

Exploratory Data Analysis with Python or R - Tell Me About a Dataset

Pamulapati Harika (001266981)

Saint Louis University

ORES-5160 Data Management

October 26 2023

Introduction

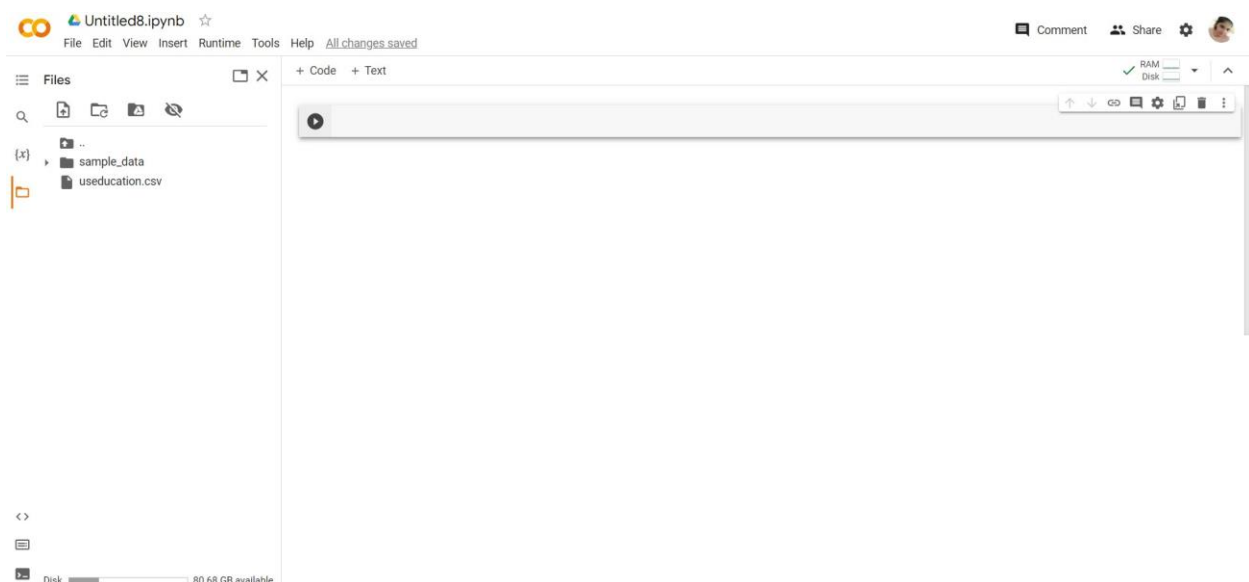
The steps involved in performing an Exploratory Data Analysis (EDA) on the “**US education**” dataset are outlined in this study. The dataset appears to be historical in nature and is organized to offer insights into how financing for education, enrollment, and other factors have changed over time in various U.S. states. In the field of education, this research is frequently utilized for policy analysis, decision-making, and research on education.

About the dataset

The U.S. Education Dataset is a thorough compilation of information that offers insights into many facets of the American educational system across a number of years. Usually, this dataset contains statistics on academic performance, student enrollment, school finance, and other important variables. This dataset is used by researchers, policymakers, educators, and analysts to track student performance, spot patterns and discrepancies in the U.S. education system, and better understand how educational resources are allocated. It is an invaluable tool for assessing the impact of policies and making well-informed choices to raise the standard of education in the nation.

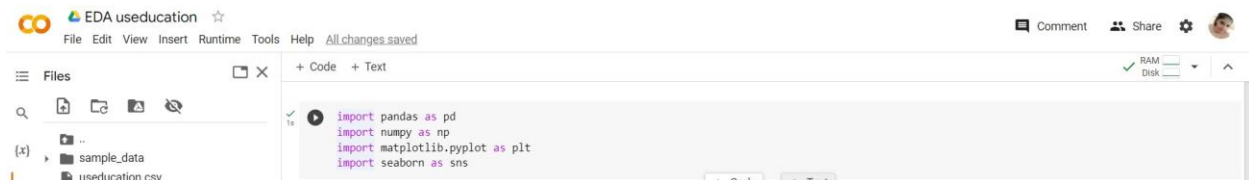
useducation dataset

I uploaded the useducation dataset into the Google colab that we utilised for the group project



2. Loading the dataset

Then I started by importing the libraries that are required for data manipulation and visualization, such as pandas, numpy, matplotlib, and seaborn, in order to start our investigation. These libraries enable us to create intelligent visualizations and operate efficiently with our dataset.

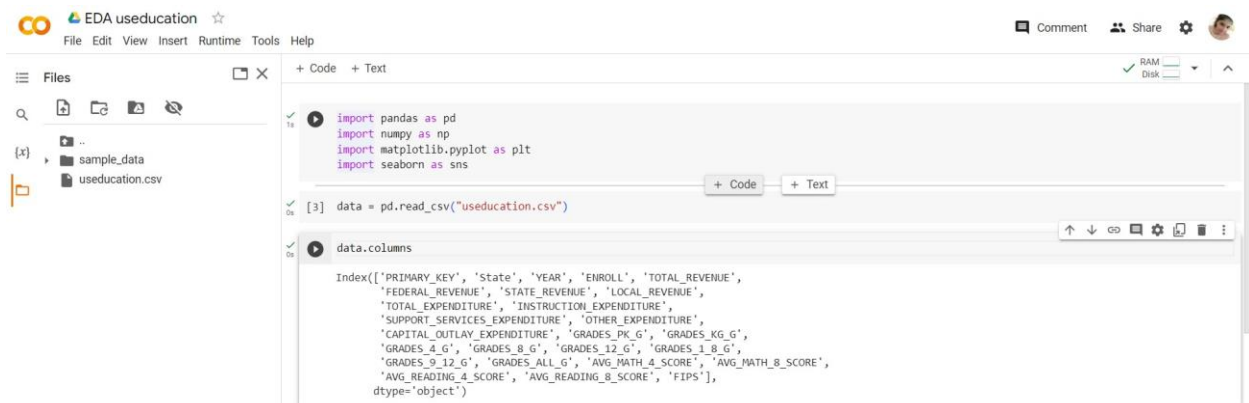


3. Once the file was read into Google Colab, I labeled it "data" and used the code below.

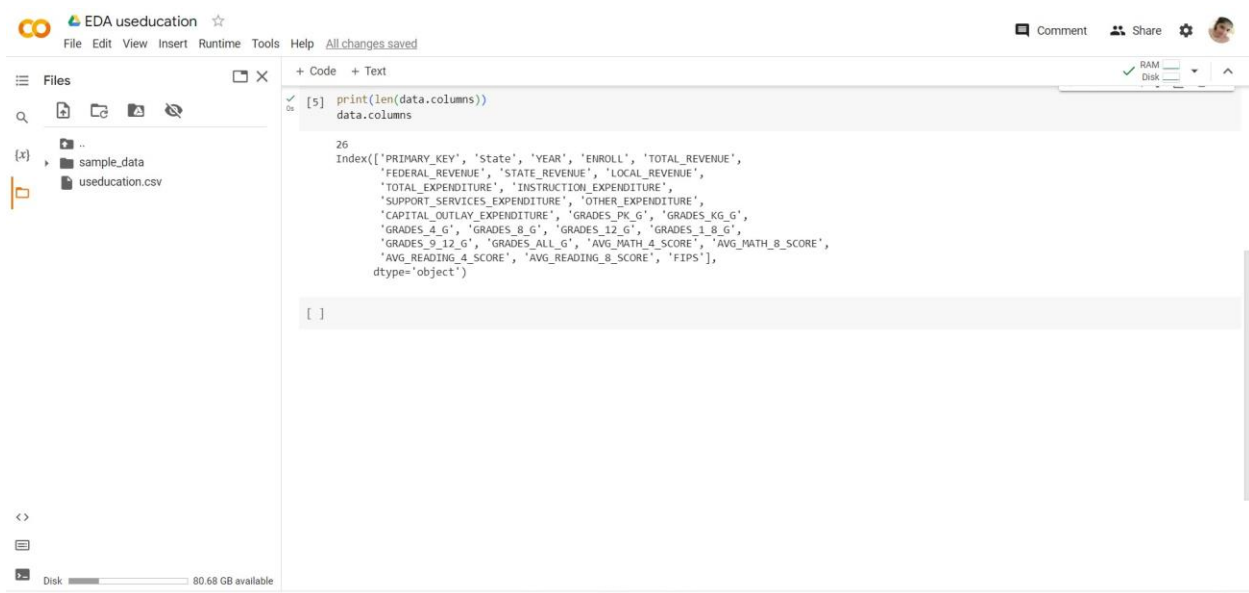
data = pd.read_csv("useducation.csv")



4. I used the code below to extract each column from the provided dataset.



5. The useducation dataset contains total **26** columns and the code is **data.columns** used to get all of them.



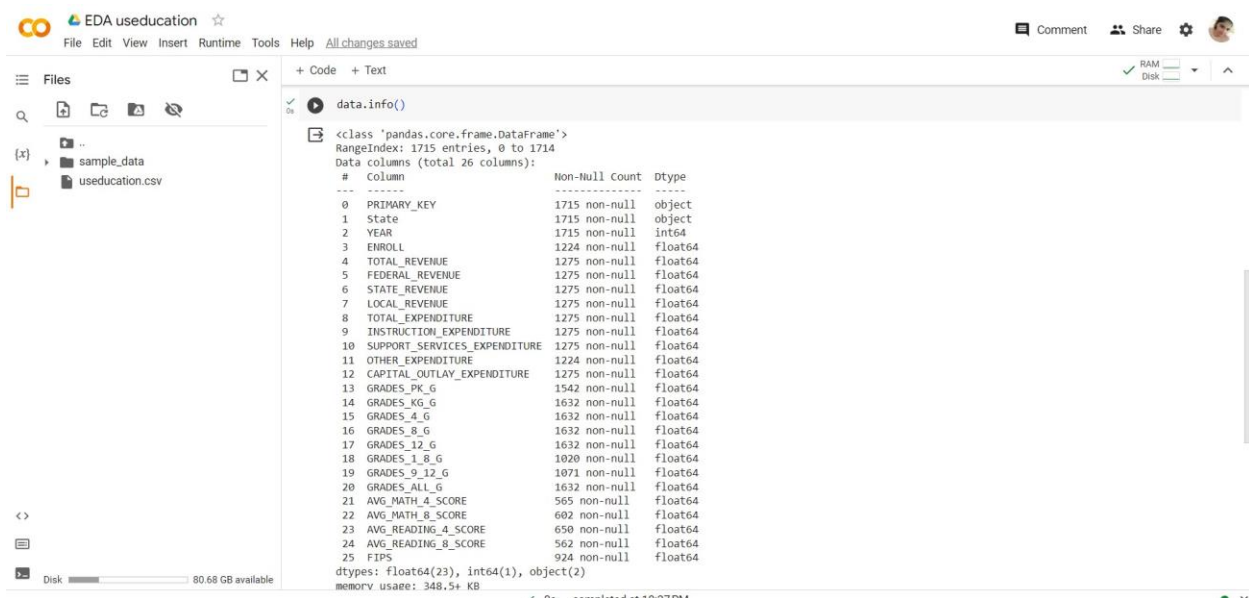
The screenshot shows the EDA useducation interface. The code editor contains the following code:

```
[5] print(len(data.columns))
data.columns
```

The output shows the length of the columns list is 26, followed by the list of column names:

```
Index(['PRIMARY_KEY', 'State', 'YEAR', 'ENROLL', 'TOTAL_REVENUE',
      'FEDERAL_REVENUE', 'STATE_REVENUE', 'LOCAL_REVENUE',
      'TOTAL_EXPENDITURE', 'INSTRUCTION_EXPENDITURE',
      'SUPPORT_SERVICES_EXPENDITURE', 'OTHER_EXPENDITURE',
      'CAPITAL_OUTLAY_EXPENDITURE', 'GRADES_PK_G', 'GRADES_KG_G',
      'GRADES_4_G', 'GRADES_8_G', 'GRADES_12_G', 'GRADES_1_8_G',
      'GRADES_9_12_G', 'GRADES_ALL_G', 'AVG_MATH_4_SCORE', 'AVG_MATH_8_SCORE',
      'AVG_READING_4_SCORE', 'AVG_READING_8_SCORE', 'FIPS'],
      dtype='object')
```

6. I was able to get the data information out of the file by using the code below. The **data.info()** domain is used to generate general information about the useducation data.



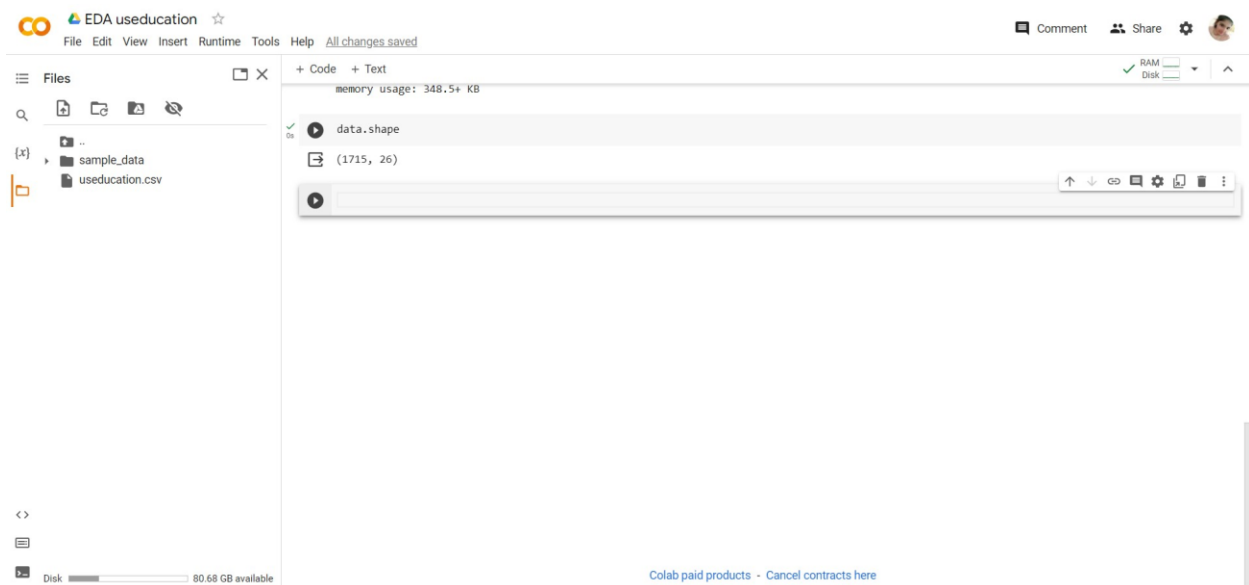
The screenshot shows the EDA useducation interface. The code editor contains the following code:

```
data.info()
```

The output shows the data information:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1715 entries, 0 to 1714
Data columns (total 26 columns):
#   Column                                Non-Null Count  Dtype
---  ---                                -
0   PRIMARY_KEY                          1715 non-null   object
1   State                                1715 non-null   object
2   YEAR                                 1715 non-null   int64
3   ENROLL                               1224 non-null   float64
4   TOTAL_REVENUE                        1275 non-null   float64
5   FEDERAL_REVENUE                      1275 non-null   float64
6   STATE_REVENUE                        1275 non-null   float64
7   LOCAL_REVENUE                        1275 non-null   float64
8   TOTAL_EXPENDITURE                    1275 non-null   float64
9   INSTRUCTION_EXPENDITURE              1275 non-null   float64
10  SUPPORT_SERVICES_EXPENDITURE          1275 non-null   float64
11  OTHER_EXPENDITURE                     1224 non-null   float64
12  CAPITAL_OUTLAY_EXPENDITURE            1275 non-null   float64
13  GRADES_PK_G                           1542 non-null   float64
14  GRADES_KG_G                           1632 non-null   float64
15  GRADES_4_G                            1632 non-null   float64
16  GRADES_8_G                            1632 non-null   float64
17  GRADES_12_G                           1632 non-null   float64
18  GRADES_1_8_G                          1020 non-null   float64
19  GRADES_9_12_G                         1071 non-null   float64
20  GRADES_ALL_G                          1632 non-null   float64
21  AVG_MATH_4_SCORE                      565 non-null   float64
22  AVG_MATH_8_SCORE                      602 non-null   float64
23  AVG_READING_4_SCORE                   650 non-null   float64
24  AVG_READING_8_SCORE                   562 non-null   float64
25  FIPS                                  924 non-null   float64
dtypes: float64(23), int64(1), object(2)
```

7. The total number of rows and columns in the file can be found using the following code. There are 1715 rows and 26 columns in the dataset.



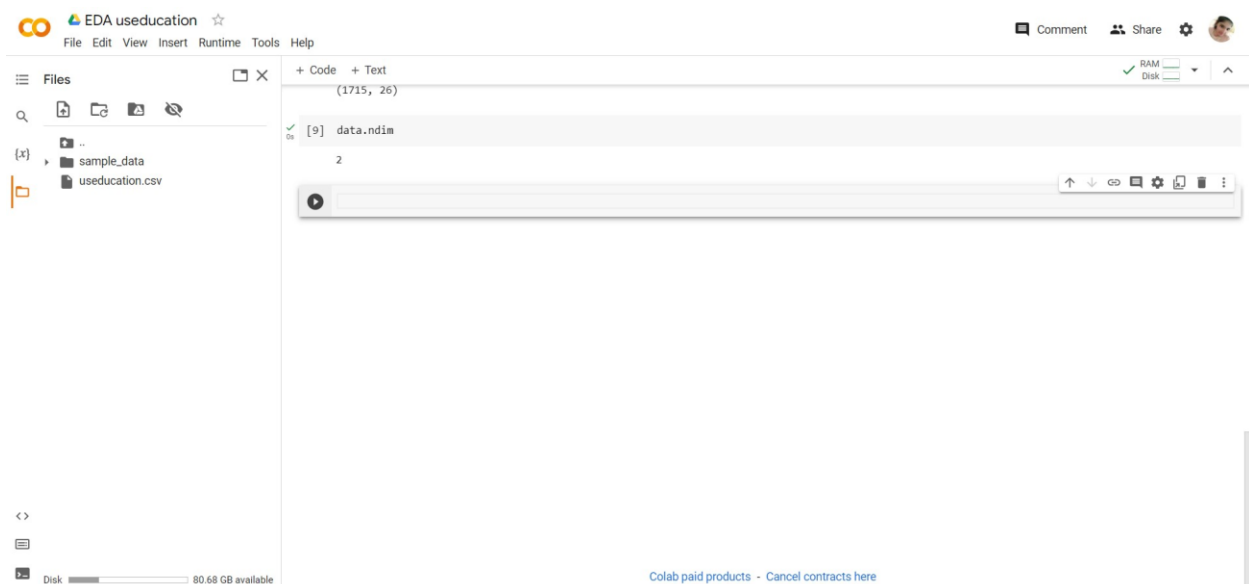
The screenshot shows the Google Colab interface for a notebook titled "EDA useducation". The left sidebar displays the file explorer with a folder named "sample_data" containing a file "useducation.csv". The main code editor area shows a code cell with the following content:

```
+ Code + Text
memory usage: 348.5+ KB

data.shape
(1715, 26)
```

The output of the code cell is displayed below the code, showing the shape of the data as a tuple (1715, 26). The interface includes a top menu bar with options like File, Edit, View, Insert, Runtime, Tools, and Help. The bottom status bar indicates "Disk" usage and "80.68 GB available".

8. To determine the number of dimensions or axes in the data structure. The dataset has two dimensions and **data.ndim** is used to determine the number of axes or dimensions in the data structure.



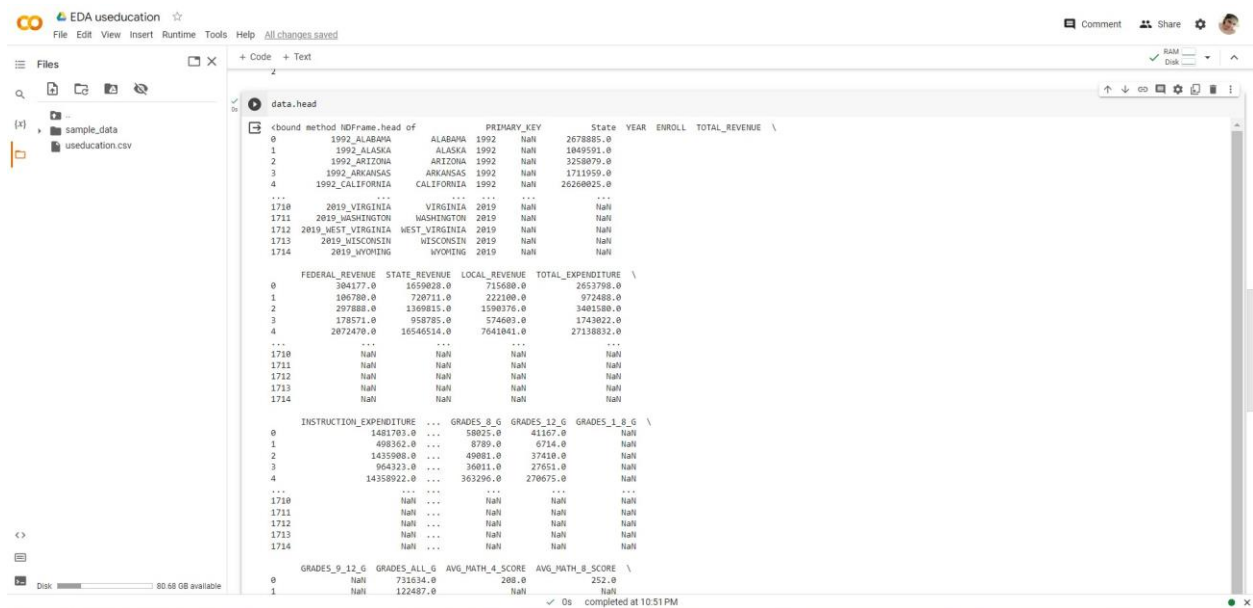
The screenshot shows the Google Colab interface for a notebook titled "EDA useducation". The left sidebar displays the file explorer with a folder named "sample_data" containing a file "useducation.csv". The main code editor area shows a code cell with the following content:

```
+ Code + Text
(1715, 26)

[9] data.ndim
2
```

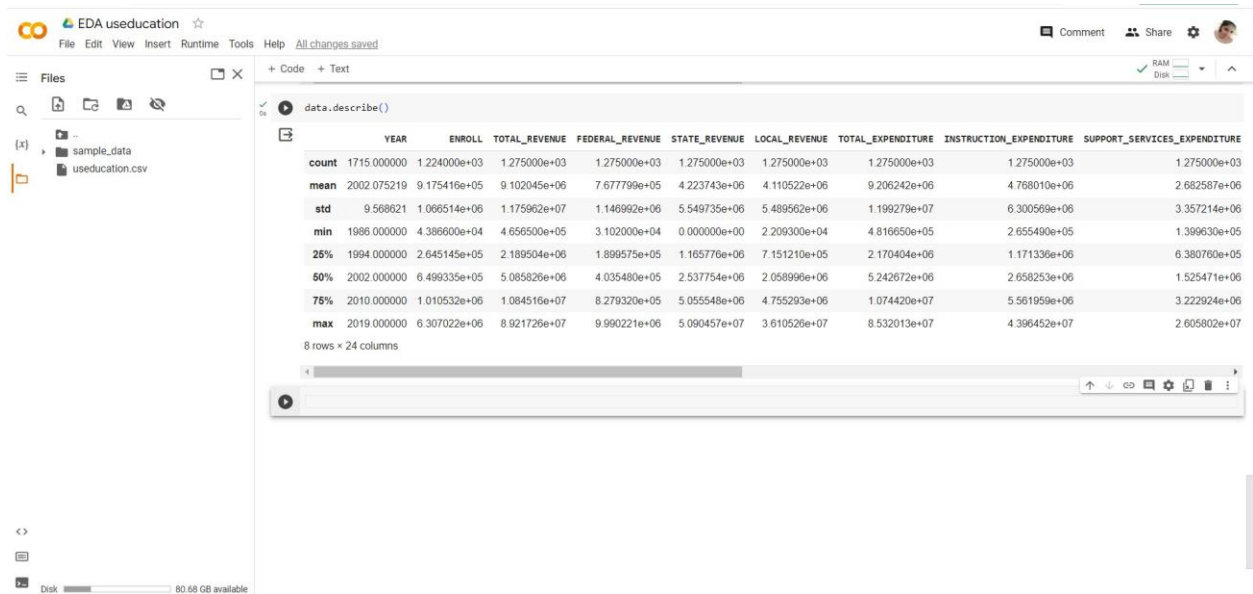
The output of the code cell is displayed below the code, showing the number of dimensions of the data as 2. The interface includes a top menu bar with options like File, Edit, View, Insert, Runtime, Tools, and Help. The bottom status bar indicates "Disk" usage and "80.68 GB available".

9. The `data.head` extension is used to display the real values in the rows and columns.



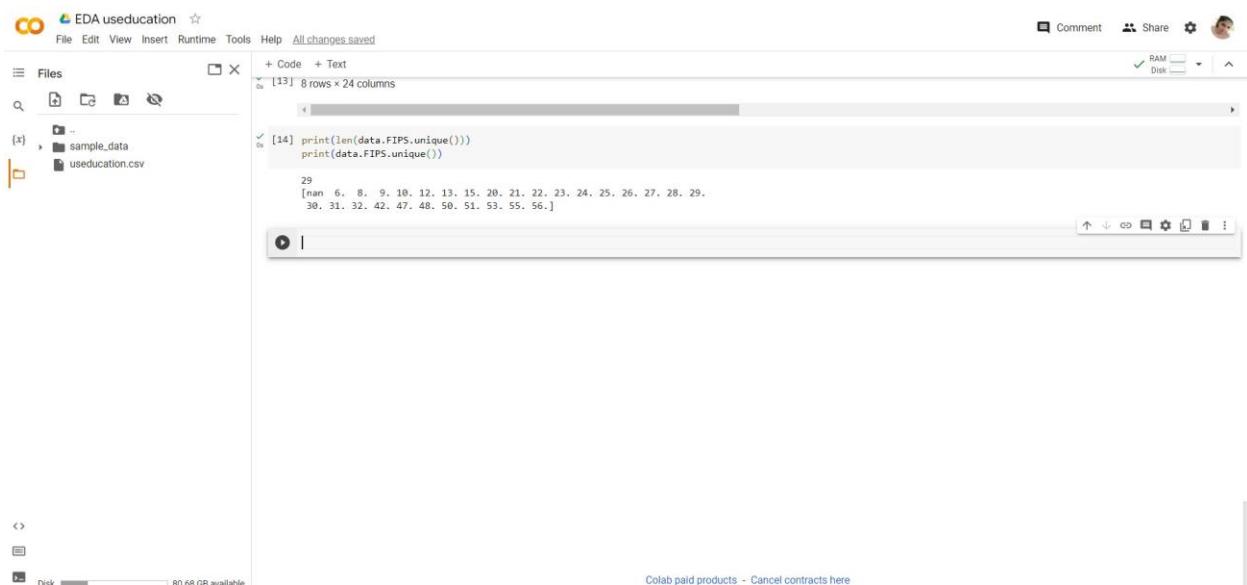
The screenshot shows the EDA useducation interface with the `data.head()` command executed. The output displays the first few rows of the dataset, showing columns for `PRIMARY_KEY`, `State`, `YEAR`, `ENROLL`, `TOTAL_REVENUE`, `FEDERAL_REVENUE`, `STATE_REVENUE`, `LOCAL_REVENUE`, `TOTAL_EXPENDITURE`, `INSTRUCTION_EXPENDITURE`, `SUPPORT_SERVICES_EXPENDITURE`, and `GRADES_12_G`. The data is organized into sections separated by ellipses, showing the first 5 rows, rows 1710-1714, and the last 5 rows of the dataset.

10. An summary of my dataset that is descriptive was produced using techniques such as `data.describe()`.



The screenshot shows the EDA useducation interface with the `data.describe()` command executed. The output displays a summary of the dataset, showing columns for `YEAR`, `ENROLL`, `TOTAL_REVENUE`, `FEDERAL_REVENUE`, `STATE_REVENUE`, `LOCAL_REVENUE`, `TOTAL_EXPENDITURE`, `INSTRUCTION_EXPENDITURE`, and `SUPPORT_SERVICES_EXPENDITURE`. The summary includes statistics such as `count`, `mean`, `std`, `min`, `25%`, `50%`, `75%`, and `max`. The data is organized into sections separated by ellipses, showing the first 5 rows, rows 1710-1714, and the last 5 rows of the dataset.

11. The code to find the total number of unique FIPS is provided below. A total of 29 FIPS with unique names are displayed.



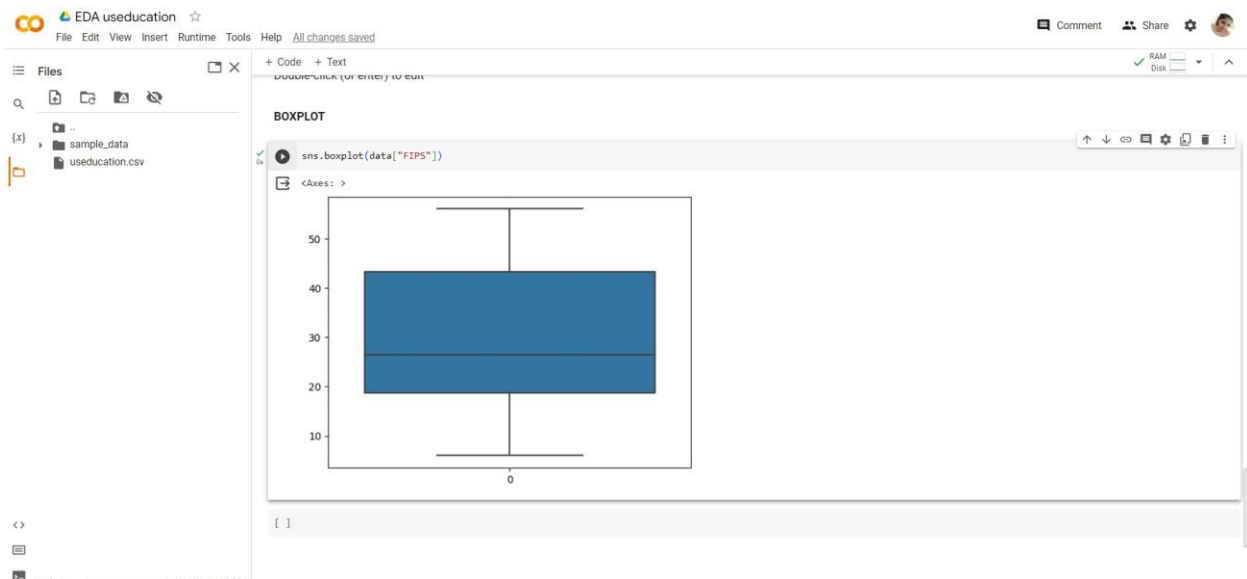
The screenshot shows a Google Colab interface with the following code and output:

```
[13] 8 rows x 24 columns
```

```
[14] print(len(data.FIPS.unique()))
      print(data.FIPS.unique())
```

```
29
[nan  6.  8.  9. 10. 12. 13. 15. 20. 21. 22. 23. 24. 25. 26. 27. 28. 29.
 30. 31. 32. 42. 47. 48. 50. 51. 53. 55. 56.]
```

12. To use Boxplot Python, you must import Matplotlib. Underneath the boxplot is the combination of the FIPS.



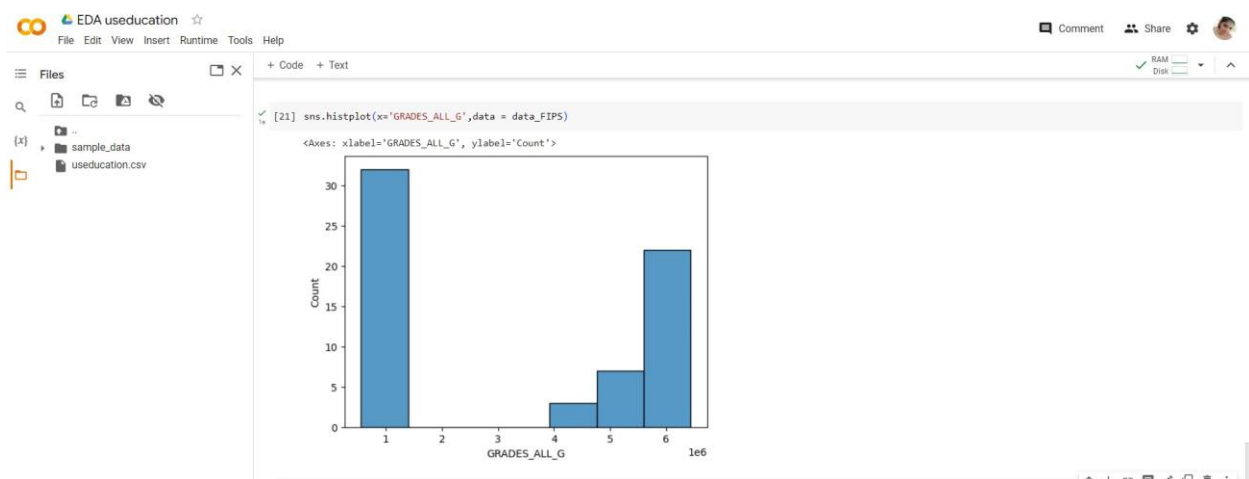
Interpretation:

The displayed boxplot provides insights into the distribution of the "FIPS" data. The central rectangle represents the interquartile range (IQR) with its lower and upper bounds indicating the 25th (Q1) and 75th (Q3) percentiles, respectively. The

horizontal line within the box signifies the median value, which offers a measure of central tendency. The whiskers, or the lines extending from the box, show the range within which the bulk of values fall, with outliers typically lying beyond these extremities. Based on the plot, the data appears to be roughly symmetric, given that the median is centrally located within the IQR. The absence of any points outside the whiskers suggests that there aren't noticeable outliers in the dataset. The overall range seems to span from a value slightly above 0 to about 50.

HISTOGRAM

In Python EDA, a histplot is a visual aid that displays data about a numerical variable's distribution, central trends, and spread by fusing a histogram and a kernel density plot. It's a great tool for figuring out patterns in data.



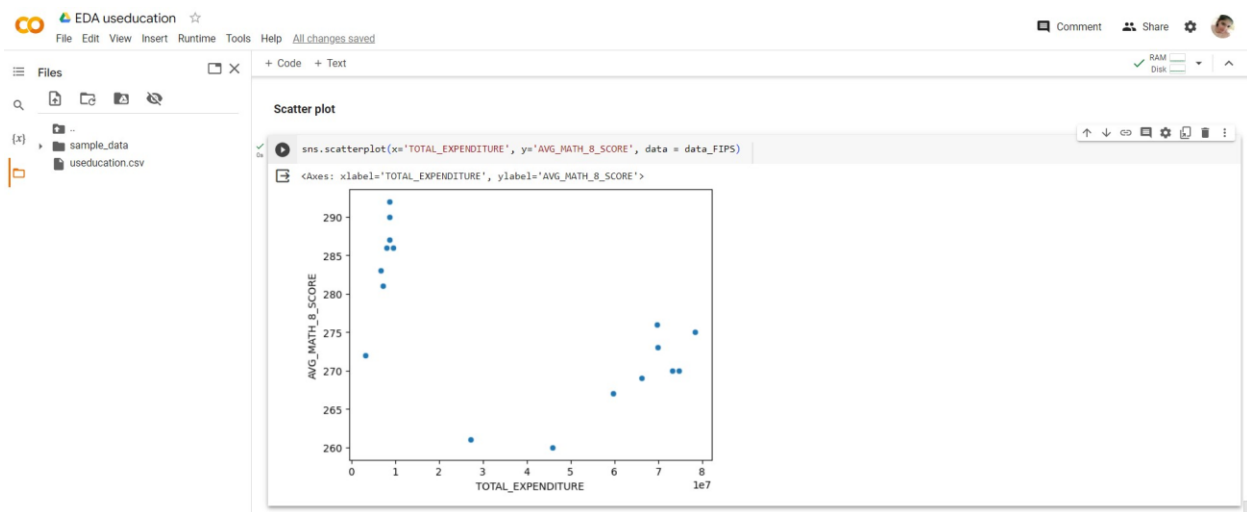
Interpretation:

The histogram showcases the distribution of data for the "GRADES_ALL_G" variable, which seems to represent the total grades for all students. The x-axis, labeled "GRADES_ALL_G," spans from around 2 million to slightly above 6 million. The y-axis, labeled "Count," indicates the frequency of observations. Most prominently, there's a significant peak at the leftmost side, suggesting that a substantial number of observations have grade totals close to 2 million. A smaller peak, or secondary mode, is noticeable between 5 million to slightly above 6 million. The region between these two peaks has fewer observations, with a particularly notable dip around the 4 million mark. In summary, the distribution

appears bimodal, with a significant number of entities having grade totals near 2 million and another notable group clustered just above 5 million.

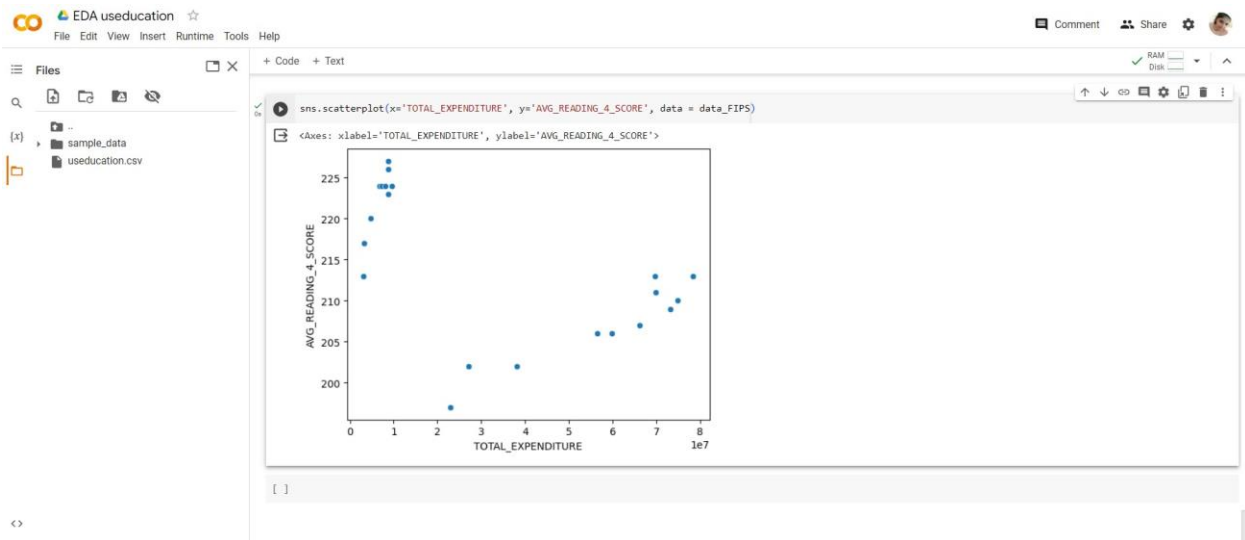
Scater Plot

In Python's EDA, a scatterplot is a graphical display that shows discrete data points as dots on a two-dimensional graph. It helps in the analysis of data patterns and correlations by revealing the relationship between two numerical variables.



Interpretation

The scatter plot displays the relationship between "TOTAL_EXPENDITURE" and "AVG_MATH_8_SCORE." The x-axis represents "TOTAL_EXPENDITURE" and spans from around 0 to slightly above 8×10^7 . On the y-axis, the "AVG_MATH_8_SCORE" fluctuates between approximately 260 and 290. The plot suggests a mild positive correlation between the variables: as the "TOTAL_EXPENDITURE" increases, the "AVG_MATH_8_SCORE" generally seems to rise. However, this relationship is not very tight, as the data points are somewhat dispersed. While the plot implies that greater expenditure might be associated with higher math scores for 8th graders, the wide distribution of data points indicates other influencing factors or the possibility of outliers. In summary, there might be a link between expenditure and math scores, but the data suggests that other factors might also play a significant role in determining scores.



Interpretation

The scatter plot visualizes the relationship between "TOTAL_EXPENDITURE" and "AVG_READING_4_SCORE." On the x-axis, "TOTAL_EXPENDITURE" seems to be represented in a range that spans from approximately 0 to slightly above 8×10^7 . On the y-axis, the "AVG_READING_4_SCORE" appears to range from around 200 to slightly above 225. There seems to be a modest positive correlation between the two variables: as the "TOTAL_EXPENDITURE" increases, the "AVG_READING_4_SCORE" generally appears to increase as well. However, the relationship is not perfectly linear, as there are data points scattered throughout, indicating potential outliers or other factors at play. Overall, while there's a suggestion that higher expenditure might be associated with better reading scores for 4th graders, the scatter of data points indicates that it's essential to consider other influencing factors or investigate the relationship further.

Citations:

1. U.S Education Datasets: <https://www.kaggle.com/datasets/noriuk/us-education-datasets-unification-project>

2. Google colab file

<https://colab.research.google.com/drive/111YU50NgpZXope4yIrtA49k6KqvsYTfp#scrollTo=cmqDoXL4Wxbw>

3. ORES_Group_dataset-GoogleDrive.(n.d.).

https://drive.google.com/drive/folders/173tJ7JkxJeu9Pi62k9t00J40alAZfPP0?usp=share_link