# College Analytics

**A Visual storytelling of Colleges in US**

## Data Analysis Capstone Project

-Shilpa Rao

# Checkpoint 1 : Selecting your business issue and dataset

The data set I am using gives a detailed listing of top American colleges. This project will allow prospective students and families to get a better understanding of the opportunities in front of them. College searching can be very difficult to navigate especially for immigrant families or those who have not gone through this process before. This visual storytelling is an attempt to make this process easier for all
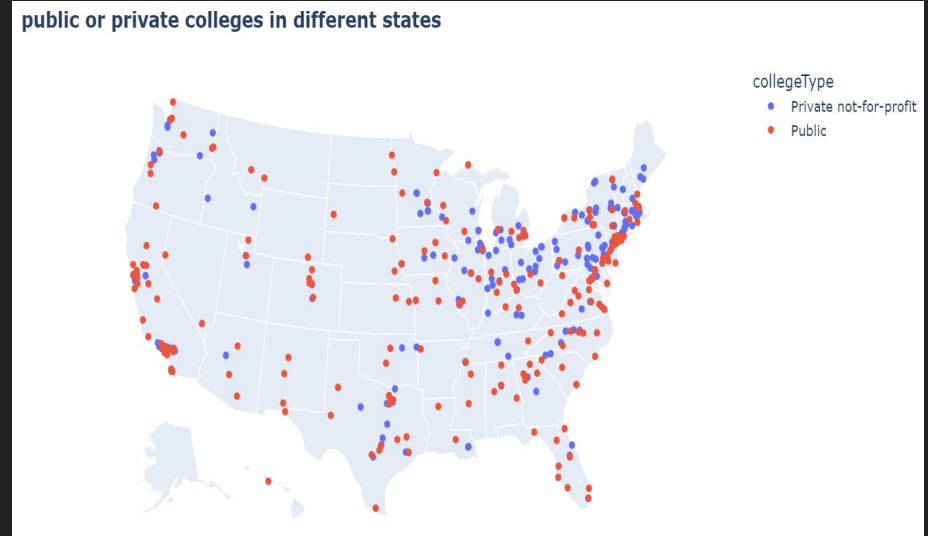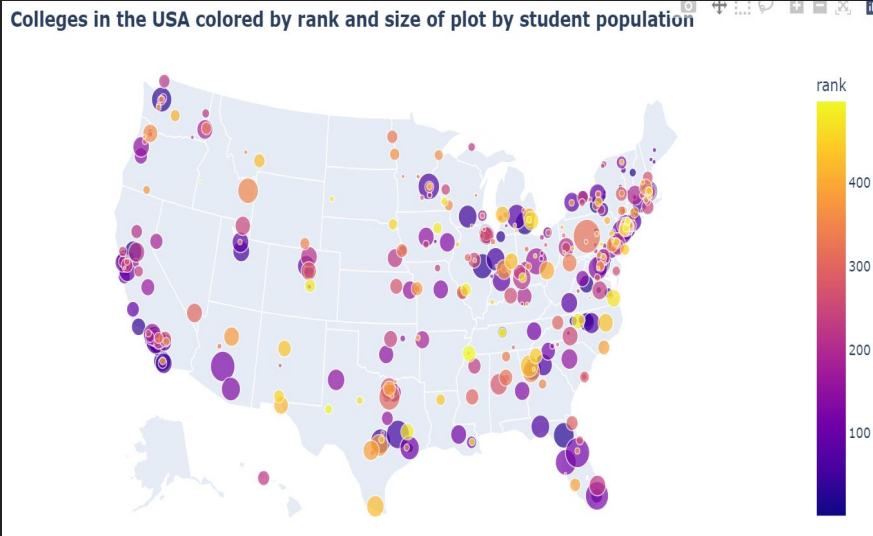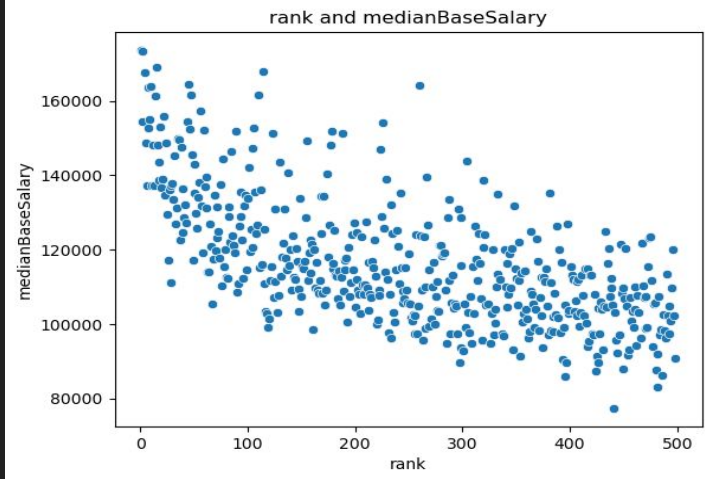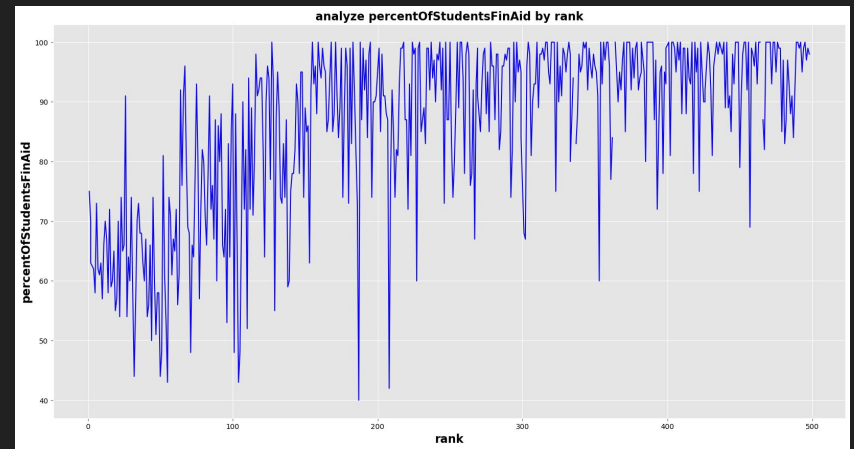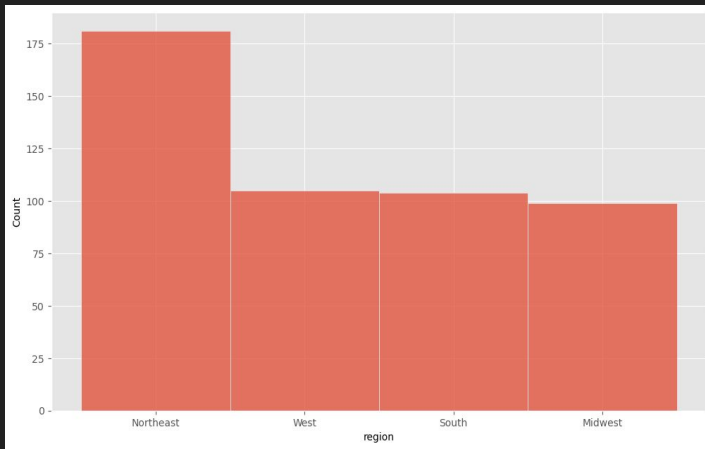Problem : Navigating search for Colleges in US

❖ Solution : Provide a visual analysis of Colleges in US for future students
❖ Feature : Give a comparison of rank of colleges to median base salary
❖ Feature : Give a visual map of public and private schools
❖ Feature: Show the opportunity of grants and financial aids
❖ Full Project Tech stack: Python, Tableau, Jupyter notebook

Dataset I used : https://www.kaggle.com/datasets/kabhishm/top-american-colleges-2022

# Checkpoint 2 : Exploratory Data Analysis(EDA)

❖ Understand the data set by checking info, description, first few rows, shape..
❖ Created multiple visualizations to understand the dataset



Colleges in the USA colored by rank and size of plot by student population



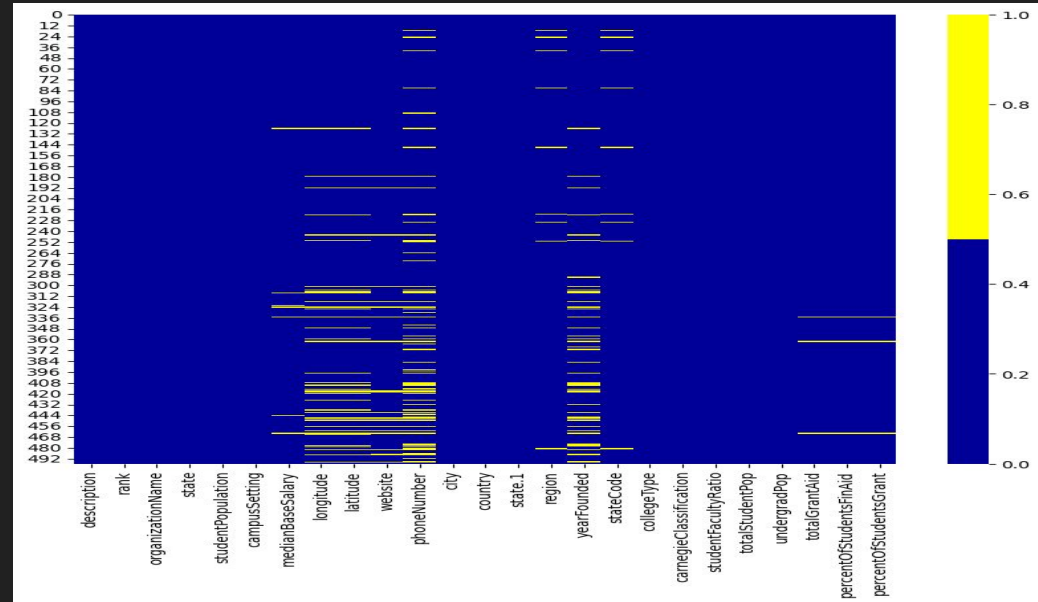public or private colleges in different states

- ❖ colleges are concentrated in coastal locations(especially northeast) and lesser in midwest
- ❖ looks like higher ranked and higher population schools are connected and are higher in number along the coast
- ❖ missing values, and repeated columns need handling
- ❖ unnecessary columns such as website, phone number, and country name(since all colleges here are in US) can be removed

# Checkpoint 3 : Cleaning Data

❖ Identified missing data and visualize it

❖ Dropped website, phone Number - lots of missing data and irrelevant

❖ Filled missing medianBaseSalaray, totalGrantAid with mean

❖ Imputed the missing latitude and longitude values - not the best options

❖ Drop studentPopulation as it has same values as totalStudentPop - unnecessary data

Graph of missing data

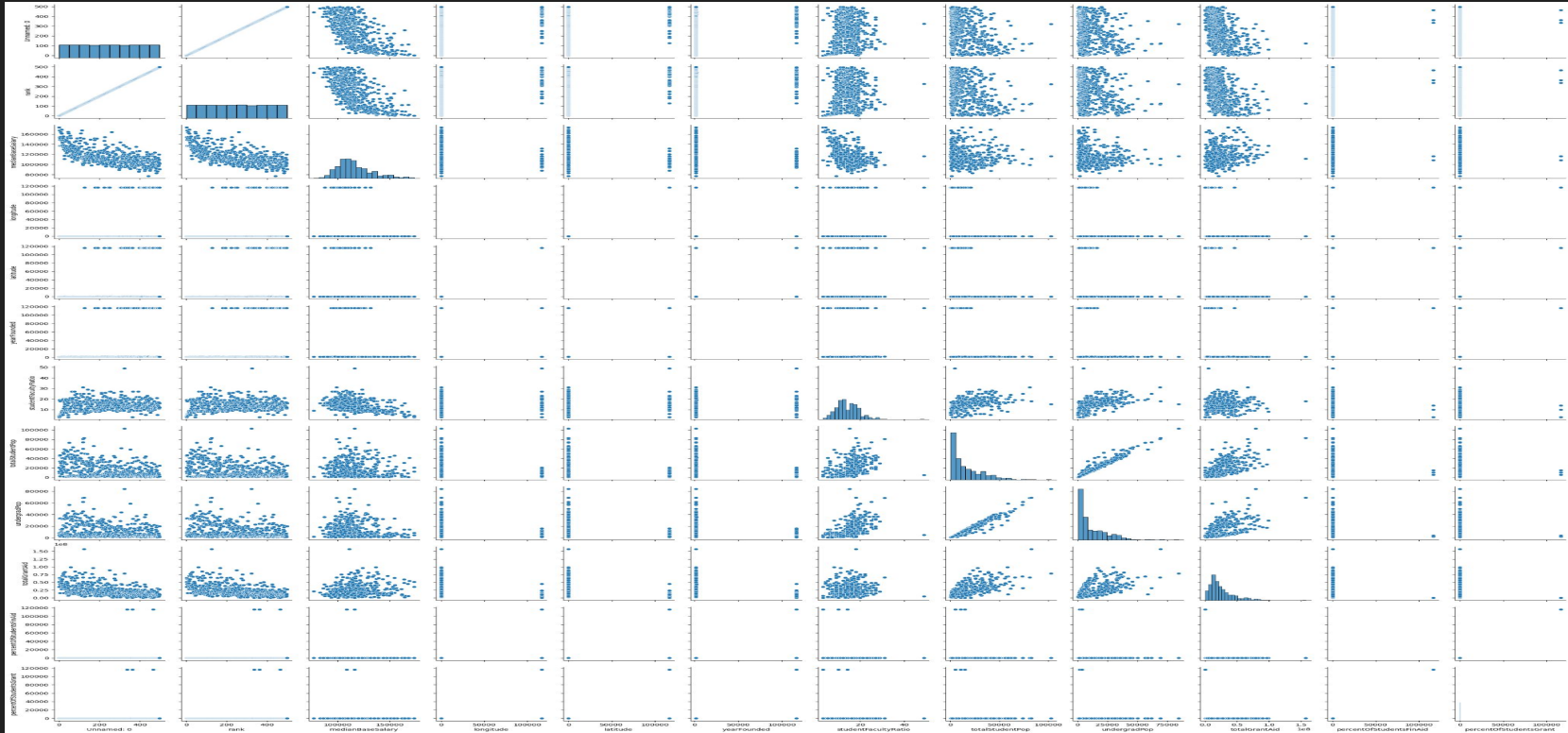Tableau story reflects the findings from EDA after the Cleaning process

Here is the link to my Tableau story:

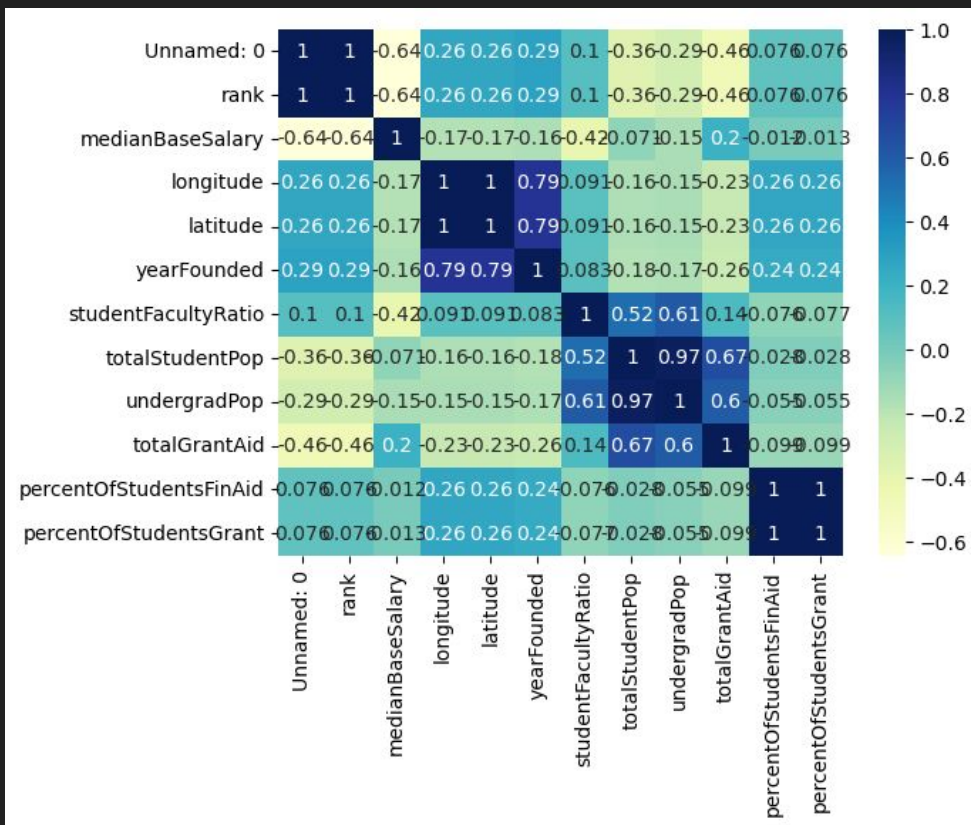https://public.tableau.com/app/profile/shilpa.rao5350/viz/CollegeAnalytics-Story/Story1

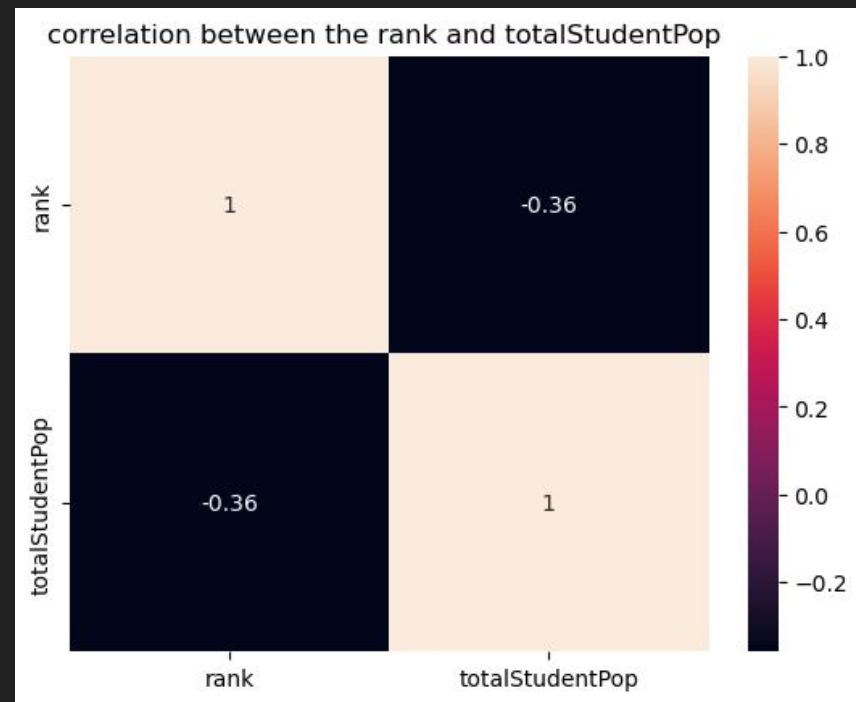Here is a pdf version of the tableau story:

# Checkpoint 5 : Modeling Data

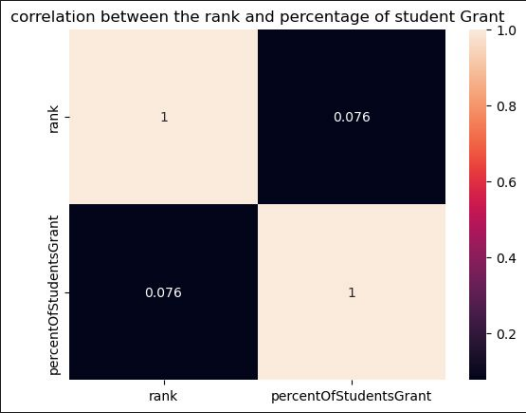correlation - pairplot for all columns
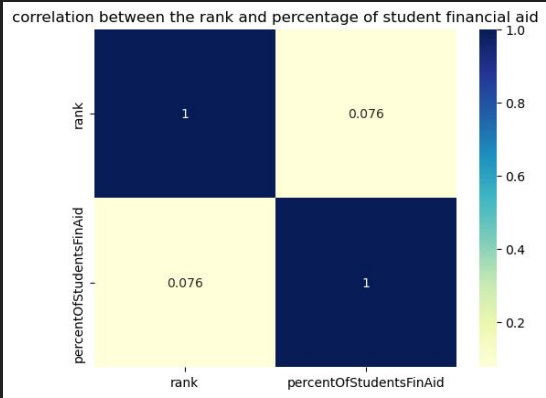
## correlation - heatmap for all columns

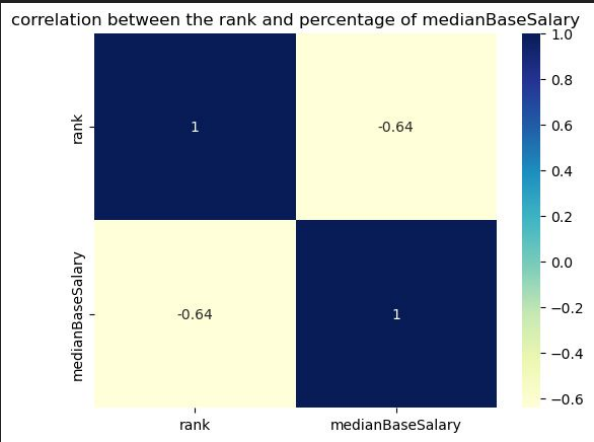correlation between the rank and totalStudentPop

heat map for correlation between the rank of the colleges vs percent Of Students Grant

heat map for correlation between the rank of the colleges vs percentOfStudentsFinAid



correlation between the rank and percentage of student Grant



correlation between the rank and percentage of student financial aid

correlation between the rank of the colleges vs medianBaseSalary



correlation between the rank and percentage of medianBaseSalary

# Conclusion

- ❖ colleges are concentrated in coastal locations, especially northeast and lesser in midwest
- ❖ looks like higher ranked and higher population schools are connected and are higher in number along the coast
- ❖ Higher the rank of the college, higher the median base salary. This means that after graduating from higher ranked colleges, possibility of getting higher salary is more likely
- ❖ lower the rank shows higher percentage of financial aids and higher rank shows lower percentage of financial aid
- ❖ lower the rank shows higher percentage of grants and higher rank shows lower percentage of grants