

SUMMER BOOTCAMP 2024

Data Science Project KESHAV BISHT Education Post 12th Data Analysis

S No.	Topic	Page No.
1	Cover Page	1
2	Index	2
3	List of Tables	3
4	List of Figures	3
5	Problem Statement	4
6	Data Exploration	5-7
7	Basic EDA	8-16
8	Application and Enrollment Analysis	17-21
9	Academic Excellence	22-24
10	Student Demographics	25-27
11	Cost of Spendings	28-31
12	Faculty Qualification	32-34
13	Student-Faculty Interaction	35
14	Alumni Engagement	36-37
15	Graduation Rates	38-40
16	Overall Insights	41-42

PROBLEM STATEMENT / OBJECTIVE

The objective of this analysis is to gain insights into the characteristics of colleges and answer key questions related to the educational landscape. By understanding the data, we aim to inform strategies for improving the quality of education and enhancing the overall college experience. The analysis will provide valuable insights and recommendations for stakeholders in the education sector.

IMPORTING THE LIBRARIES

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

IMPORTING THE DATA SET

```
In [2]: df = pd.read_csv("Users/keshavbisht/Downloads/1-Education_Post_12th_Standard_New.csv")
```

BASIC EXPLORATION

1. FIRST FIVE ROWS

```
In [15]: df.head()

Out[15]:
```

	Names	Apps	Accept	Enroll	Top10perc	Top25perc	F.Undergrad	P.Undergrad	Outstate	Room.Board	Books	Personal	PhD	Terminal	S.F.Ratio	perc.alumni	Expend	Grad.Rate
0	Abilene Christian University	1660.0	1232	721.0	23.0	52	2885	537	7440	3300	450	2200.0	70	78	18.1	12	7041	60
1	Adelphi University	2186.0	1924	512.0	16.0	29	2683	1227	12280	6450	750	1500.0	29	30	?	16	10527	56
2	Adrian College	1428.0	1097	336.0	22.0	50	1036	99	11250	3750	400	1185.0	63	66	12.9	30	8735	54
3	Agnes Scott College	417.0	349	NaN	60.0	89	510	63	12960	5450	450	875.0	92	97	7.7	37	19016	59
4	Alaska Pacific University	193.0	146	55.0	16.0	44	249	869	7560	4120	800	1500.0	76	72	11.9	2	10922	15

2. INFO

```
In [36]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 777 entries, 0 to 776
Data columns (total 18 columns):
 #   Column                Non-Null Count  Dtype  
---  --
 0   Names                 777 non-null    object  
 1   Apps                 775 non-null    float64 
 2   Accept              777 non-null    int64   
 3   Enroll              775 non-null    float64 
 4   Top10perc           773 non-null    float64 
 5   Top25perc           777 non-null    int64   
 6   F.Undergrad         777 non-null    int64   
 7   P.Undergrad         777 non-null    int64   
 8   Outstate            777 non-null    int64   
 9   Room.Board          777 non-null    int64   
10   Books               777 non-null    int64   
11   Personal            774 non-null    float64 
12   PhD                777 non-null    int64   
13   Terminal            777 non-null    int64   
14   S.F.Ratio           777 non-null    object  
15   perc.alumni         777 non-null    int64   
16   Expend              777 non-null    int64   
17   Grad.Rate           777 non-null    int64   
dtypes: float64(4), int64(12), object(2)
memory usage: 109.4+ KB
```

3. LAST 5 ROWS

```
In [37]: df.tail()

Out[37]:
```

	Names	Apps	Accept	Enroll	Top10perc	Top25perc	F.Undergrad	P.Undergrad	Outstate	Room.Board	Books	Personal	PhD	Terminal	S.F.Ratio	perc.alumni	Expend	Grad.Rate
772	Worcester State College	2197.0	1515	543.0	4.0	26	2689	2029	6797	3900	500	1200.0	60	70	21	14	4469	40
773	Xavier University	1959.0	1805	695.0	24.0	47	2849	1107	11520	4960	600	1250.0	73	75	13.3	31	9189	83
774	Xavier University of Louisiana	2097.0	1915	695.0	34.0	61	2793	166	6900	4200	617	791.0	67	75	14.4	20	8323	49
775	Yale University	10705.0	2453	1317.0	95.0	99	5217	83	19840	6510	630	2115.0	96	96	5.8	49	40386	99
776	York College of Pennsylvania	2989.0	1855	691.0	28.0	63	2988	1726	4990	3560	500	1250.0	75	75	18.1	28	4509	99

4. SHAPE

```
In [39]: df.shape
Out[39]: (777, 18)
```

5. STATISTICAL SUMMARY

```
In [20]: df.describe()

Out[20]:
```

	Apps	Accept	Enroll	Top10perc	Top25perc	F.Undergrad	P.Undergrad	Outstate	Room.Board	Books	Personal	PhD	Terminal	S.F.Ratio	perc.alumni	Expend	Grad.Rate
count	775.000000	777.000000	775.000000	773.000000	777.000000	777.000000	777.000000	777.000000	777.000000	777.000000	774.000000	777.000000	777.000000	777.000000	777.000000	777.000000	777.000000
mean	3007.592258	2018.804376	780.961290	27.620957	55.796654	3699.907336	855.298584	10440.669241	4357.526384	547.875161	1601.507752	72.660232	79.702703	22.743987	966	14.722359	12.391901
std	3873.414660	2451.113971	930.077779	17.645470	19.804778	4850.420531	1522.431887	4023.016484	1096.696416	167.426237	7369.594038	16.328195	14.722359	1.891901	512	1.000000	0.000000
min	81.000000	72.000000	35.000000	1.000000	9.000000	139.000000	1.000000	2340.000000	1780.000000	0.000000	50.000000	62.000000	71.000000	13.000000	675	0.000000	0.000000
25%	778.000000	604.000000	242.500000	15.000000	41.000000	992.000000	95.000000	7320.000000	3597.000000	465.000000	855.000000	62.000000	71.000000	13.000000	675	0.000000	0.000000
50%	1561.000000	1110.000000	434.000000	23.000000	54.000000	1707.000000	353.000000	9990.000000	4200.000000	500.000000	1200.000000	75.000000	82.000000	21.000000	837	0.000000	0.000000
75%	3635.000000	2424.000000	902.500000	35.000000	69.000000	4005.000000	967.000000	12925.000000	5050.000000	600.000000	1687.500000	85.000000	92.000000	31.000000	1083	0.000000	0.000000
max	48094.000000	26330.000000	6392.000000	96.000000	100.000000	31643.000000	21836.000000	21700.000000	8124.000000	2340.000000	205500.000000	103.000000	100.000000	64.000000	5623	0.000000	0.000000

6. NULL VALUES

```
In [21]: df.isnull()

Out[21]:
```

	Names	Apps	Accept	Enroll	Top10perc	Top25perc	F.Undergrad	P.Undergrad	Outstate	Room.Board	Books	Personal	PhD	Terminal	S.F.Ratio	perc.alumni	Expend	Grad.Rate
0	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
3	False	False	True	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
...
772	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
773	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
774	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
775	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
776	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False
777	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False	False

777 rows x 18 columns

7. DUPLICATE VALUES

```
In [22]: df.duplicated()

Out[22]:
```

0	False
1	False
2	False
3	False
4	False
...	...
772	False
773	False
774	False
775	False
776	False
Length:	777, dtype: bool

8. OUTLIERS AND THEIR AUTHICITY

```
In [23]: plt.figure(figsize=(10,3))
sns.boxplot(data = df)

Out[23]:
```

<Axes: >

9. ANOMALIES OR WRONG ENTRY

```
In [5]: df.describe(include='all')

Out[5]:
```

	Names	Apps	Accept	Enroll	Top10perc	Top25perc	F.Undergrad	P.Undergrad	Outstate	Room.Board	Books	Personal	PhD	Terminal	S.F.Ra
count	777	775.000000	777.000000	775.000000	773.000000	777.000000	777.000000	777.000000	777.000000	777.000000	777.000000	774.000000	777.000000	777.000000	777.000000
unique	777	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
top	Abilene Christian University	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
freq	1	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
mean	NaN	3007.592258	2018.804376	780.961290	27.620957	55.796654	3699.907336	855.298584	10440.669241	4357.526384	547.875161	1601.507752	72.660232	79.702703	966
std	NaN	3873.414660	2451.113971	930.077779	17.645470	19.804778	4850.420531	1522.431887	4023.016484	1096.696416	167.426237	7369.594038	16.328195	14.722359	512
min	NaN	81.000000	72.000000	35.000000	1.000000	9.000000	139.000000	1.000000	2340.000000	1780.000000	0.000000	50.000000	62.000000	71.000000	675
25%	NaN	778.000000	604.000000	242.500000	15.000000	41.000000	992.000000	95.000000	7320.000000	3597.000000	465.000000	855.000000	62.000000	71.000000	675
50%	NaN	1561.000000	1110.000000	434.000000	23.000000	54.000000	1707.000000	353.000000	9990.000000	4200.000000	500.000000	1200.000000	75.000000	82.000000	837
75%	NaN	3635.000000	2424.000000	902.500000	35.000000	69.000000	4005.000000	967.000000	12925.000000	5050.000000	600.000000	1687.500000	85.000000	92.000000	1083
max	NaN	48094.000000	26330.000000	6392.000000	96.000000	100.000000	31643.000000	21836.000000	21700.000000	8124.000000	2340.000000	205500.000000	103.000000	100.000000	5623

APPLICATION AND ENROLLMENT ANALYSIS

1. WHAT IS THE AVERAGE NUMBER OF APPLICATIONS RECEIVED BY COLLEGES?

```
In [7]: avg_apps = df['Apps'].mean()
avg_apps

Out[7]: 3007.592258064516
```

2. WHAT PERCENTAGE OF APPLICATIONS ARE ACCEPTED ON AVERAGE ACROSS ALL COLLEGES?

```
In [8]: acceptance_rate = (df['Accept'] / df['Apps']).mean() * 100
acceptance_rate

Out[8]: 74.662674988846
```

3. What is the average enrollment rate (number of students enrolled divided by number of applications accepted)?

```
In [9]: enrollment_rate = (df['Enroll'] / df['Accept']).mean() * 100
enrollment_rate

Out[9]: 41.19383221625694
```

4. Which college has the highest number of applications received?

```
In [31]: highest_apps_college = df[df['Apps'] == df['Apps'].max()][['Names']].values[0]
highest_apps_college

Out[31]: 'Rutgers at New Brunswick'
```

5. Average percentage of new students from the top 10% of their high school class

```
In [32]: avg_top10perc = df['Top10perc'].mean()
avg_top10perc

Out[32]: 27.620957395184995
```

ACADEMIC EXCELLENCE

1. Average percentage of new students from the top 10% of their high school class

```
In [35]: avg_top10perc = df['Top10perc'].mean()
avg_top10perc

Out[35]: 27.620957395184995
```

2. Average percentage of new students from the top 25% of their high school class

```
In [36]: avg_top25perc = df['Top25perc'].mean()
avg_top25perc

Out[36]: 55.796657966538
```

3. Correlation between the percentage of students from the top 10% and the top 25% of their high school class

```
In [37]: correlation_top10_top25 = df['Top10perc'].corr(df['Top25perc'])
correlation_top10_top25

Out[37]: 0.8919819734345558
```

STUDENT DEMOGRAPHICS

1. What is the average number of full-time undergraduate students per college?

```
In [38]: avg_full_time_undergrad = df['F.Undergrad'].mean()
avg_full_time_undergrad

Out[38]: 3699.90735907336
```

2. Average number of part-time undergraduate students per college

```
In [39]: avg_part_time_undergrad = df['P.Undergrad'].mean()
avg_part_time_undergrad

Out[39]: 855.2985842985343
```

3. College with the highest number of out-of-state students

```
In [20]: highest_outstate_college = df[df['Outstate'] == df['Outstate'].max()][['Names']].values[0]
highest_outstate_college

Out[20]: 'Bennington College'
```

COST OF SPENDING

1. What is the average cost of room and board across all colleges?

```
In [21]: avg_room_board = df['Room.Board'].mean()
avg_room_board

Out[21]: 4357.526384526383
```

2. What is the average estimated book cost for a student?

```
In [22]: avg_book_cost = df['Books'].mean()
avg_book_cost

Out[22]: 547.8751608751609
```

3. What is the average estimated personal spending for a student?