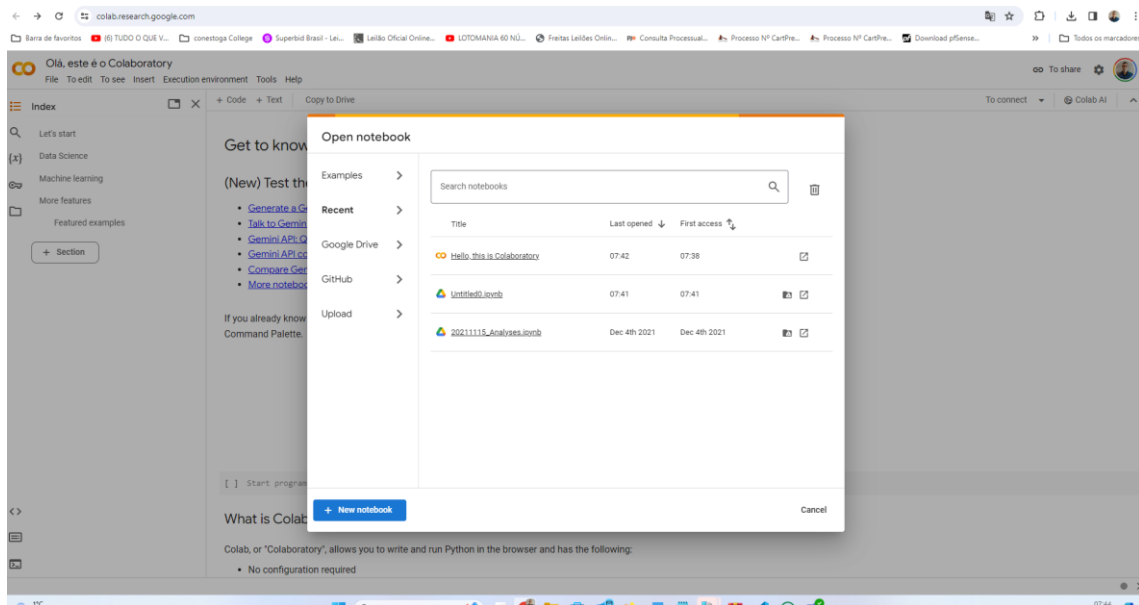IRIS



- Review the details of the IRIS Data set

The IRIS dataset is a classic machine learning and statistics dataset introduced by British biologist and statistician Ronald Fisher in 1936. It contains information about 150 iris flowers from three different species: setosa, versicolor, and virginica. Each iris sample has four features: sepal length, sepal width, petal length, and petal width, all measured in centimetres. The primary purpose of the dataset is to demonstrate various classification techniques. It is commonly used for training machine learning models to classify iris flowers based on their features. The dataset is well-structured, with each species having 50 samples, making it suitable for training and testing algorithms. Due to its simplicity and clarity, it serves as a starting point for many beginners in the field of machine learning. Fisher's work with this dataset pioneered multivariate analysis in biology. It's often utilized in tutorials, textbooks, and introductory courses to illustrate concepts like feature selection, model evaluation, and data visualization. The dataset is freely available and widely used in research and educational contexts across various disciplines, including statistics, machine learning, and botany.

+ Código    + Texto    Todas as alterações foram salvas

```python
[7] import pandas as pd
    import numpy as np
    import matplotlib.pyplot as plt
    from sklearn import datasets
```

```python
[9] iris = datasets.load_iris()
```

```python
[11] df = pd.DataFrame({
        'x': iris.data[:,0],
        'y': iris.data[:,1],
        'cluster' : iris.target
    })
```

```python
[12] df
```

|     | x   | y   | cluster |
|-----|-----|-----|---------|
| 0   | 5.1 | 3.5 | 0       |
| 1   | 4.9 | 3.0 | 0       |
| 2   | 4.7 | 3.2 | 0       |
| 3   | 4.6 | 3.1 | 0       |
| 4   | 5.0 | 3.6 | 0       |
| ... | ... | ... | ...     |
| 145 | 6.7 | 3.0 | 2       |
| 146 | 6.3 | 2.5 | 2       |
| 147 | 6.5 | 3.0 | 2       |
| 148 | 6.2 | 3.4 | 2       |
| 149 | 5.9 | 3.0 | 2       |

150 rows × 3 columns

Next steps:    Generate code with df    View recommended plots

CO Untitled0.ipynb - Colaboratory    × +
← → C 🔒 colab.research.google.com/drive/14mjsHccpxqxbs3Hi9CU1PDgbPrCqwUQk#scrollTo=mw_AP_8RFKr4&uniqifier=1
Barra de favoritos | (6) TUDO O QUE V... | conestoga College | Superbid Brasil - Lei... | Leilão Oficial Online... | LOTOMANIA 60 NÚ... | Freitas Leilões Onlin... | Consulta Processual... | Processo Nº CartPre... | Processo Nº CartPre... | Download pfSense... | » | Todos os marcadores
+ Código  + Texto   Todas as alterações foram salvas    ✓   Colab AI

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn import datasets
```

```python
[58] iris = datasets.load_iris()
```

```python
[59] df = pd.DataFrame({
        "x": iris.data[:,0],
        "y": iris.data[:,1],
        "cluster" : iris.target
    })
```

```python
[60] df
```

|     | x   | y   | cluster |
|-----|-----|-----|---------|
| 0   | 5.1 | 3.5 | 0       |
| 1   | 4.9 | 3.0 | 0       |
| 2   | 4.7 | 3.2 | 0       |
| 3   | 4.6 | 3.1 | 0       |
| 4   | 5.0 | 3.6 | 0       |
| ... | ... | ... | ...     |
| 145 | 6.7 | 3.0 | 2       |
| 146 | 6.3 | 2.5 | 2       |
| 147 | 6.5 | 3.0 | 2       |
| 148 | 6.2 | 3.4 | 2       |
| 149 | 5.9 | 3.0 | 2       |

150 rows × 3 columns

Next steps:  [Generate code with df]  [View recommended plots]

```python
[61] centroids = {}
for i in range(3):
    result_list = []
    result_list.append(df.loc[df['cluster'] == i ]['x'].mean())
    result_list.append(df.loc[df['cluster'] == i ]['y'].mean())
    centroids[i] = result_list
```

```python
[73] centroids
```

```
{0: [5.006, 3.428],
 1: [5.936, 2.7700000000000005],
 2: [6.587999999999998, 2.974]}
```

CO Untitled0.ipynb - Colaboratory    × +
← → C 🔒 colab.research.google.com/drive/14mjsHccpxqxbs3Hi9CU1PDgbPrCqwUQk#scrollTo=mw_AP_8RFKr4&uniqifier=1
Barra de favoritos | (6) TUDO O QUE V... | conestoga College | Superbid Brasil - Lei... | Leilão Oficial Online... | LOTOMANIA 60 NÚ... | Freitas Leilões Onlin... | Consulta Processual... | Processo Nº CartPre... | Processo Nº CartPre... | Download pfSense... | » | Todos os marcadores
+ Código  + Texto   Todas as alterações foram salvas    ✓   Colab AI

```
[73] {0: [5.006, 3.428],
 1: [5.936, 2.7700000000000005],
 2: [6.587999999999998, 2.974]}
```

```python
fig = plt.figure(figsize=(5, 5))
plt.scatter(df["x"], df["y"], c=iris.target)
plt.xlabel('Speal Length', fontsize=18)
plt.ylabel('Sepal Width', fontsize=18)
```

```
Text(0, 0.5, 'Sepal Width')
```

CO Untitled0.ipynb - Colaboratory    × +
← → C 🔒 colab.research.google.com/drive/14mjsHccpxqxbs3Hi9CU1PDgbPrCqwUQk#scrollTo=mw_AP_8RFKr4&uniqifier=1
Barra de favoritos | (6) TUDO O QUE V... | conestoga College | Superbid Brasil - Lei... | Leilão Oficial Online... | LOTOMANIA 60 NÚ... | Freitas Leilões Onlin... | Consulta Processual... | Processo Nº CartPre... | Processo Nº CartPre... | Download pfSense... | » | Todos os marcadores
+ Código  + Texto   Todas as alterações foram salvas    ✓   Colab AI

```
[73] 149  5.9  3.0    2
```
150 rows × 3 columns

Next steps:  [Generate code with df]  [View recommended plots]

```python
centroids = {}
for i in range(3):
    result_list = []
    result_list.append(df.loc[df['cluster'] == i ]['x'].mean())
    result_list.append(df.loc[df['cluster'] == i ]['y'].mean())
    centroids[i] = result_list
```

```python
[73] centroids
```

```
{0: [5.006, 3.428],
 1: [5.936, 2.7700000000000005],
 2: [6.587999999999998, 2.974]}
```

```python
fig = plt.figure(figsize=(5, 5))
plt.scatter(df["x"], df["y"], c=iris.target,alpha =0.3)
colmap = {0: 'r', 1: 'g', 2: 'b'}
col = [0,1]
for i in centroids.keys():
    plt.scatter(centroids[i][0], centroids[i][1], c=colmap[i], edgecolor='k')
plt.show()
```
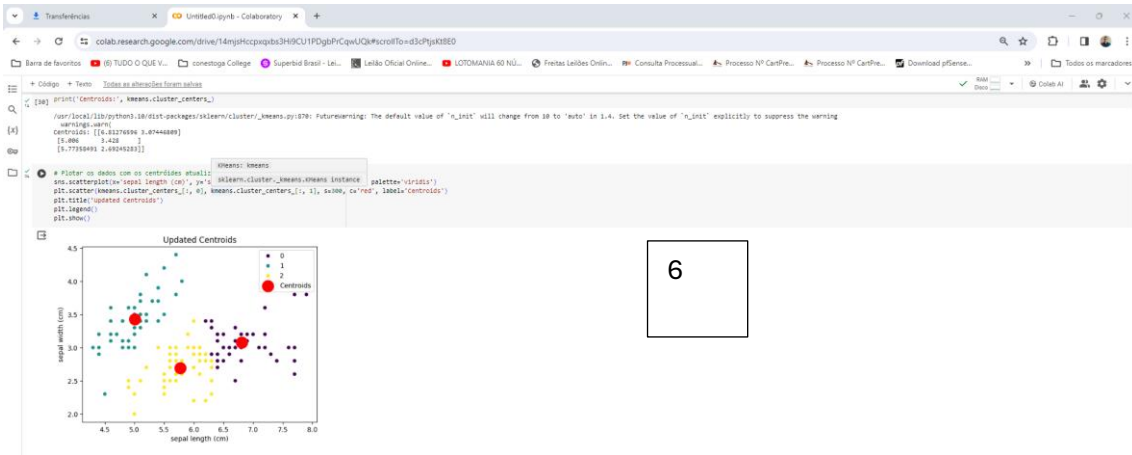
Second work.



1- Importing Libraries
2- Creating Data frame structure
3- Exploratory analysis of sepal length and sepal width for the 3 classes



5 Creating centroids

6- update cetroids



7- Final clustering result

In this assignment, I conducted an exploratory analysis of a data set that included measurements of the length and width of sepals from three different classes of flowers. After importing the libraries and enabling and structuring the data into a Data Frame, I explored the distribution of sepal measurements for each flower class.

Then, I applied the K-means clustering algorithm to group the data into clusters based on their characteristics. I created initial centroids and assigned the data to the closest clusters. Iteratively, I was able to update the centroids and recalculate the assignment of the data to the clusters until the centroids stabilized.

We visualize the updated centroids with the original data to better understand the distribution of the clusters. Finally, I presented the final clustering result graphically, highlighting the clusters identified by the K-means algorithm. However, I would like to see a graphic image in another way, but I could not.

Reference

https://colab.research.google.com

https://www.kaggle.com/datasets/uciml/irishttps://medium.com/@sarakarim/k-means-clustering-for-iris-dataset-in-google-colab-30a6d78556c4

https://scikit-learn.org/stable/auto_examples/datasets/plot_iris_dataset.html