

Week08 Assignment

Exploratory Data Analysis with Python (pandas)

Venkata Sai Pallavi Pallapolu

Saint louis university

Course code: ORES-5160: data management

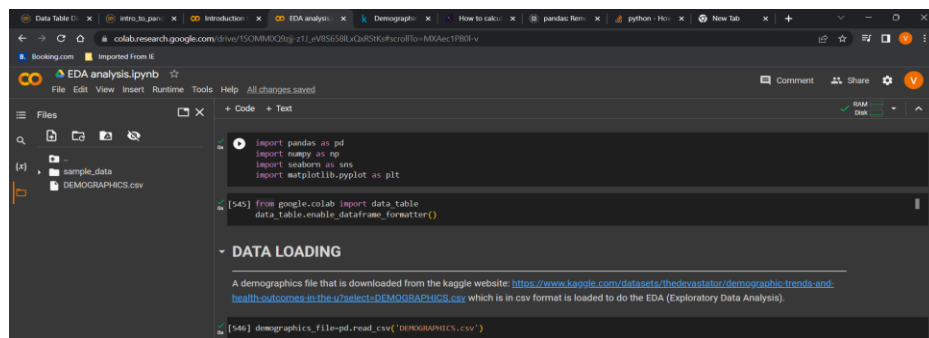
Prof. Jason Eden

October 26th, 2023

RESEARCH REPORT ON EXPLORATORY DATA ANALYSIS WITH PYTHON (PANDAS):

For Identifying outliers, trends and removing data outliers the Exploratory data analysis is performed by following the steps like: loading the data, imputation of the data, identifying the outliers, data visualization and generating insights.

1. **Data loading:** demographics dataset is downloaded on local computer as a csv (comma-separated value) file and then loaded into the google colab platform by clicking on the upload option inside the file icon. Then, to perform the Exploratory Data Analysis on the loaded dataset pandas, NumPy, Seaborn, and Matplotlib libraries are imported in the code. The code is as follows:



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

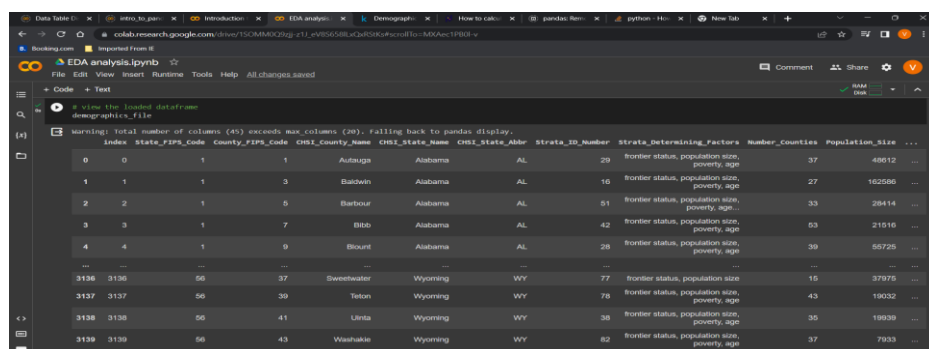
from google.colab import data_table
data_table.enable_dataframe_formatter()

# DATA LOADING

A demographics file that is downloaded from the kaggle website: https://www.kaggle.com/datasets/thedevastator/demographic-trends-and-health-outcomes-in-the-u/s?select=DEMOGRAPHICS.csv which is in csv format is loaded to do the EDA (Exploratory Data Analysis).

demographics_file=pd.read_csv("DEMOGRAPHICS.csv")
```

The loaded data Frame looks something like the picture below indicating all the series (columns) and their values. The bottom of the output cell indicates the total number of rows and columns which is 3141 rows and 45 columns in this case.



	index	State_FIPS_Code	County_FIPS_Code	OnMI_County_Name	OnMI_State_Name	OnMI_State_Abb	Strata_ID_Number	Strata_Determining_Factors	Number_Counties	Population_Size	...
0	0	1	1	Autauga	Alabama	AL	29	frontier status, population size, poverty, age	37	48612	...
1	1	1	3	Baldwin	Alabama	AL	16	frontier status, population size, poverty, age	27	162586	...
2	2	1	5	Barbour	Alabama	AL	51	frontier status, population size, poverty, age	33	28414	...
3	3	1	7	Beale	Alabama	AL	42	frontier status, population size, poverty, age	53	21516	...
4	4	1	9	Blount	Alabama	AL	28	frontier status, population size, poverty, age	39	55725	...
...
3136	3136	56	37	Sweetwater	Wyoming	WY	77	frontier status, population size	15	37975	...
3137	3137	56	39	Teton	Wyoming	WY	78	frontier status, population size, poverty, age	43	19032	...
3138	3138	56	41	Uinta	Wyoming	WY	38	frontier status, population size, poverty, age	35	19039	...
3139	3139	56	43	Washakie	Wyoming	WY	82	frontier status, population size, poverty, age	37	7933	...

2. **Exploring the Data:** To know what the loaded dataframe contains, it is explored by using several functions:

- The number of rows and columns: by using `dataframe.shape()` function
- Dimensions: by using `dataframe.ndim()` function- gives the number of dimensions in the dataframe
- Initial rows of the dataframe: by using `dataframe.head()` function- produces the first 5 rows of dataframe.

```

# describes all the statistic values for every column like distinct values, frequency, standard deviation, minimum, maximum, etc.
demographics_file.describe(include="all")

```

Index	state_FIPS_Code	County_FIPS_Code	Cnty_Name	Cnty_State_Abbrev	State_ID_Number	State_Determining_Factor	Number_Counties	Population_Size
3141.000000	3141.000000	3141.000000	3141	3141	3141	3141.000000	3141	31410000e+03
NaN	NaN	NaN	1847	81	81	NaN	4	NaN
NaN	NaN	NaN	Washington	Texas	TX	NaN	frontier status, population size, poverty, age	NaN
NaN	NaN	NaN	32	254	254	NaN	1702	NaN
1870.000000	30.304890	103.716891	NaN	NaN	NaN	44.08275	NaN	38.486151
908.872824	15.134423	107.980484	NaN	NaN	NaN	25.118134	NaN	10.280185
5.000000	1.000000	1.000000	NaN	NaN	NaN	1.000000	NaN	15.000000
795.000000	15.000000	35.000000	NaN	NaN	NaN	23.000000	NaN	32.000000
1570.000000	25.000000	79.000000	NaN	NaN	NaN	44.000000	NaN	37.000000
2355.000000	45.000000	133.000000	NaN	NaN	NaN	66.000000	NaN	45.000000
3140.000000	55.000000	840.000000	NaN	NaN	NaN	88.000000	NaN	62.000000

This above picture indicates the usage of `describe()` function to produce the output showing the statistics of each series that includes the count, distinct values, top value (if it's a category), frequency, mean, standard deviation, minimum value, 25th percent value, 50th percent value, 75th percent value, maximum value.

```

demographics_file.dropna(how="all", axis=0)
demographics_file.describe()

```

Index	Number_Counties	Population_Size	Population_Density	Poverty	Age_19_Under	Age_19_64	Age_65_84	Age_85_and_Over	White	Black
count	3141.0	3141.0	3141.0	3141.0	3141.0	3141.0	3141.0	3141.0	3141.0	3141.0
mean	38.4861509735434	94368.16427889207	249.11938872970393	12.638427252467368	24.806526583890477	60.289398280802295	12.78943011779688	2.115409105380452	87.01789239095828	8.986692136262336
std	10.290194571527541	306431.655763125	1703.041884190953	40.18646099222744	3.28177733639539	3.35056269254836	3.334034542589774	0.949118632269374	16.150478631148	14.545659099557506
min	15.0	62.0	-2222.0	-2222.0	1.4	47.6	2.1	0.1	4.7	0.0
25%	32.0	11211.0	17.0	9.8	22.7	58.3	10.7	1.5	82.8	0.5
50%	37.0	25235.0	44.0	12.6	24.6	60.3	12.5	1.9	94.1	2.1
75%	45.0	64040.0	109.0	16.2	26.4	62.3	14.7	2.6	97.6	10.3
max	62.0	9935475.0	69390.0	36.2	47.2	83.3	29.2	7.6	100.0	86.0

The above picture indicates the usage of `dropna()` function. This is a function used to remove not a number (NaN) values, null values, or missing values. If any missing values are present in the dataframe to remove them, this function has been used. Also, this code indicates 'axis=0', which means rows that is to remove any row which has missing values (Axis=1 for columns).

```
#to know the unique states and their names and number of unique counties in each state
print(len(demographics_file.GSI_State_Name.unique()))
print(demographics_file.GSI_State_Name.unique())
print(len(demographics_file.GSI_County_Name.unique()))
```

```
51
['Alabama' 'Alaska' 'Arizona' 'Arkansas' 'California' 'Colorado'
 'Connecticut' 'Delaware' 'District of Columbia' 'Florida' 'Georgia'
 'Hawaii' 'Idaho' 'Illinois' 'Indiana' 'Iowa' 'Kansas' 'Kentucky'
 'Louisiana' 'Maine' 'Maryland' 'Massachusetts' 'Michigan' 'Minnesota'
 'Mississippi' 'Missouri' 'Montana' 'Nebraska' 'Nevada' 'New Hampshire'
 'New Jersey' 'New Mexico' 'New York' 'North Carolina' 'North Dakota'
 'Ohio' 'Oklahoma' 'Oregon' 'Pennsylvania' 'Rhode Island' 'South Carolina'
 'South Dakota' 'Tennessee' 'Texas' 'Utah' 'Vermont' 'Virginia'
 'Washington' 'West Virginia' 'Wisconsin' 'Wyoming']
1847
```

The skewness of the dataset can be known by using skew () function. The dataframe that has been loaded is highly skewed as the columns indicated the skewness as 27.56 as highest and -54.80 as least so this dataset can be thought as highly skewed one (2).

```
skew=demographics_file.skew(axis=0)
print(skew)
```

```
Number_Countries    0.382221
Population_Size     15.273352
Population_Density   27.566085
Poverty             -54.805956
Age_19_Under        0.724751
Age_19_54           0.382267
Age_65_84           0.530835
Age_85_and_Over     1.155888
White               -1.999811
Black               2.244748
Native_American     7.862679
Asian              10.770755
Hispanic            3.543323
dtype: float64
<ipython-input-881-488881710220>:1: FutureWarning: the default value of numeric_only in DataFrame.skew is deprecated. In a future version, it will default to False. In addition, specify skew=demographics_file.skew(axis=0)
```

The above picture also indicates a code to print the distinct number of states from the dataframe, their names, and the number of distinct counties in those states.

```
# to know the relationship between population size and poverty
demographics_file[["Population_Size","Poverty"]].describe()
```

Index	Population_Size	Poverty
count	3141.0	3141.0
mean	94368.16427889207	12.638427752467368
std	306431.655763125	40.10646090222744
min	62.0	-2222.2
25%	11211.0	9.8
50%	25235.0	12.6
75%	64040.0	16.2
max	9935475.0	36.2

The above picture shows the relationship between population size and poverty. That is for an average of 94368.16 population the poverty rate is 12.63. It also shows other descriptive statistical values.

Removing the unnecessary columns from the dataframe:

```

demographics_file.drop(['Index', 'State_FIPS_Code', 'County_FIPS_Code', 'CHS_State_Abb', 'strata_ID_Number', 'strata_Determining_Factors', 'Min_Population_Size', 'Max_Population_Size'], axis=1, inplace=True)

demographics_file.info()

```

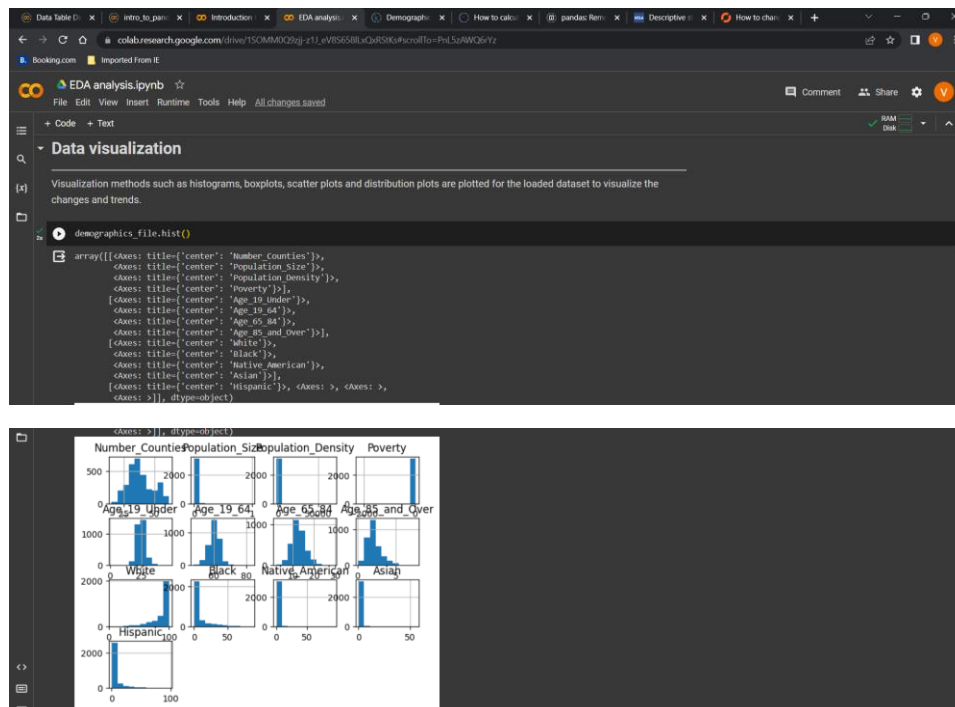
```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1141 entries, 0 to 1140
Data columns (total 14 columns):
 #   Column              Non-Null Count  Dtype  
---  -
 0   CHS_County_Name     1141 non-null  object  
 1   CHS_State_Name      1141 non-null  object  
 2   Number_Counties     1141 non-null  int64   
 3   Population_Size     1141 non-null  int64   
 4   Population_Density  1141 non-null  float64  
 5   Poverty             1141 non-null  float64  
 6   Age_19_Under       1141 non-null  float64  
 7   Age_19_64         1141 non-null  float64  
 8   Age_65_84         1141 non-null  float64  
 9   Age_85_and_Over    1141 non-null  float64  
10   White              1141 non-null  float64  
11   Black              1141 non-null  float64  
12   Native_American    1141 non-null  float64  
13   Asian              1141 non-null  float64  
14   Hispanic            1141 non-null  float64  
dtypes: float64(10), int64(3), object(2)
memory usage: 104.2+ KB

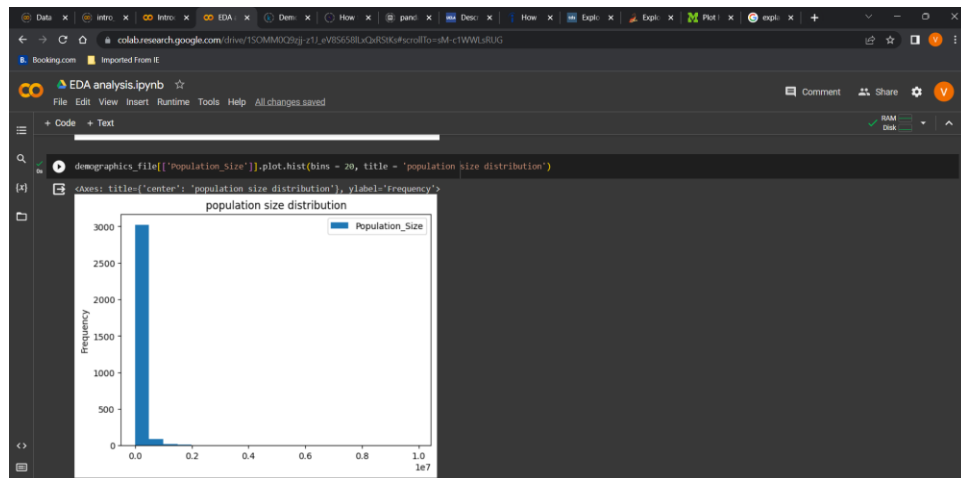
```

This above picture indicates the names of series and their datatypes after the removal of unnecessary columns.

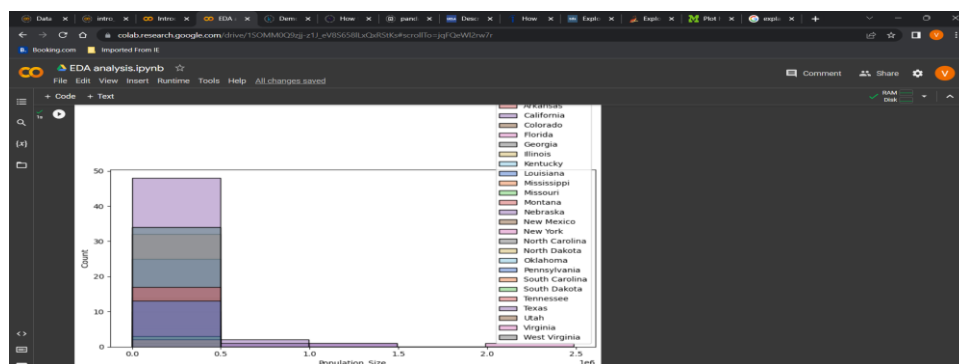
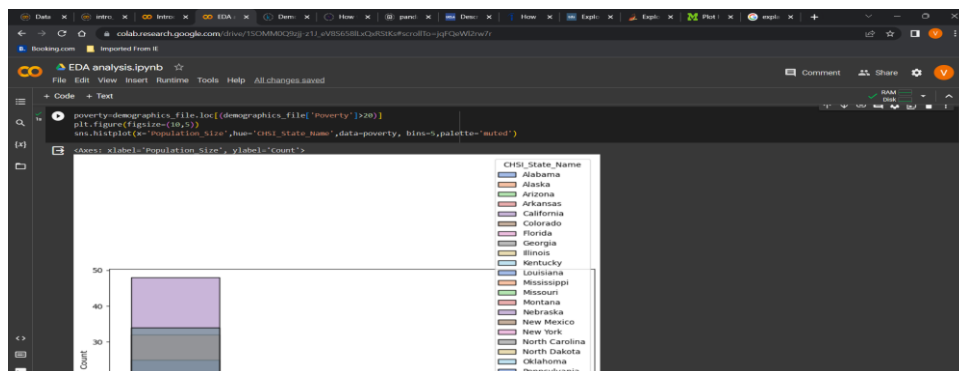
3. **Data visualization:** Visualizations can provide valuable insights into your data that might not be apparent from looking at the raw data alone (3). The various data visualization methods used in pandas for EDA are histograms, boxplots, scatterplots, pie charts, distribution charts and many more.
 - **HISTOGRAMS:** Histograms are a great way to visualize the distribution of your data. They can provide insights into the central tendency, variability, and skewness of the data. These histograms divide the numerical values into bins and count the number of observations that fall into each bin (4).



The above pictures indicate the histograms for all the useful columns in the dataframe by writing the code dataframe.hist() using hist () function.

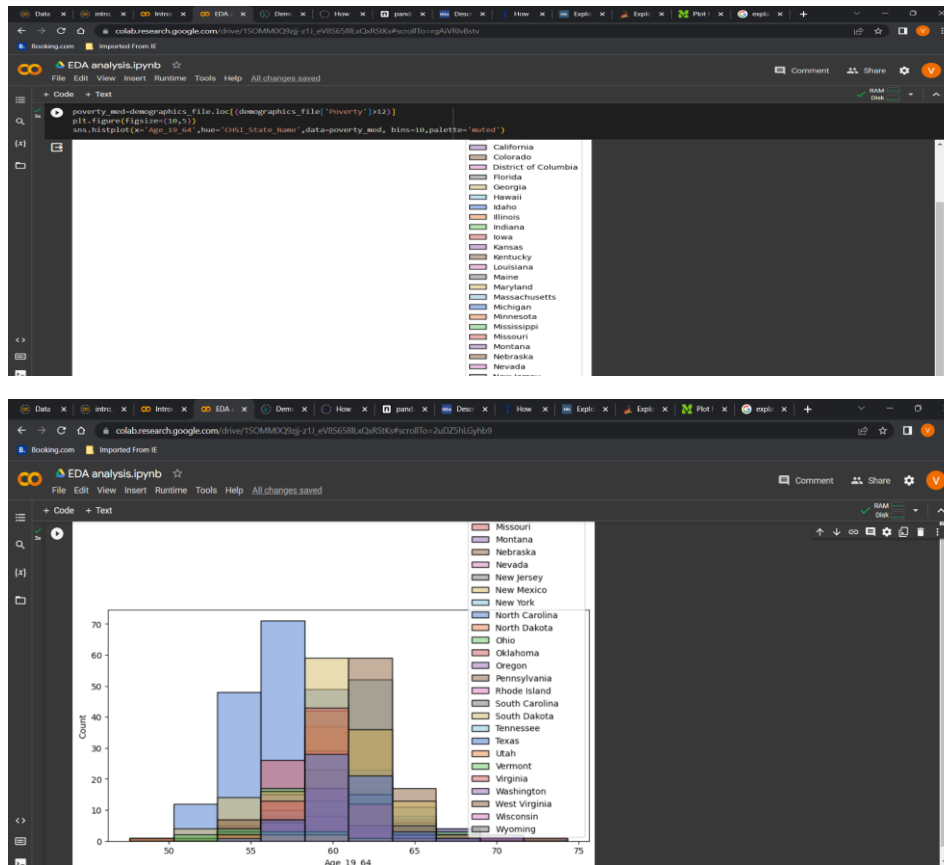


The above picture indicates the population size distribution that is the data that has been aggregated based on different population size intervals. The height of each bar represents the frequency of occurrences within that range. This histogram shows the asymmetric distribution suggesting highly skewness and uneven spread. From this histogram, population size with range 0.005 -0.1 can be identified as the most frequently occurring range (central tendency).



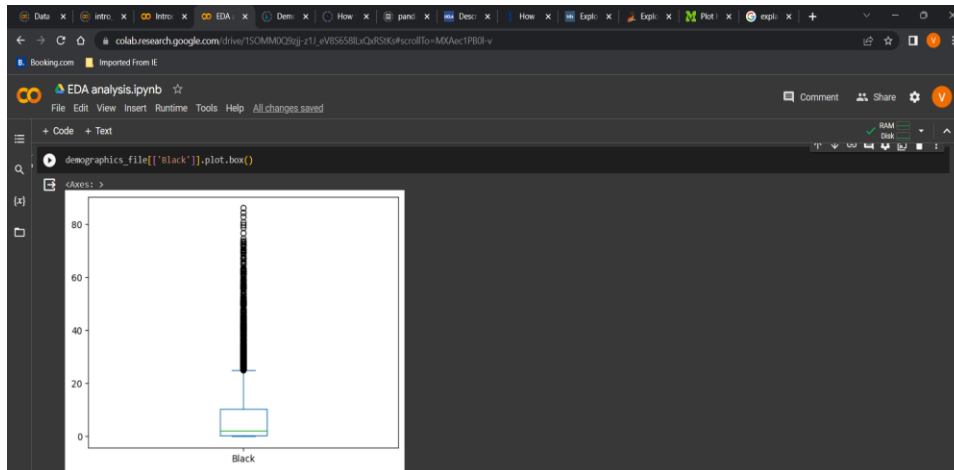
The above pictures produce the histogram that shows the distribution of population size with poverty greater than 20 categorized by different states by using seaborn library. X-axis represents the population size with state name as hue that means histogram will display different colors for each state and y-axis represents the count of counties falling within each population size range and the number of intervals or bins are 5 for this histogram. By looking at the histogram, it is clear

that 'Nebraska' has a greater number of counties with poverty greater than 20 with population size under 500000.

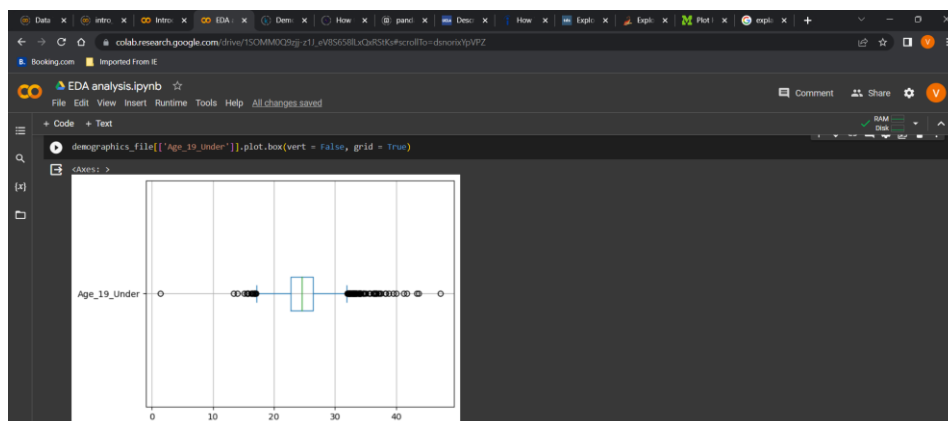


The above pictures indicate the distribution of age 19-64 people with poverty greater than 12 categorized by different states. The code also used a filter condition that is poverty greater than 12 and replaced it in the data. In the code `palette='muted'` indicates the color palette. From the histogram, it can be said that 'Alabama', 'Texas' have the highest number of people in age groups 19-64 .

- **Boxplots:** boxplots can tell you about outliers and what their values are. They can also tell you if a data is symmetrical, how tightly the data is grouped and if and how the data is skewed (5).

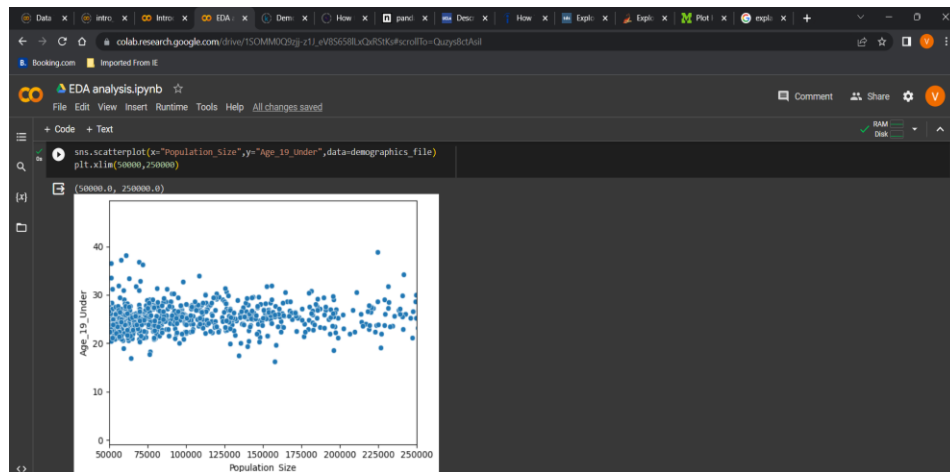


The above picture shows the code by using `plot.box()` function to get the distribution, central tendency and spread of 'Black' variable (race) from demographics file that is the count of population within a state. The line in the middle of the box indicates the median value of the data. The bottom of the box represents the first quartile, and the top indicates the third quartile while the box indicates the interquartile range. The points beyond the black line (Whiskers) are considered as outliers.

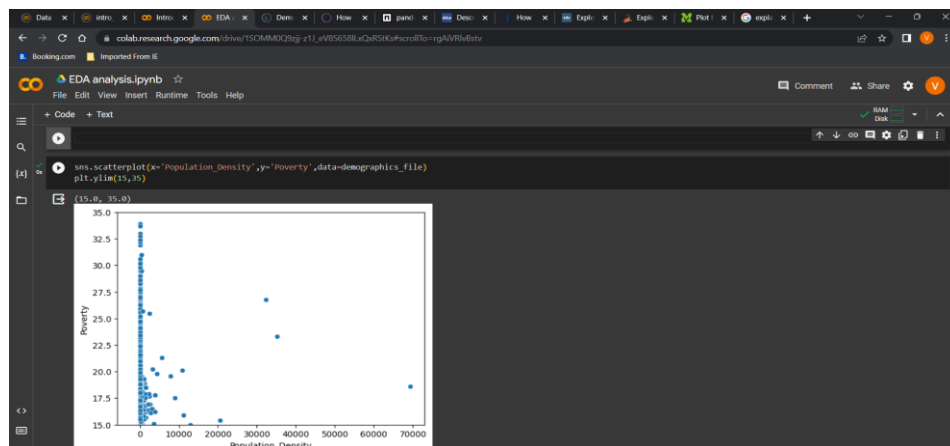


The above picture represents another format for box plot (6) to represent the plot in horizontal direction by using `plot.box(vert=False and grid=True)` meaning the plot should be in grid format not vertical. The box plot that resulted contains the distribution of 'age 19 and under' variable from demographics file, its median, minimum value, maximum value and outliers. From the boxplot, it is clear that there are a median of 23 people in each county under age 19.

- **Scatter plots:** A scatter plot uses dots to represent values for two different numeric variables. Each dot represents a single tree; each point's horizontal position indicates that tree's diameter (in centimeters) and the vertical position indicates that tree's height (in meters) (7).



The above picture indicates the scatter plot between two variables they are: 'age 19 under' and 'population size'. The code is written by using seaborn library and gives x-axis a parameter: population size and on y-axis: age 19 under is represented. And also, the limit for population size is set between 50000 and 250000. 20-30 range in age 19 under is the frequent range for population size under age 19 for this scatter plot.



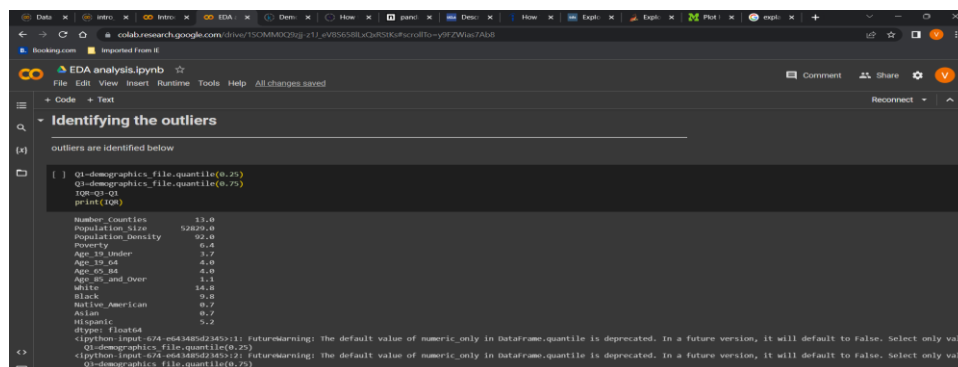
The above picture indicates the scatter plot which on x-axis have population density of various states living in 15-35 range of poverty and on y-axis poverty is represented. It also has a limit of 15 to 35. This scatterplot explains the relationship between poverty and population density.

- **Bar plots:** A bar plot or bar chart is a graph that represents the category of data with rectangular bars with lengths and heights that is proportional to the values which they represent (8).

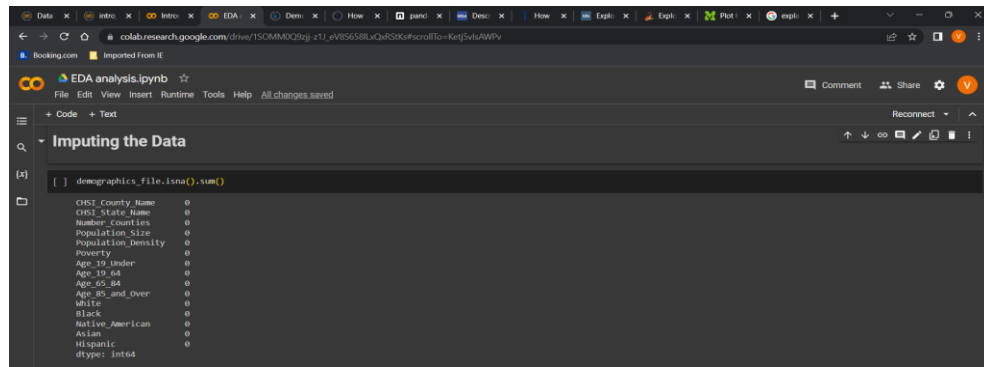


The above picture indicates the bar plot by using matplotlib, seaborn libraries which indicates the age 65-84 on x-axis and black column on y-axis and hue is state name the condition here is poverty greater than 30. Distribution of black with age 65-84 from various states having poverty more than 30 is indicated in this bar chart.

4. **Identifying the outliers:** The outliers can be identified by calculating the inter quartile range of the demographics file which is the difference between third quartile and first quartile. The below picture indicates the interquartile range for all the columns and their values are as follows:



5. **Data imputation:** In statistics, imputation refers to the process of substituting alternate values for missing data. It is known as "unit imputation" when replacing a data point and "item imputation" when replacing a data point's component (8).



```
[ ] demographics_file.isna().sum()

CHS1_County_Name      0
CHS1_State_Name        0
Number_Counties       0
Population_Size       0
Population_Density     0
Poverty               0
Age_19_Under          0
Age_19_64             0
Age_65_84             0
Age_85_and_Over       0
White                 0
Black                 0
Native_American       0
Asian                 0
Hispanic              0
dtype: int64
```

This picture indicates that there are no null values in the data by using the function `.isna().sum()` so, there is no need to replace the null values.

6. Gathering insights: on performing the Exploratory Data Analysis, this dataset has provided with some insights like understanding the distribution and characteristics of counties within different states, poverty rate of various states with different population size, different ethnic composition like white, black, native American, Hispanic, Asian to understand the ethnic diversity and distribution across various states and also proportions of various age groups and their distribution.

Conclusion:

By performing the above-mentioned steps such as data loading, exploring the data, data visualization, data imputation, identifying the outliers, and gathering insights, valuable insights into the demographic and socioeconomic characteristics of the population across various states and counties are drawn. The study revealed trends in population size, population density, age distribution, poverty rates, and ethnic diversity. Several visualizations, such as histograms, boxplots, scatter plots, and bar plots aided in the identification of patterns, outliers, and potential correlations in the data.

This EDA provides a full overview of the demographics dataset's important aspects and informs possibilities for more specialized study and decision-making procedure.

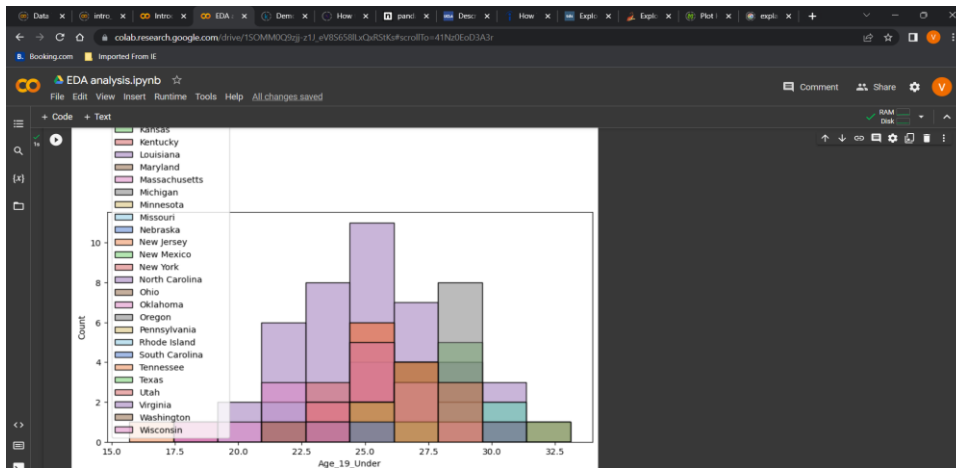
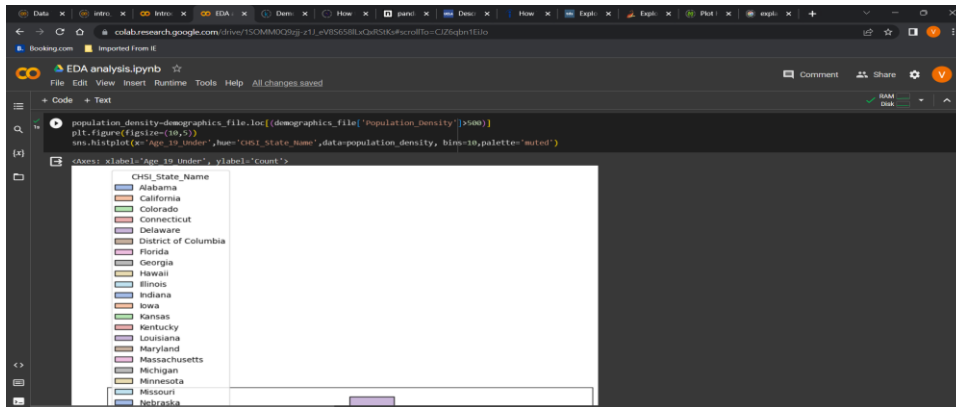
References:

1. <https://saturncloud.io/blog/pandas-tips-change-column-type/>
2. [https://medium.com/@atanudan/kurtosis-skew-function-in-pandas-aa63d72e20de#:~:text=skewness\(\)%20function%20in%20pandas,present%20in%20the%20DataFrame%20object.](https://medium.com/@atanudan/kurtosis-skew-function-in-pandas-aa63d72e20de#:~:text=skewness()%20function%20in%20pandas,present%20in%20the%20DataFrame%20object.)
3. <https://docs.kanaries.net/articles/exploratory-data-analysis-python-pandas>
4. https://mode.com/example-gallery/python_histogram/#:~:text=A%20histogram%20divides%20the%20values,of%20values%20within%20a%20variable.
5. <https://builtin.com/data-science/boxplot>
6. <https://towardsdatascience.com/exploratory-data-analysis-eda-visualization-using-pandas-ca5a04271607>

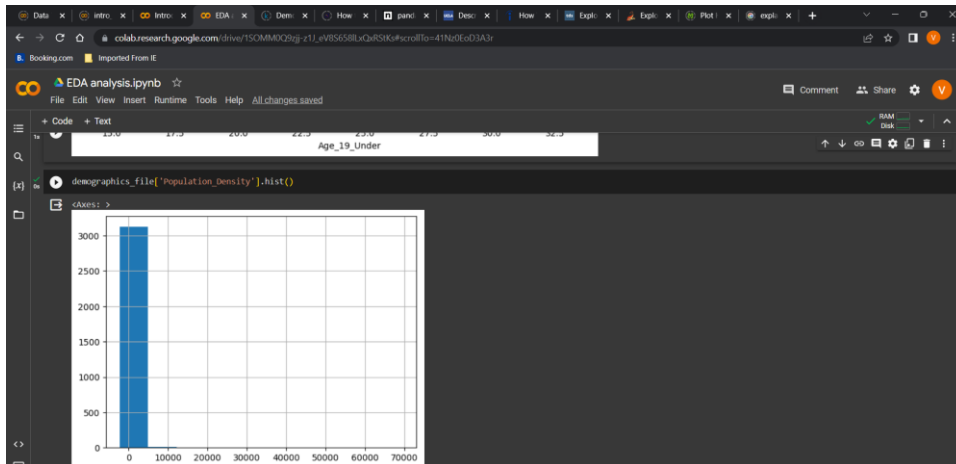
7. <https://chartio.com/learn/charts/what-is-a-scatter-plot/#:~:text=What%20is%20a%20scatter%20plot,to%20observe%20relationships%20between%20variables.>
8. <https://www.simplilearn.com/data-imputation-article#:~:text=ProgramExplore%20Program-,What%20Is%20Data%20Imputation%3F,from%20a%20dataset%20each%20time.>

APPENDIX:

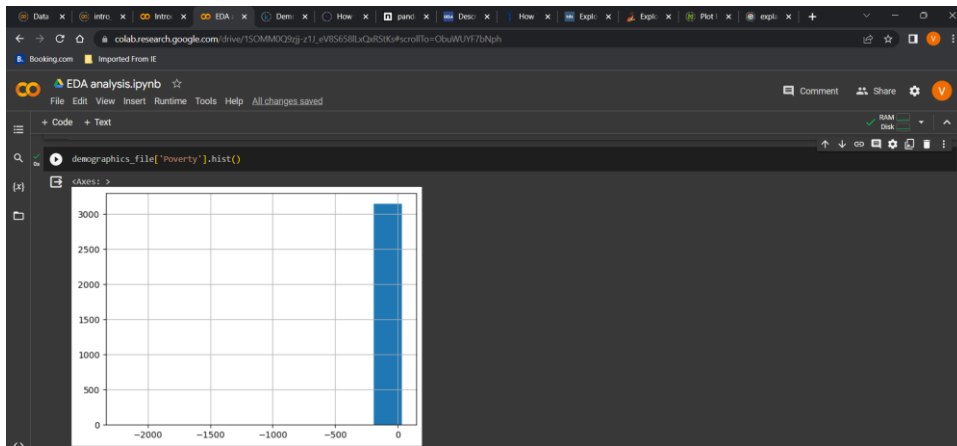
- Histogram indicating the distribution of age 19 under with population density greater than 500 over different states.



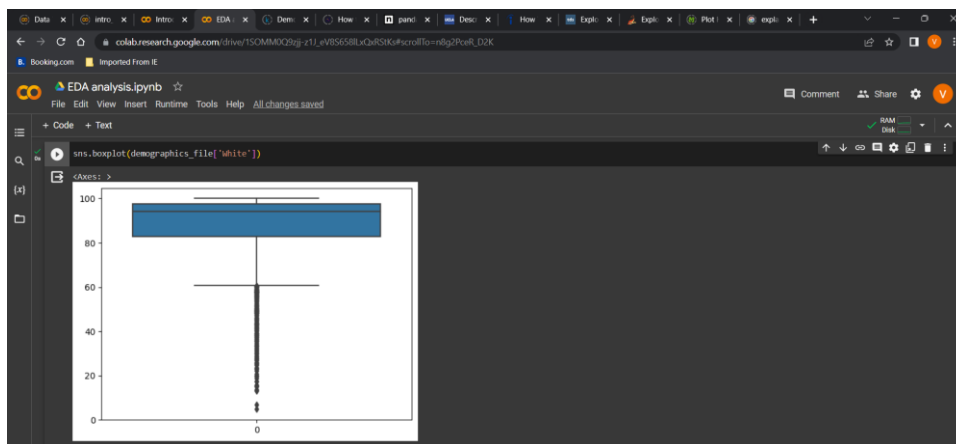
- Histogram indicating the distribution of population density



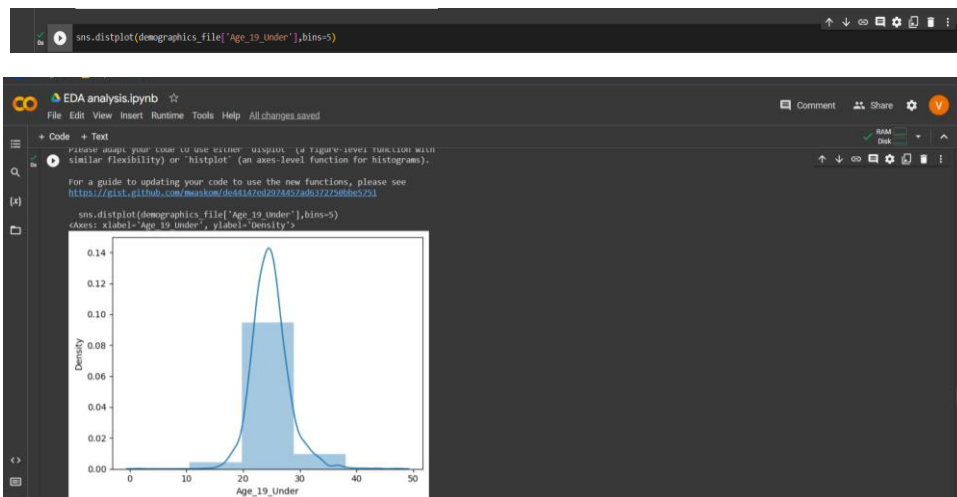
- Histogram indicating the distribution of poverty



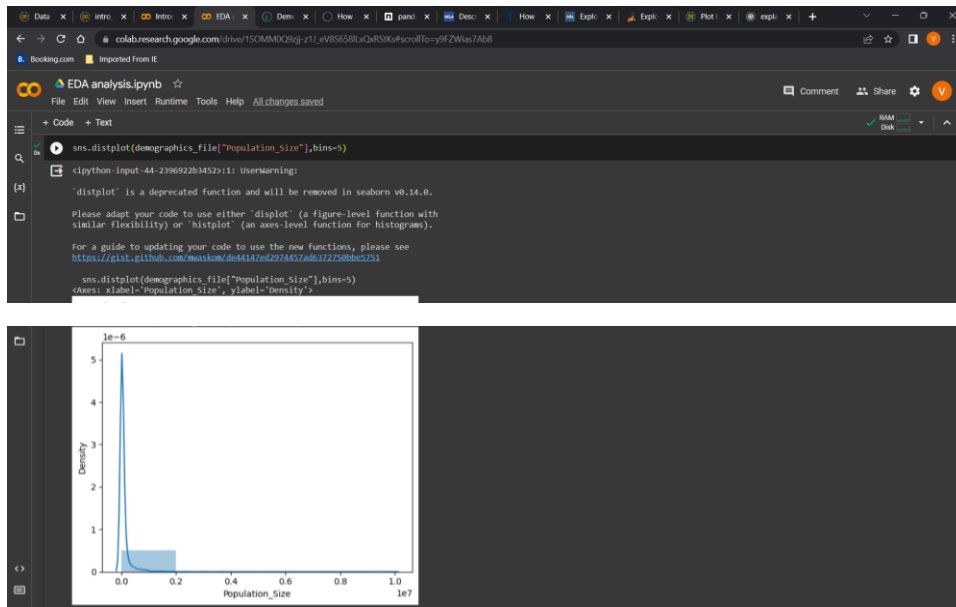
- Box plot indicating the distribution of white column from dataframe



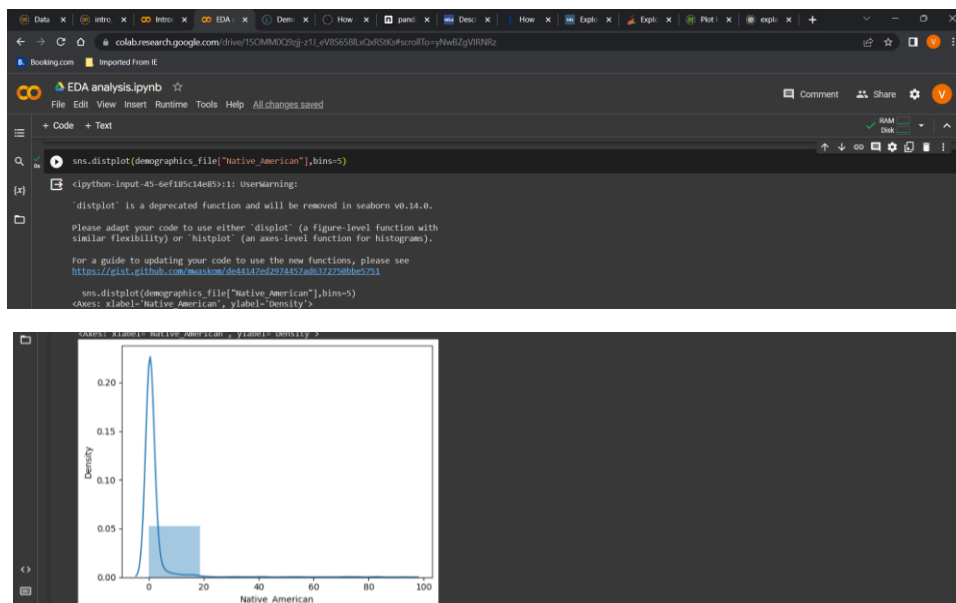
- Distribution plot indicating the distribution of age 19 under



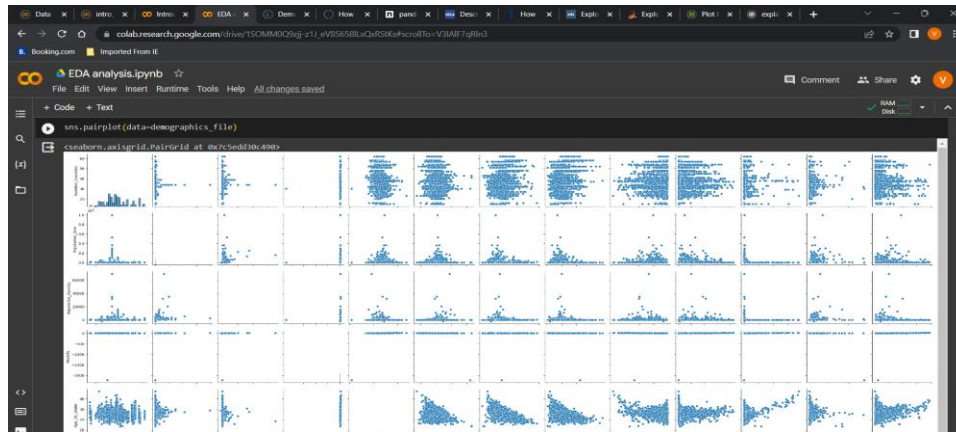
- Distribution plot indicating the distribution of population size column



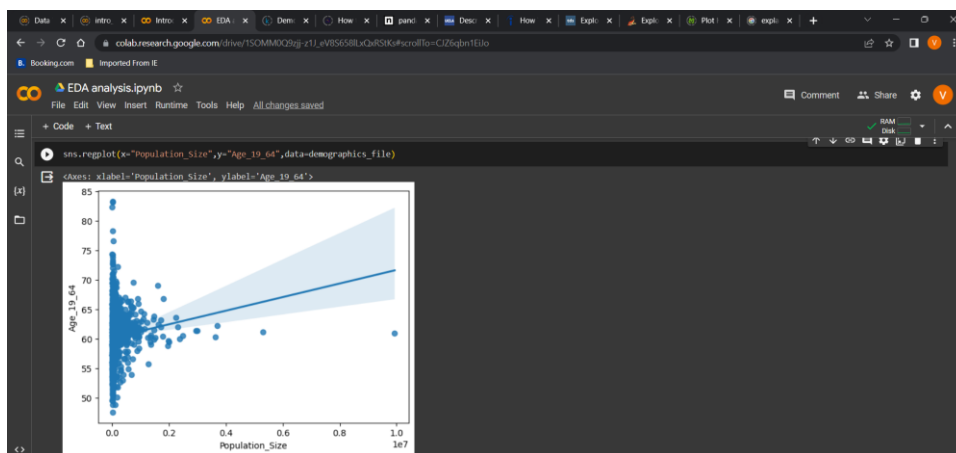
- Distribution plot indicating the distribution of native american column from dataframe



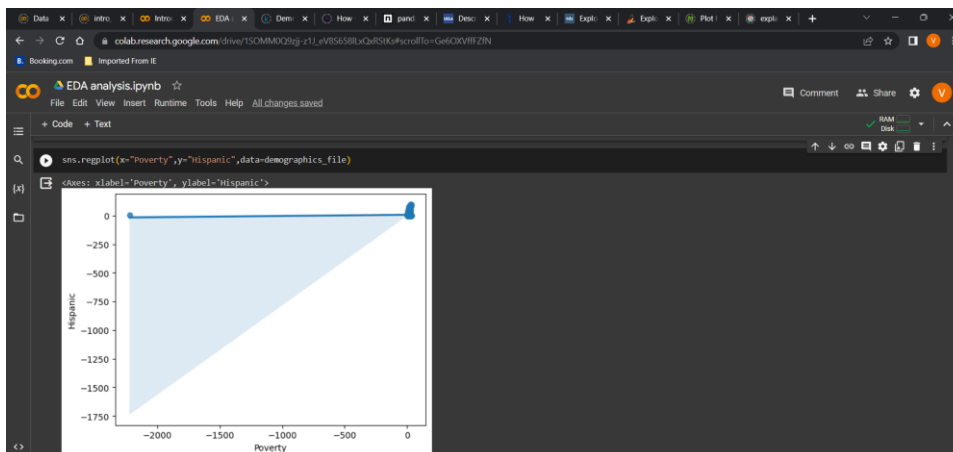
- Pairplot of all the columns



- Regression plot of population size and age 19-64 on x-axis and y-axis respectively



- Regression plot for distribution of poverty and hispanic on x-axis and y-axis respectively



- Correlation plot for population size, age 19-64 and poverty

EDA analysis.ipynb

the below method is referred from this website <https://towardsdatascience.com/exploratory-data-analysis-eda-visualization-using-pandas-ca5894771607>

```
[54]: corr = demographics_file[['Population_Size', 'Age_19_64', 'Poverty']].corr()
      corr.style.background_gradient(cmap='bwr').set_precision(2)
```

<ipython-input-54-38758c9ce841>:2: FutureWarning: this method is deprecated in favour of 'Styler.format(precision=...)'
corr.style.background_gradient(cmap='bwr').set_precision(2)

	Population_Size	Age_19_64	Poverty
Population_Size	1.00	0.10	-0.00
Age_19_64	0.10	1.00	-0.14
Poverty	-0.00	-0.14	1.00