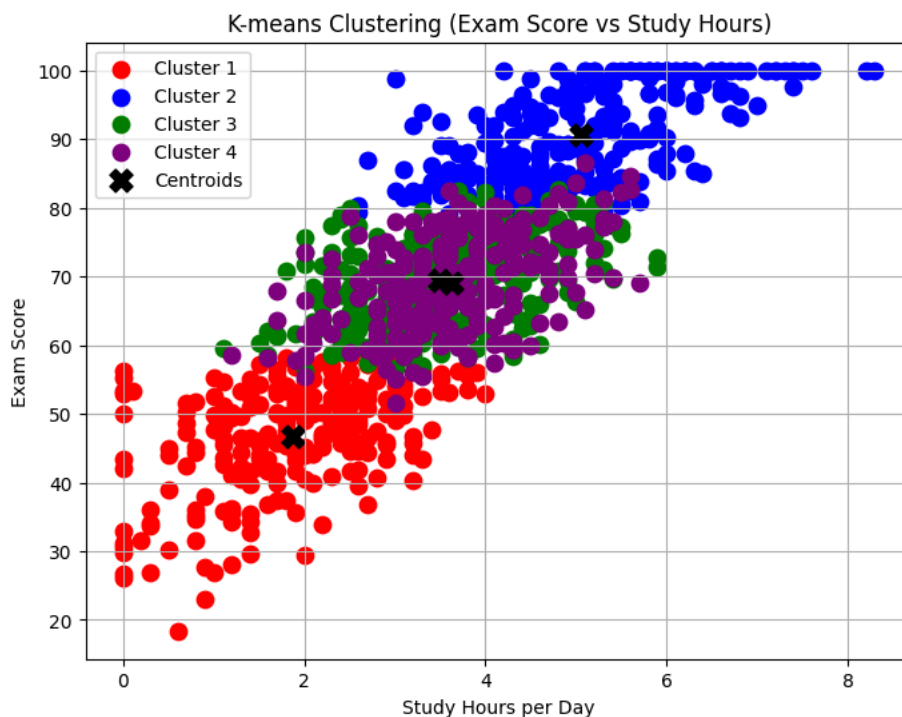```
 1 from google.colab import files
 2 uploaded = files.upload()
 3
 4 import pandas as pd
 5 import numpy as np
 6 import matplotlib.pyplot as plt
 7 from sklearn.cluster import KMeans
 8
 9
10 df = pd.read_csv("student_habits_performance.csv")
11
12 # - 'student_id' is an identifier
13 # - Categorical variables are not used here because they require special encoding (outside current scope)
14
15 numeric_data = df.select_dtypes(include=[np.number])
16
17 data_points = numeric_data.to_numpy()
18
19 K = 4 #The value of k was set to 4 based on visual inspection of the scatter plot and the assumption that students can be grouped into fo
20
21 # use scikit-learn to calculate the centroids
22 kmeans = KMeans(n_clusters=K, random_state=0, n_init=10)
23 labels = kmeans.fit_predict(data_points)
24 centroids = kmeans.cluster_centers_
25
26 plt.figure(figsize=(8, 6))
27 colors = ['red', 'blue', 'green', 'purple']
28 for k in range(K):
29     cluster = data_points[labels == k]
30     plt.scatter(cluster[:, 1], cluster[:, -1], color=colors[k], label=f'Cluster {k+1}', s=80)
31
32 plt.scatter(centroids[:, 1], centroids[:, -1], color='black', marker='X', s=150, label='Centroids')
33 plt.xlabel('Study Hours per Day')
34 plt.ylabel('Exam Score')
35 plt.title('K-means Clustering (Exam Score vs Study Hours)')
36 plt.legend()
37 plt.grid(True)
38 plt.show()
39
```

⤓  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 (7).csv



K-means Clustering (Exam Score vs Study Hours)

1. Do you think these centers could be representative of the data? Why?

Indeed, they seem quite representative. Each of them is located around the center of its respective group of data points, showing that the algorithm successfully grouped students with similar characteristics and test scores.

2. How do you get the value of k for use? The value of k = 4 was likely based on the decisions made;

- Visual inspection of data distribution.
- Prior knowledge of how many natural clusters can be expected.
- Alternatively, the elbow method is used to define where adding more clusters no longer significantly improves the fit.

3. Would the centers use more if they used more? Less value?

- Most of the value can be obtained in smaller, more connected groups, which can overwhelm the data and reduce interpretation.
- A lower value of k can confuse different samples into larger groups, lacking meaningful distinctions.

So there is a trade-off: higher means more granular, which means more general. The corporeality depends on the actual data structure.

4. How far apart are the centers? Are any close to others?

From the plots;

- Some centroids are well separated, as in the red and blue clusters.
- Others, especially green and purple, seem to be closer together, which indicates the behavior of overlapping clusters with insufficient differentiation.

5. What would happen to the centers if we had many outliers in the box-and-whisker analysis?

Outliers would pull the centroids toward extreme values, potentially misplacing the center of the true cluster.

6. What can you say about the data based on the centers?

There are distinct groups of students based on their study hours and exam performance. Students who study more tend to score higher, as shown by the upward trend in the plot. Some students may study less but still perform well, possibly due to other factors (like natural aptitude).

https://github.com/SvpremeLemus/Act