# Project Report

Name: Akshata More

Topic: College Score Cards

Data Source: https://data.world/xprizeai-edu/college-scorecard/workspace/file?filename=merged_2013_PP.csv  (data.world)

Data Set Link:https://query.data.world/s/55ejgjtervq7xijexg3r74rvdc3xac?dws=00000

Analysis Technique Used: Principal Component Analysis (PCA) + Cluster Analysis

Problem Statement: Using cluster analysis to identify the groups of characteristically similar schools in the College Scorecard dataset.

# 1. INTRODUCTION:

- This data set of "College Scorecard" was primarily released by U.S. Department of Education which is now available on https://data.world/.

- **Motto:** The Department of Education is focused on ensuring that parents, students, and policymakers are able to use its publicly available data to take to take better college decisions.

- Principal Component Analysis or PCA is a statistical technique that is used to reduce dimensionality of data and to obtain suitable principal components, which explains the maximum variance.

- To select the appropriate data from the vast data set to perform analyses like Cluster Analysis, PCA is used as a prerequisite step.

- Cluster Analysis is a statistical technique used to group similar objects into categories or clusters.

# 2. Steps Involved:

**a.** I chose the College Scorecard data set of the year 2013 in csv format



**b.** Loaded the csv into Google Colab and Imported Libraries.

**c.** Computed description of the data set.

**d.** Data Preprocessing: Data Cleaning and Assessing.

    i. I have selected the following variables, and created a new dataframe , "df_clean":

UNITID:

UNITID which is the unique identification number

INSTNM:

institution's name


CITY, STABBR:


NUMBRANCH:

the number of branch campuses


HIGHDEG:

Highest award


PREDDEG:

Predominant undergraduate award


CONTROL:

institution's governance structure is public, private nonprofit, or private for-profit.


DISTANCEONLY:

distance education-


TUITFTE:

The net tuition revenue per full-time equivalent (FTE)


AVGFACSAL:

average faculty salary


ADM_RATE_ALL:

Fall admissions rate, represents the admissions rate across all campuses


SATVR25, SATVR75, SATMT25, SATMT75, ACTCM25, ACTCM75:

The files include the 25th and 75th percentiles of SAT


UGDS:

the number of degree/certificate-seeking undergraduates


PCTFLOAN:

the share of undergraduate students who received federal loans in a given year.
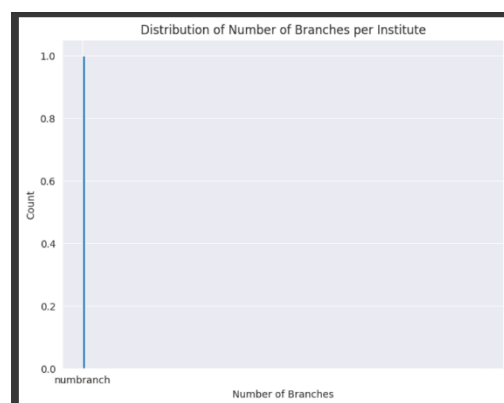
CDR3:

The three-year cohort default rate (CDR3) represents a snapshot in time.

e. Checked for null values and cleaned the na values.
f. Many Categorical variables have been encoded as Numeric, as evident from the data set. I converted those back for following variables:
   i. STABBR: converted column to categorical type
   ii. HIGHDEG and PREDDEG: mapped the numerical value to categorical value
   iii. CONTROLD
   iv. DISTANCEONLY
   v. UGDS
g. Created some graphs/charts for data visualization to answer questions like:
   i. How are institutes spread out across the US?
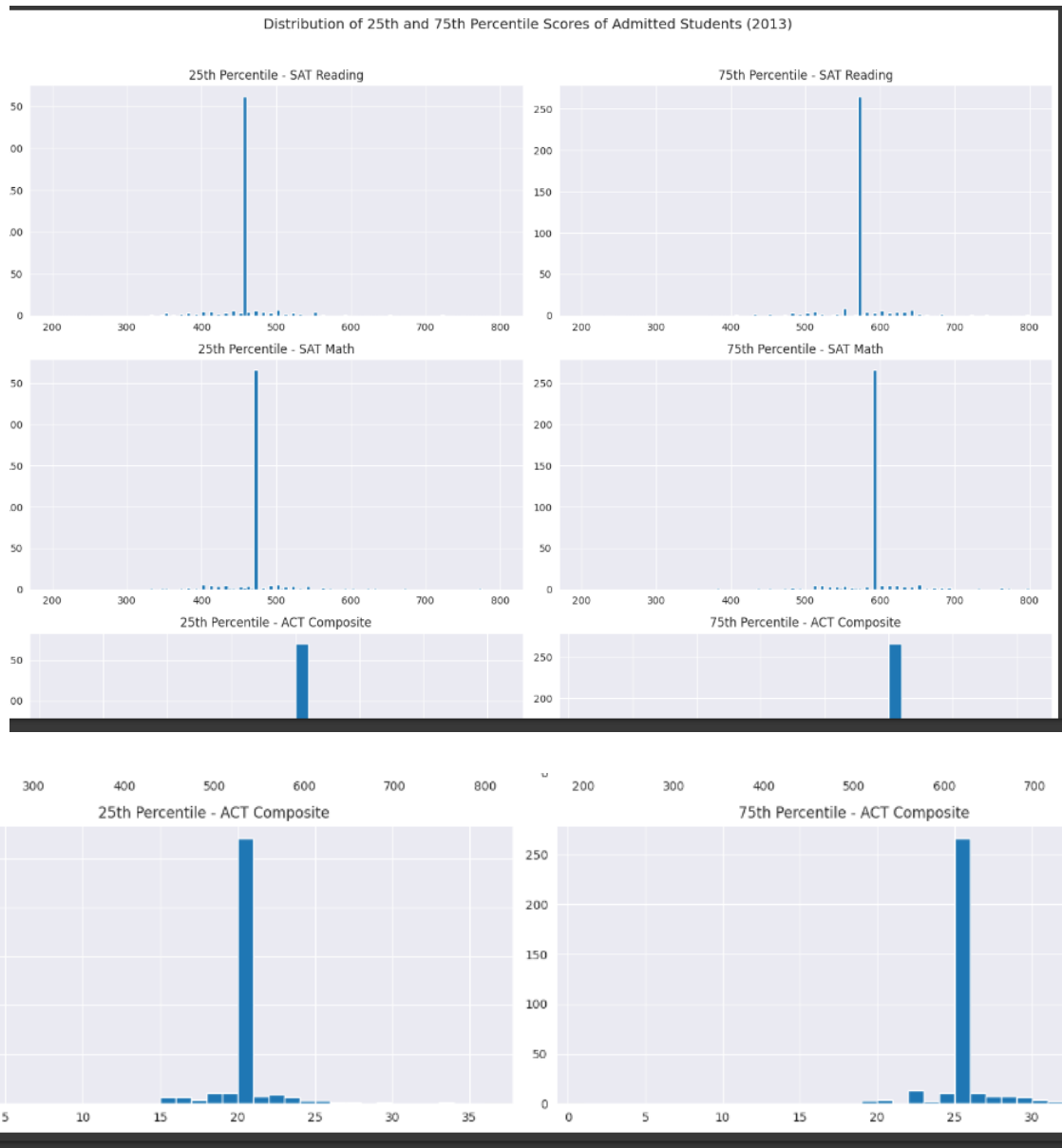


State-wise Distribution of Institutes

1. The state with code 140 has by far the greatest number of institutes.
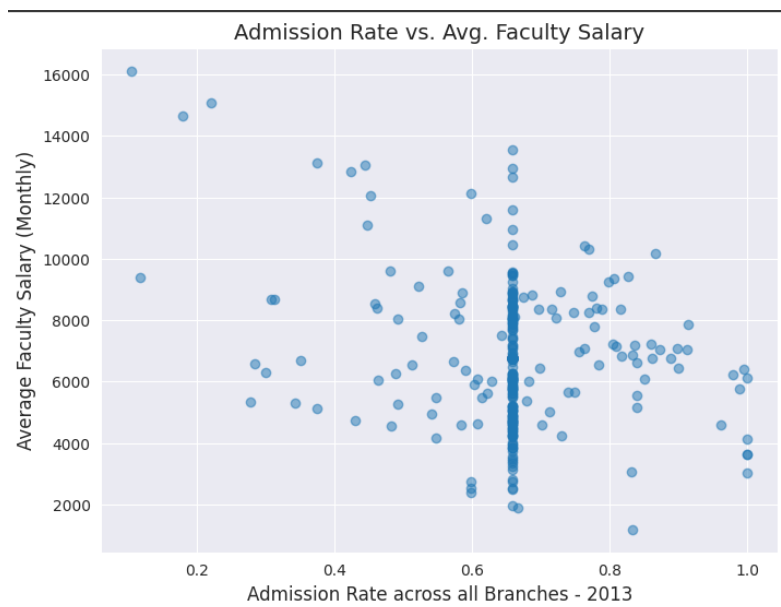2. The states with code 55 to 60 have closely competing numbers.

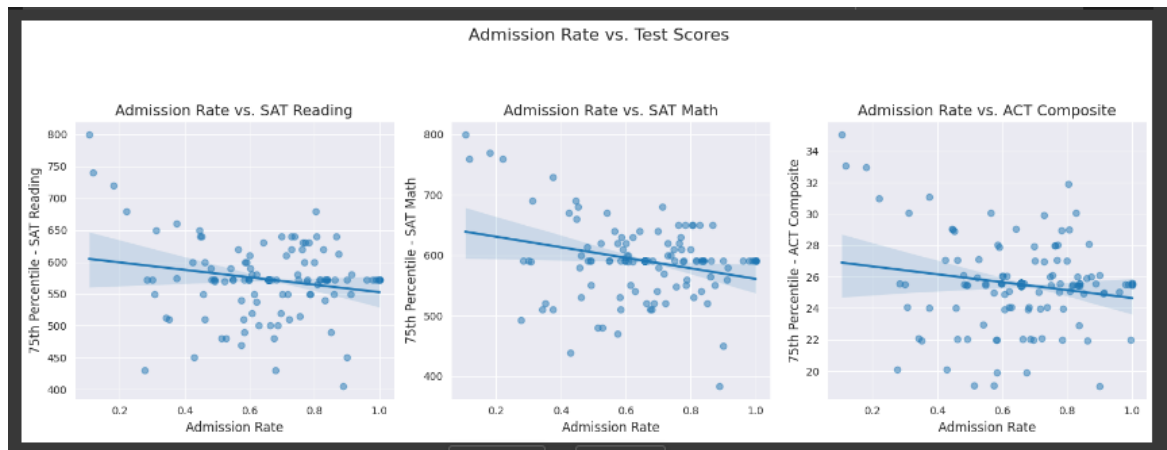   ii. What is the median number of branches per Institute?



Distribution of Number of Branches per Institute

iii. What makes for a good or bad SAT or ACT score? : Distribution of 25th and 75th percentile Scores of Admitted Students



Distribution of 25th and 75th Percentile Scores of Admitted Students (2013)



iv. How does admission rate of an institute affect the average faculty salary?
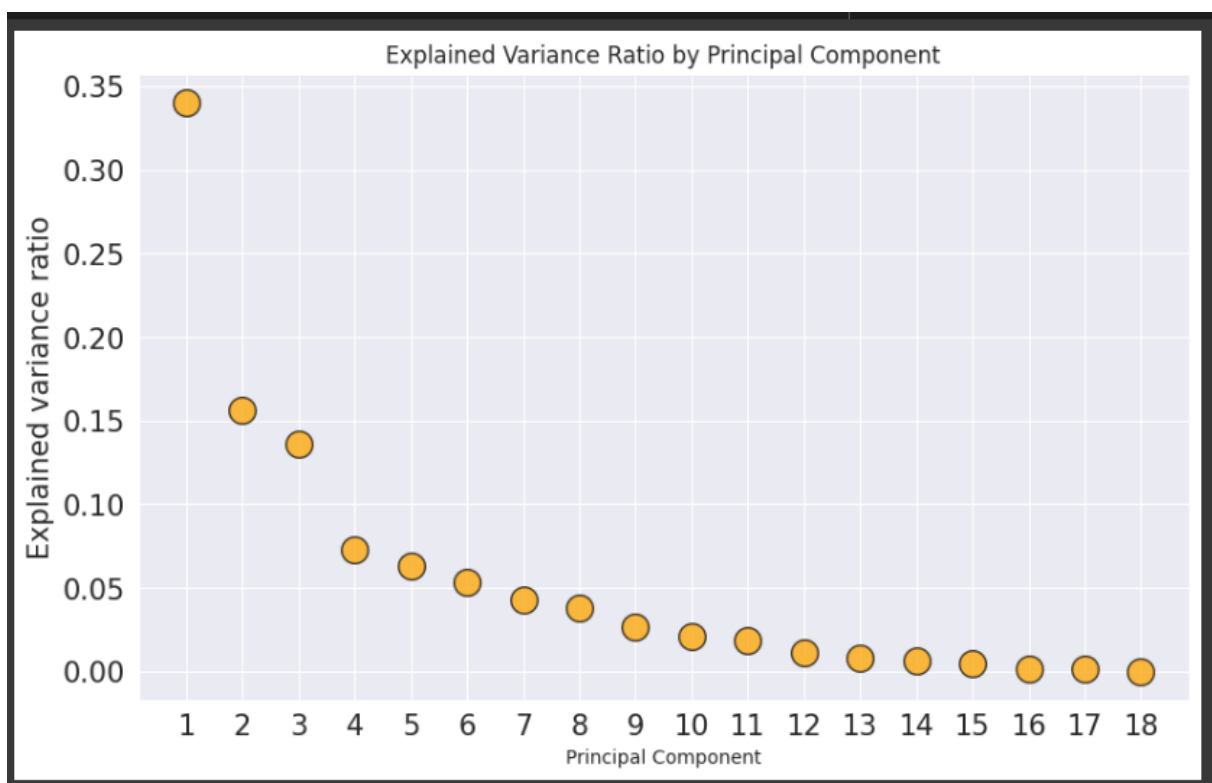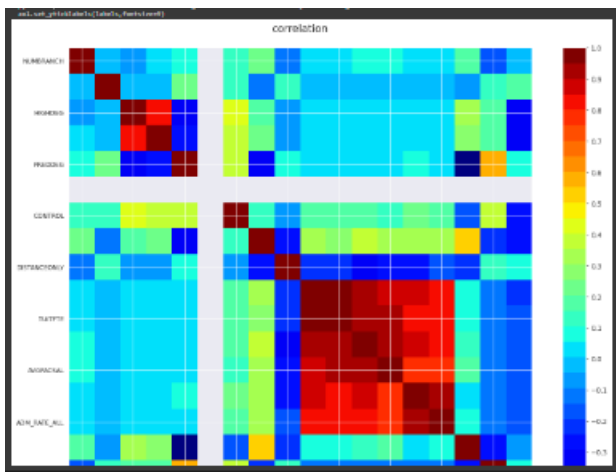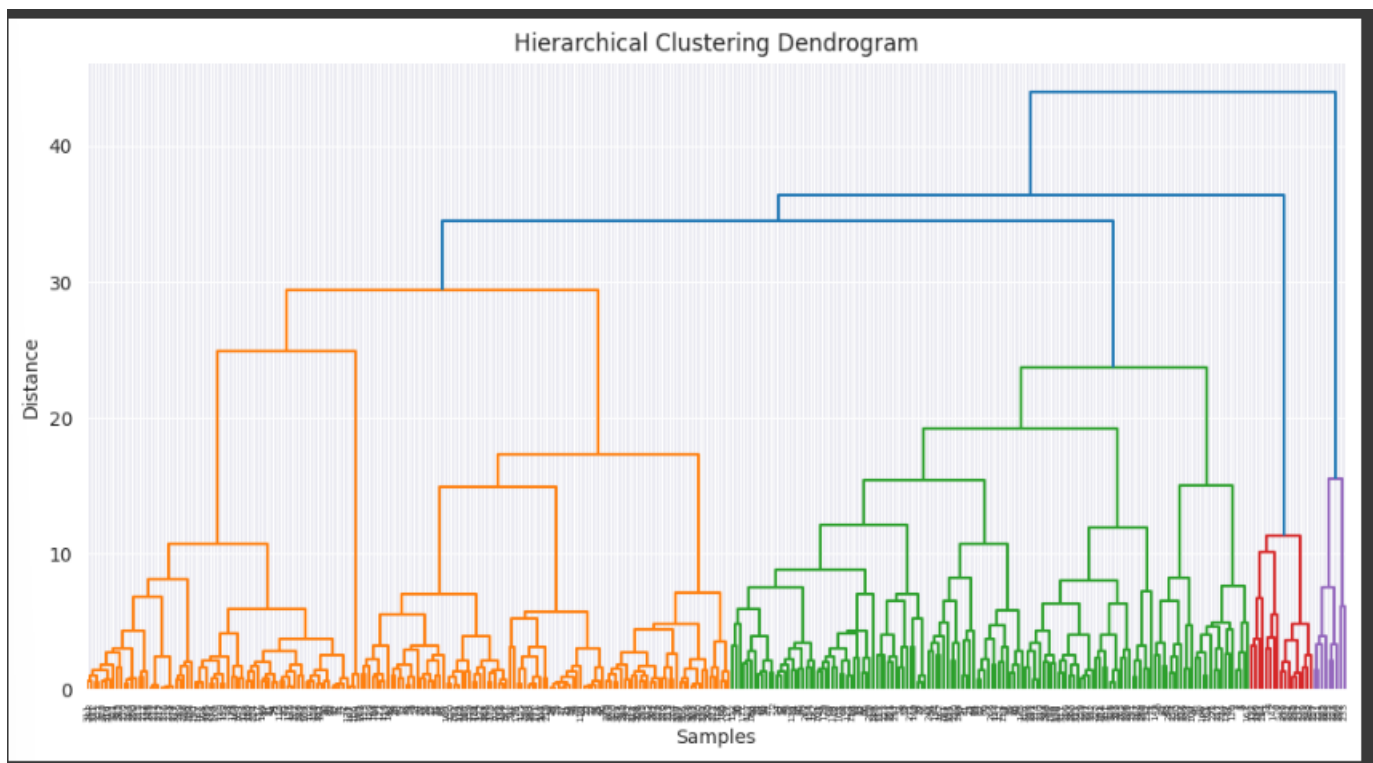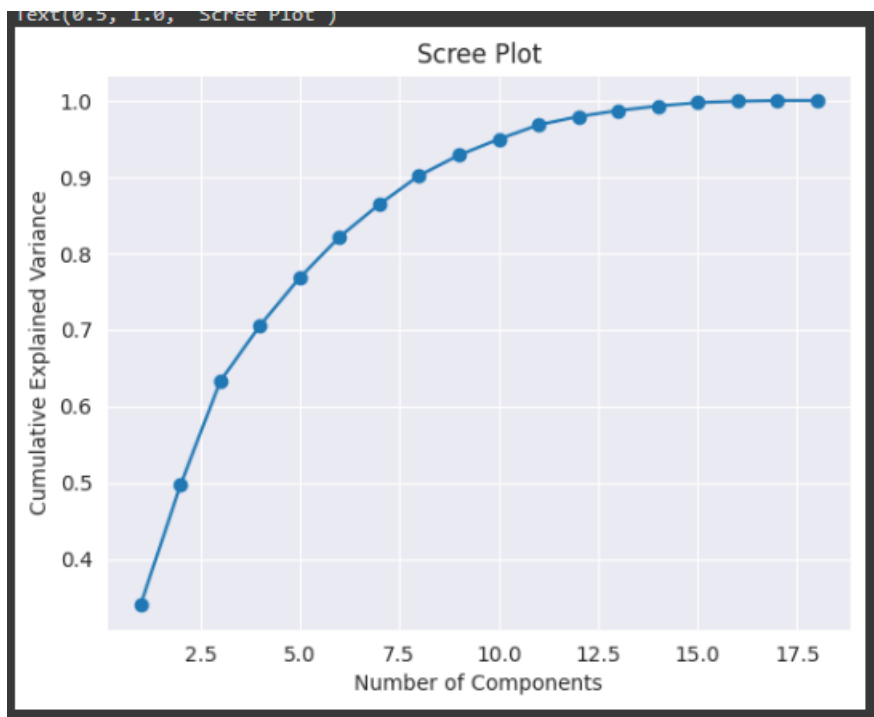


Admission Rate vs. Avg. Faculty Salary

v. One might expect that institutes having a high 75th percentile of test scores for admitted students, would have low admission rates. Does the data support this assumption?



h. Plotted covariance matrix – to check independency

i. Performed PCA

j. Performed Cluster Analysis

**3. Results :**

Scree Plot



Hierarchical Clustering Dendrogram

4. Conclusion:

Hence, here using PCA and Clustering you can draw conclusions as to which college is better depending on many varibales chosen before like act scores, act scores, distance only, etc. According to the cumulative explained variance we retain 18 components and reject other for clustering.