

Center for Undergraduate Research Opportunities

Project 1: Milestone Report

This project depicts an end-to-end Data Science workflow aimed at exploring the student retention rate (whether or not a student graduates from a higher education institution). Attrition results in a severe loss of resources by society, by students, and by colleges that spend to provide programs and services to help retain and graduate students. When a student leaves college prematurely, any debt incurred must be repaid, despite the inability to graduate. In addition, the students may become turned off to the educational system in general, unexpected to benefit from education opportunities that may have aided in job attainment or advancement efforts. Administrators at historically black colleges and universities (HBCUs) report low income as a main reason for student departure among such institutions.

With the escalating cost of higher education tuition and fees, the federal and state government's have increased expenditures to higher education through financial aid programs, retention has become a salient issue. Although large public state universities offer generous aid for students, there are a number of environmental issues that can affect whether or not a university retains their students. My client would be any federal administration that demonstrates an ongoing commitment to advance diversity and opportunity in higher education, such as the U.S. Department of Education. This analysis would help them look at other factors that may contribute to student retention or departures. Is there another way the education industry can retain students? Are there other factors that play into student retention? To what extent does income really matter when it comes to retaining undergraduates? Do certain institutions offer formal or informal student experiences that contribute to academic integration?

I believe these are important questions to ask since the rapidly growing post-secondary education industry relies on young people essentially acting as patrons of the institutions. It is important to find ways to maintain high student retention and graduation rates.

```
# Native
import datetime
import matplotlib.pyplot as plt
import os
import sys

# Non-Native
import numpy as np
import pandas as pd
import seaborn as sns
sns.set(style='whitegrid')
```

Data Acquisition

The datasets used in this analysis is derived from the National Center for Education Statistics' (NCES). The dataset is available to download from the Integrated Postsecondary Education Data System. Data reporting is mandatory for institutions that participate in or are applicants for participation in any federal student financial aid programs authorized by Title IV of the Higher Education Act of 1965.

I downloaded the data in three separate CSV files due to the nature of the data contained. The first dataframe reports institutional characteristics, the second reports 12 month enrollment counts, and the last contains statistics regarding graduation completion counts.

These 3 datasets contain about()variables that may contribute to student success including under: financial aid status, graduation rates, and demographic breakdowns of the undergraduate cohort, and student-faculty interaction.

This dataset contains approximately 5900 higher education institutions. I read the csv with `pd.read_csv()` into a pandas data frame.

```
df = pd.read_csv("CSV_142024-95.csv")
df2 = pd.read_csv("CSV_142024-46.csv")
df3 = pd.read_csv("CSV_142024-393.csv")
```

Data Cleaning & Pre-Processing

The dataset was already clean with few missing values. Given the dataframes containing the student demographic and achievement data, clean them up by:

- Converting the HBCU Status into a binary variable (0/1 for No and Yes, respectively).
- Convert variable names into shorter, and useable names for each data item.
- Handle any missing/ NaN values
- Ensure reported undergraduate counts and percentages are equivalent to total reported statistics for the 2021-2022 cohort.

```
df.rename(columns= {"institution name": "institution",
                    "HD2022.OPE Title IV eligibility indicator code" : "TitleIV_Indicator", ##drop th
                    "IC2022.Full time first-time degree/certificate-seeking undergraduate students en
                    "SFA2122.Percent of undergraduate students awarded federal, state, local, institu
                    "DRVGR2022.Graduation rate, total cohort" : "TotalGradRate",
                    "DRVGR2022.Graduation rate, American Indian or Alaska Native" : "AmericanIndianG
                    "DRVGR2022.Graduation rate, Asian" : "AsianGradRate",
                    "DRVGR2022.Graduation rate, Native Hawaiian or Other Pacific Islander" : "Pacific
                    "DRVGR2022.Graduation rate, Black, non-Hispanic" : "BlackGradRate",
                    "DRVGR2022.Graduation rate, Hispanic" : "HispanicGradRate",
```

```

        "DRVGR2022.Graduation rate, White, non-Hispanic" : "WhiteGradRate",
        "DRVGR2022.Graduation rate, two or more races" : "MultiRacialGradRate",
        "DRVGR2022.Graduation rate, Race/ethnicity unknown" : "UnreportedRaceGradRate",
        "DRVGR2022.Transfer-out rate - Bachelor cohort" : "TransferOutRate",
        "HD2022.Historically Black College or University" : "HBCUStatus",
        "HD2022.Postsecondary and Title IV institution indicator" : "PostSecondaryIndic",
        "DRVEF122022.Undergraduate 12-month unduplicated headcount" : "UndergradCount",
        "DRVEF122022.Percent of undergraduate 12-month unduplicated headcount that are Am
"DRVEF122022.Percent of undergraduate 12-month unduplicated headcount that are Asian" : "Asi
"DRVEF122022.Percent of undergraduate 12-month unduplicated headcount that are Black or Afr
"DRVEF122022.Percent of undergraduate 12-month unduplicated headcount that are Hispanic/Lati
"DRVEF122022.Percent of undergraduate 12-month unduplicated headcount that are Native Hawai
"DRVEF122022.Percent of undergraduate 12-month unduplicated headcount that are White" : "Wh
"DRVEF122022.Percent of undergraduate 12-month unduplicated headcount that are two or more r
"DRVEF122022.Percent of undergraduate 12-month unduplicated headcount that are race/ethnicit

df.drop(columns=["TitleIV_Indicator", "PostSecondaryIndic", "DegreeSeeking"], inplace=True)

df["TransferOutCount"] = df["UndergradCount"] * df["TransferOutRate"]
df["TotalGrantedAid"] = df["PercentGrantedAid"] * df["UndergradCount"]

df["unitid"].is_unique

True

df["HBCUStatus"].replace(to_replace = dict(Yes = 1, No = 0), inplace=True)

df.dropna(inplace=True)

df.isnull().sum()

unitid                0
institution            0
year                  0
PercentGrantedAid     0
TotalGradRate         0
AmericanIndianGradRate 0
AsianGradRate         0
PacificIslanderGradRate 0
BlackGradRate         0
HispanicGradRate      0
WhiteGradRate         0
MultiRacialGradRate   0
UnreportedRaceGradRate 0
TransferOutRate       0
HBCUStatus            0
UndergradCount        0
AmericanIndianUndergradPercent 0
AsianUndergradPercent  0
BlackUndergradPercent  0

```

```

HispanicUndergradPercent      0
PacificIslanderUndergradPercent  0
WhiteUndergradPercent         0
MultiRacialUndergradPercent   0
UnreportedRaceUndergradPercent 0
TransferOutCount              0
TotalGrantedAid               0
dtype: int64

df2.rename(columns= {"EFFY2022.Level and degree/certificate-seeking status of student" : "StudentLevel",
                     "EFFY2022.Grand total" : "TotalUndergradCount",
                     "EFFY2022.American Indian or Alaska Native total" : "AmericanIndianUndergradCount",
                     "EFFY2022.Asian total" : "AsianUndergradCount",
                     "EFFY2022.Black or African American total" : "BlackUndergradCount",
                     "EFFY2022.Hispanic or Latino total" : "HispanicUndergradCount",
                     "EFFY2022.Native Hawaiian or Other Pacific Islander total" : "PacificIslanderUndergradCount",
                     "EFFY2022.White total" : "WhiteUndergradCount",
                     "EFFY2022.Two or more races total" : "MultiRacialUndergradCount",
                     "EFFY2022.Race/ethnicity unknown total" : "UnreportedRaceUndergradCount"},
            inplace=True)

df2.drop(columns=["StudentLevel", "IDX_E12"], inplace=True)

df2.isnull().sum()

unitid      0
institution name  0
year        0
TotalUndergradCount  0
AmericanIndianUndergradCount  0
AsianUndergradCount  0
BlackUndergradCount  0
HispanicUndergradCount  0
PacificIslanderUndergradCount  0
WhiteUndergradCount  0
MultiRacialUndergradCount  0
UnreportedRaceUndergradCount  0
dtype: int64

df3.rename(columns= {"C2022_C.Award Level code" : "DegreeType",
                    "C2022_C.Grand total" : "TotalStudentsGraduating",
                    "C2022_C.American Indian or Alaska Native total" : "AmericanIndianGraduatingCount",
                    "C2022_C.Asian total" : "AsianGraduatingCount",
                    "C2022_C.Black or African American total" : "BlackGraduatingCount",
                    "C2022_C.Hispanic or Latino total" : "HispanicGraduatingCount",
                    "C2022_C.Native Hawaiian or Other Pacific Islander total" : "PacificIslanderGraduatingCount",
                    "C2022_C.White total" : "WhiteGraduatingCount",
                    "C2022_C.Two or more races total" : "MultiRacialGraduatingCount",
                    "C2022_C.Race/ethnicity unknown total" : "UnreportedRaceGraduatingCount"},
            inplace=True)

```

```

, inplace=True)

df3.drop(columns = ["DegreeType", "IDX_C"], inplace=True)

df3.isnull().sum()

unitid                                0
institution name                      0
year                                  0
TotalStudentsGraduating              0
AmericanIndianGraduatingCount        0
AsianGraduatingCount                 0
BlackGraduatingCount                 0
HispanicGraduatingCount              0
PacificIslanderGraduatingCount       0
WhiteGraduatingCount                 0
MultiRacialGraduatingCount           0
UnreportedRaceGraduatingCount        0
dtype: int64

```

One discrepancy found during the data cleaning process is the reported total count of undergraduate students. When adding up the counts of undergraduate demographics, I found that the reported count did not match the breakdown of the student demographics. - For example, for Alabama A&M, the total undergraduate count was reported as 5663. When adding up the race demographic counts, it totaled as 5610. This is likely due to multiple categorizations of race, and the addition of Unreported or Multi Racial categories for students that identify with two or more, or chose to leave the field unreported. - For this, I created a column in the dataframe to calculate the Actual Undergraduate Count

```
df2['ActualUndergrad_Count'] = df2["AmericanIndianUndergradCount"]+ df2["AsianUndergradCount"]
```

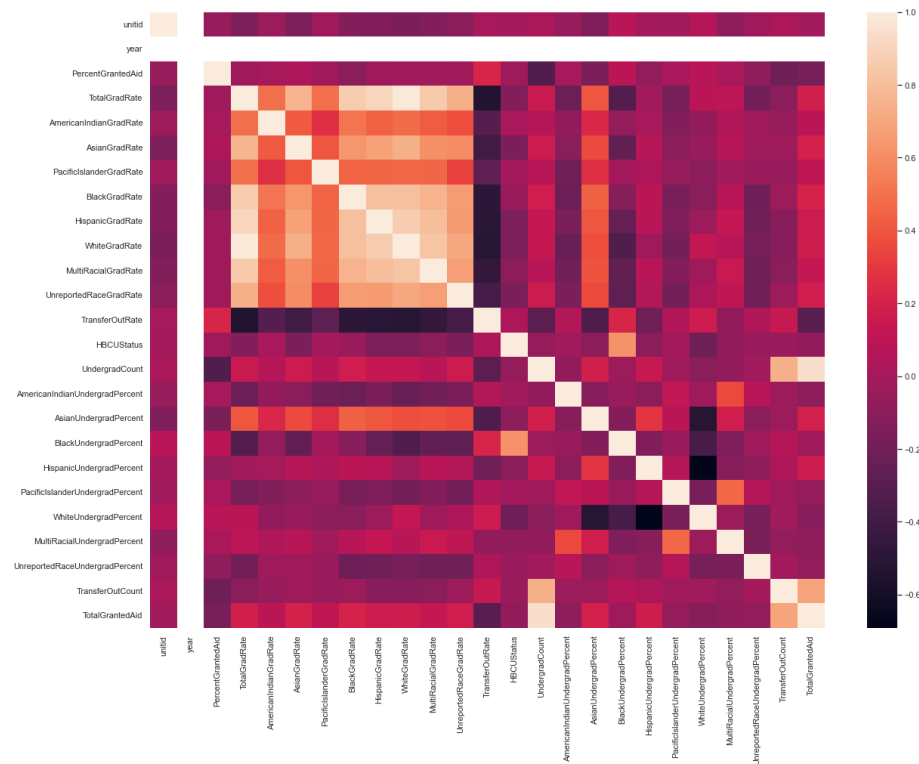
Exploratory Data Analysis

```

corr1 = df.corr()
plt.subplots(figsize=(20,15))
sns.heatmap(corr1,
            xticklabels=corr1.columns.values,
            yticklabels=corr1.columns.values)

<AxesSubplot:>

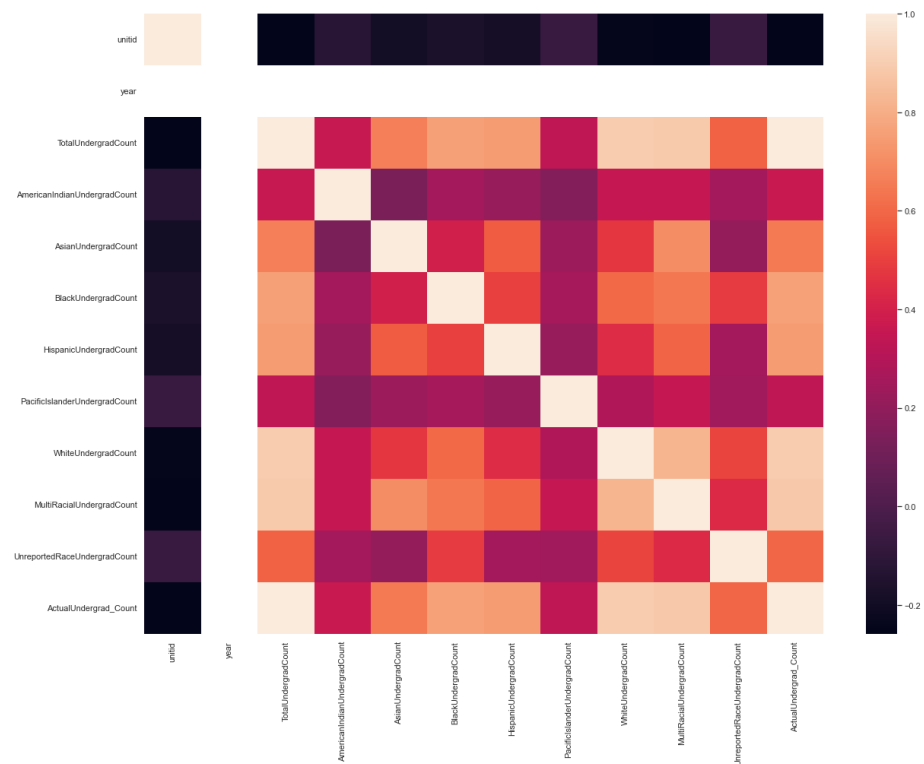
```



I created a correlation plot to see if any of the institutional characteristic variables were correlated with each other. I found that very few variables were correlated. This correlation plot shows that each demographic's graduation rate and undergraduate percent are all correlated with each other. The total graduation rate is also correlated with the white grad rate, multi-racial grad rate, hispanic grad rate, and black grad rate.

```
corr2 = df2.corr()
plt.subplots(figsize=(20,15))
sns.heatmap(corr2,
             xticklabels=corr2.columns.values,
             yticklabels=corr2.columns.values)
```

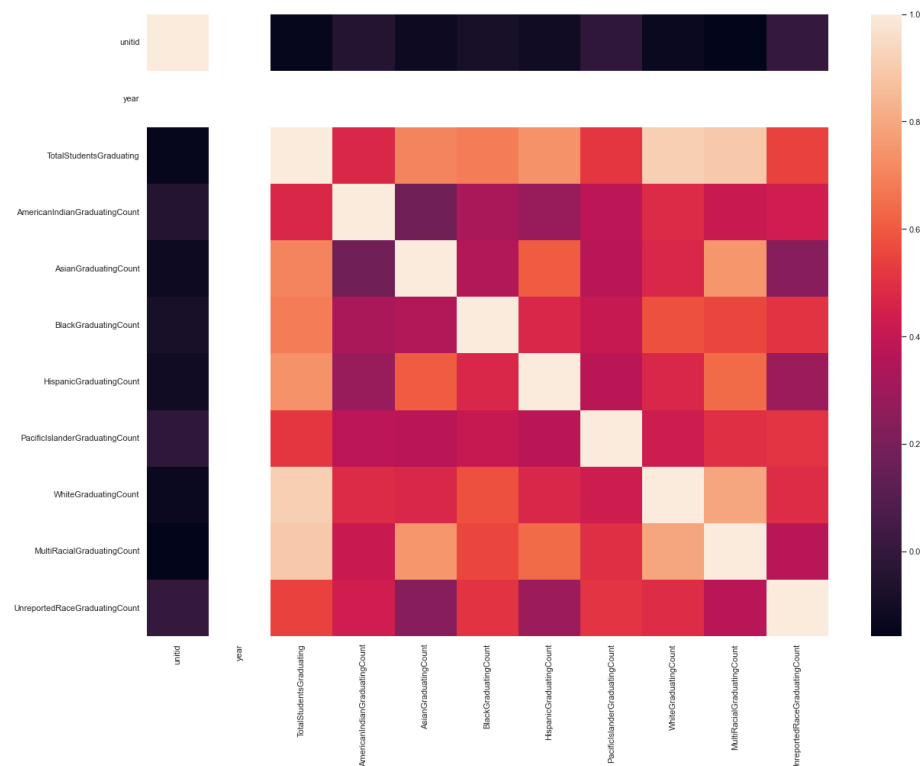
<AxesSubplot:>



This correlation plot shows if any of the undergraduate 12-month enrollment variables were correlated with each other. I found that very few variables were correlated. This correlation plot shows that the white and multiracial undergraduate count is correlated with the total undergraduate count.

```
corr3 = df3.corr()
plt.subplots(figsize=(20,15))
sns.heatmap(corr3,
             xticklabels=corr3.columns.values,
             yticklabels=corr3.columns.values)
```

<AxesSubplot:>



This correlation plot to see if any of the variables reporting graduation counts were correlated with each other. I found that very few variables were correlated. This correlation plot shows that each white graduation counts and multiracial graduation counts are correlated. Additionally, the Asian graduation count is correlated with multiracial graduating count.

The analysis of the correlation plots reiterates the importance of drawing insights from subgroups, to gain a picture of what compositional diversity looks like in American schools.

Correlation Tests

Tests whether two samples have a linear relationship.

```
from scipy.stats import pearsonr
stat, p = pearsonr(df['TotalGradRate'], df['WhiteGradRate'])
print('stat= %.3f, p= %.3f' % (stat, p))
if p > 0.05:
    print('Total Graduation Rate and White Graduation Rate are probably independent')
else:
    print('Total Graduation Rate and White Graduation Rate are probably dependent')

stat= 0.981, p= 0.000
```


Total Graduation Rate and White Graduation Rate are probably dependent

```
stat, p = pearsonr(df['TotalGradRate'], df['BlackGradRate'])
print('stat= %.3f, p= %.3f' % (stat, p))
if p > 0.05:
    print('Total Graduation Rate and Black Graduation Rate are probably independent')
else:
    print('Total Graduation Rate and Black Graduation Rate are probably dependent')
```

stat= 0.859, p= 0.000

Total Graduation Rate and Black Graduation Rate are probably dependent

```
stat, p = pearsonr(df['TotalGradRate'], df['AsianGradRate'])
print('stat= %.3f, p= %.3f' % (stat, p))
if p > 0.05:
    print('Total Graduation Rate and Asian Graduation Rate are probably independent')
else:
    print('Total Graduation Rate and Asian Graduation Rate are probably dependent')
```

stat= 0.757, p= 0.000

Total Graduation Rate and Asian Graduation Rate are probably dependent

```
stat, p = pearsonr(df['TotalGradRate'], df['HispanicGradRate'])
print('stat= %.3f, p= %.3f' % (stat, p))
if p > 0.05:
    print('Total Graduation Rate and Hispanic Graduation Rate are probably independent')
else:
    print('Total Graduation Rate and Hispanic Graduation Rate are probably dependent')
```

stat= 0.902, p= 0.000

Total Graduation Rate and Hispanic Graduation Rate are probably dependent

```
stat, p = pearsonr(df['TotalGradRate'], df['MultiRacialGradRate'])
print('stat= %.3f, p= %.3f' % (stat, p))
if p > 0.05:
    print('Total Graduation Rate and Mutli Racial Graduation Rate are probably independent')
else:
    print('Total Graduation Rate and Mutli Racial Graduation Rate are probably dependent')
```

stat= 0.848, p= 0.000

Total Graduation Rate and Mutli Racial Graduation Rate are probably dependent

```
stat, p = pearsonr(df['TotalGradRate'], df['MultiRacialGradRate'])
print('stat= %.3f, p= %.3f' % (stat, p))
if p > 0.05:
    print('Total Graduation Rate and Mutli Racial Graduation Rate are probably independent')
else:
    print('Total Graduation Rate and Mutli Racial Graduation Rate are probably dependent')
```

stat= 0.848, p= 0.000

Total Graduation Rate and Mutli Racial Graduation Rate are probably dependent

```

stat, p = pearsonr(df['TotalGradRate'], df['UnreportedRaceGradRate'])
print('stat= %.3f, p= %.3f' % (stat, p))
if p > 0.05:
    print('Total Graduation Rate and Unreported Race Graduation Rate are probably independent')
else:
    print('Total Graduation Rate and Unreported Race Graduation Rate are probably dependent')

stat= 0.738, p= 0.000
Total Graduation Rate and Unreported Race Graduation Rate are probably dependent

stat, p = pearsonr(df['TransferOutRate'], df['WhiteUndergradPercent'])
print('stat= %.3f, p= %.3f' % (stat, p))
if p > 0.05:
    print('The transfer out rate and White Undergraduate Percent are probably independent')
else:
    print('The transfer out rate and White Undergraduate Percent are probably dependent')

stat= 0.160, p= 0.003
The transfer out rate and White Undergraduate Percent are probably dependent

stat, p = pearsonr(df['PercentGrantedAid'], df['TransferOutRate'])
print('stat= %.3f, p= %.3f' % (stat, p))
if p > 0.05:
    print('The transfer out rate and Percent Granted Aid are probably independent')
else:
    print('The transfer out rate and Percent Granted Aid are probably dependent')

stat= 0.215, p= 0.000
The transfer out rate and Percent Granted Aid are probably dependent

stat, p = pearsonr(df['TotalGradRate'], df['PercentGrantedAid'])
print('stat= %.3f, p= %.3f' % (stat, p))
if p > 0.05:
    print('Total Graduation Rate and Percent Granted Aid are probably independent')
else:
    print('Total Graduation Rate and Percent Granted Aid are probably dependent')

stat= -0.020, p= 0.706
Total Graduation Rate and Percent Granted Aid are probably independent

stat, p = pearsonr(df2['TotalUndergradCount'], df2['WhiteUndergradCount'])
print('stat= %.3f, p= %.3f' % (stat, p))
if p > 0.05:
    print('Total Undergrad Count and White Undergrad Count are probably independent')
else:
    print('Total Undergrad Count and White Undergrad Count are probably dependent')

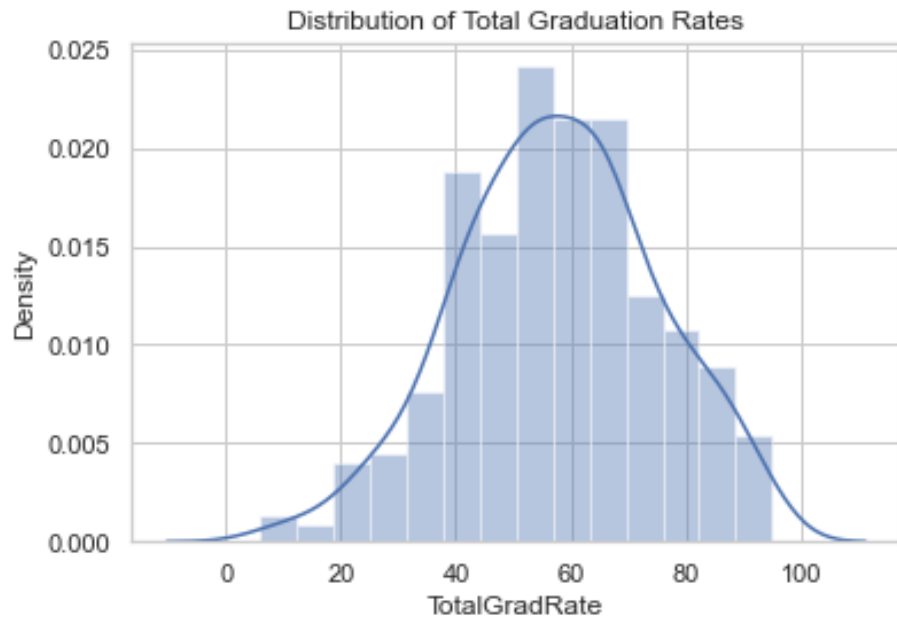
stat= 0.894, p= 0.000
Total Undergrad Count and White Undergrad Count are probably dependent

```

Visualizations

```
fig, ax = plt.subplots()
ax.set_title('Distribution of Total Graduation Rates')
sns.distplot(df.TotalGradRate)
```

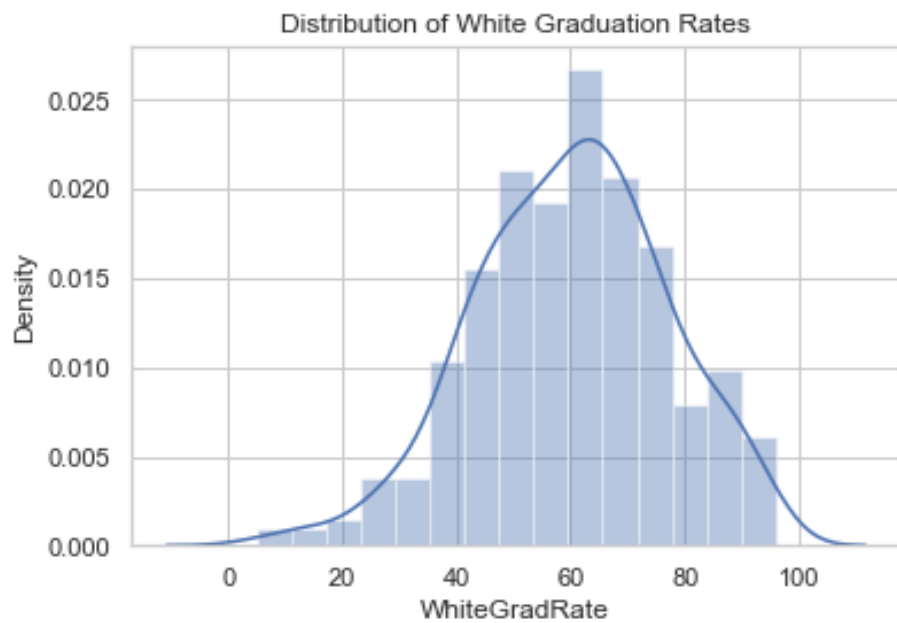
```
<AxesSubplot:title={'center': 'Distribution of Total Graduation Rates'}, xlabel='TotalGradRate'
```



This graph indicates that the average total graduation rate is approximately 50%.

```
fig, ax = plt.subplots()
ax.set_title('Distribution of White Graduation Rates')
sns.distplot(df.WhiteGradRate)
```

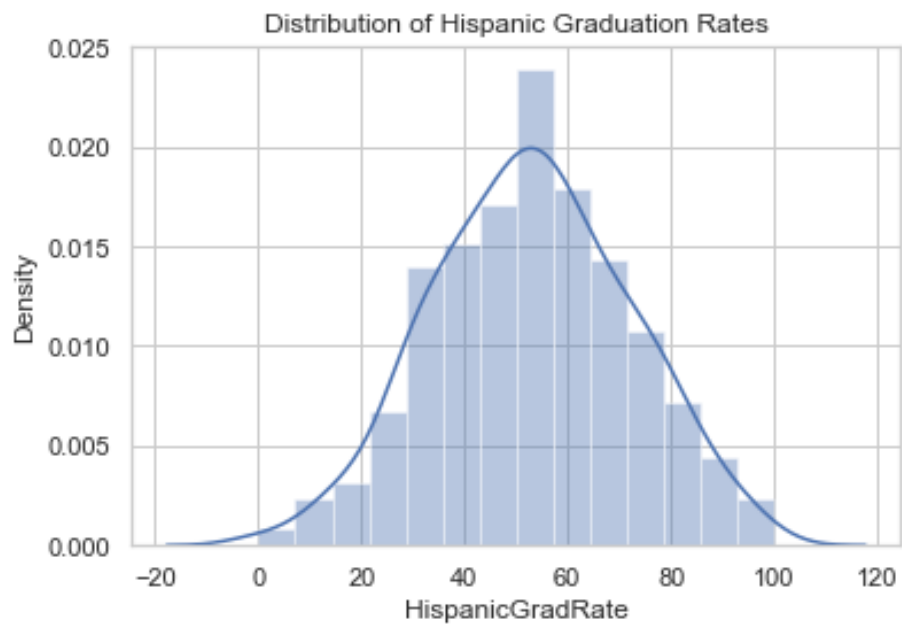
```
<AxesSubplot:title={'center': 'Distribution of White Graduation Rates'}, xlabel='WhiteGradRate'
```



This graph indicates that the average graduation rate among white undergraduates is approximately 60%.

```
fig, ax = plt.subplots()
ax.set_title('Distribution of Hispanic Graduation Rates')
sns.distplot(df.HispanicGradRate)

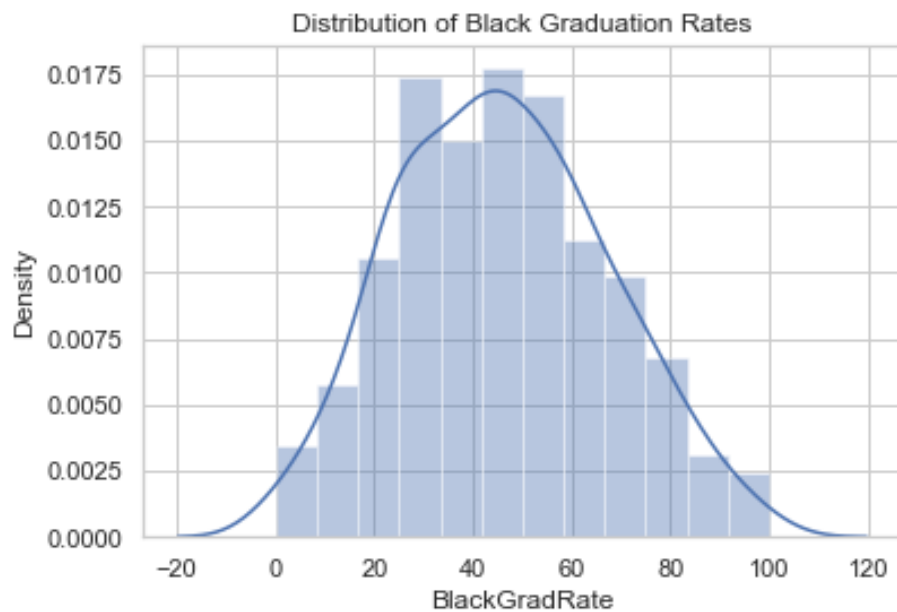
<AxesSubplot:title={'center': 'Distribution of Hispanic Graduation Rates'}, xlabel='Hispanic'
```



This graph indicates that the average graduation rate among hispanic undergraduates is approximately 50%.

```
fig, ax = plt.subplots()
ax.set_title('Distribution of Black Graduation Rates')
sns.distplot(df.BlackGradRate)
```

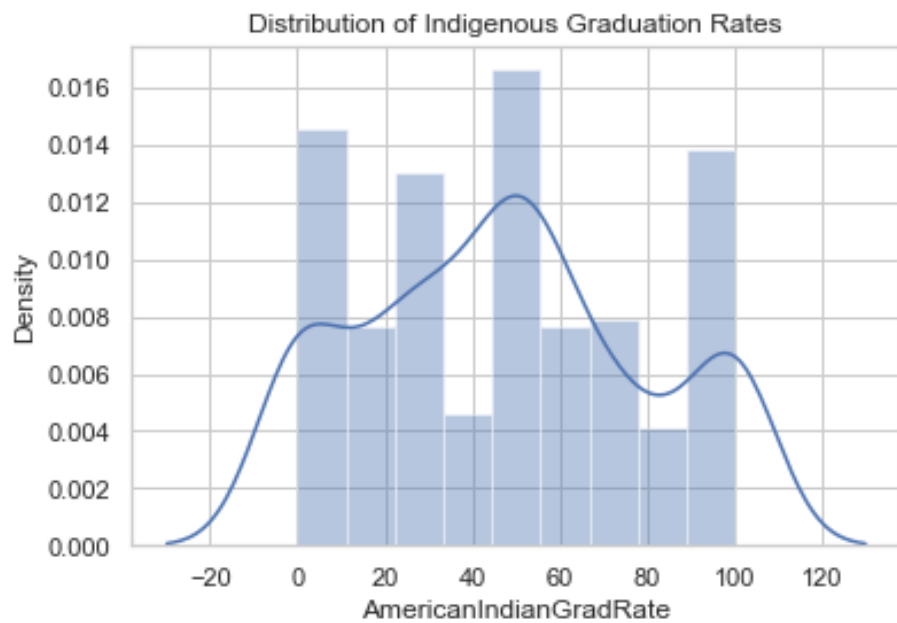
```
<AxesSubplot:title={'center': 'Distribution of Black Graduation Rates'}, xlabel='BlackGradRate'>
```



This graph indicates that the average graduation rate of black undergraduate students is approximately 50%. The graph is slightly right-skewed, indicating that the majority of students have a graduation rate of less than 50%.

```
fig, ax = plt.subplots()
ax.set_title('Distribution of Indigenous Graduation Rates')
sns.distplot(df.AmericanIndianGradRate)

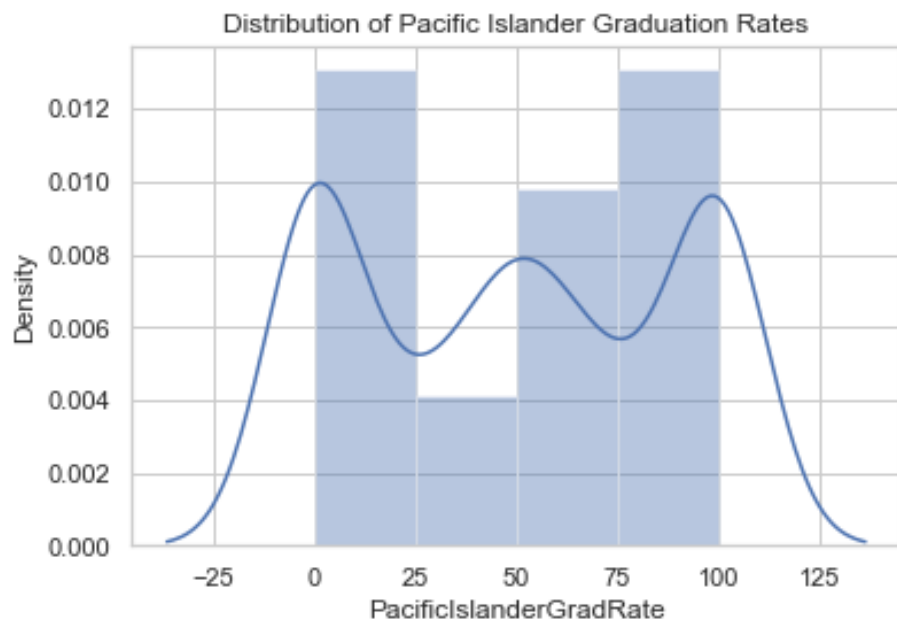
<AxesSubplot:title={'center':'Distribution of Indigenous Graduation Rates'}, xlabel='American Indian Graduation Rate', yticklabels=['Density'], yticks=[0.0, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95, 1.0], xticklabels=[0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0, 1.1, 1.2, 1.3, 1.4, 1.5, 1.6, 1.7, 1.8, 1.9, 2.0, 2.1, 2.2, 2.3, 2.4, 2.5, 2.6, 2.7, 2.8, 2.9, 3.0, 3.1, 3.2, 3.3, 3.4, 3.5, 3.6, 3.7, 3.8, 3.9, 4.0, 4.1, 4.2, 4.3, 4.4, 4.5, 4.6, 4.7, 4.8, 4.9, 5.0, 5.1, 5.2, 5.3, 5.4, 5.5, 5.6, 5.7, 5.8, 5.9, 6.0, 6.1, 6.2, 6.3, 6.4, 6.5, 6.6, 6.7, 6.8, 6.9, 7.0, 7.1, 7.2, 7.3, 7.4, 7.5, 7.6, 7.7, 7.8, 7.9, 8.0, 8.1, 8.2, 8.3, 8.4, 8.5, 8.6, 8.7, 8.8, 8.9, 9.0, 9.1, 9.2, 9.3, 9.4, 9.5, 9.6, 9.7, 9.8, 9.9, 10.0, 10.1, 10.2, 10.3, 10.4, 10.5, 10.6, 10.7, 10.8, 10.9, 11.0, 11.1, 11.2, 11.3, 11.4, 11.5, 11.6, 11.7, 11.8, 11.9, 12.0, 12.1, 12.2, 12.3, 12.4, 12.5, 12.6, 12.7, 12.8, 12.9, 13.0, 13.1, 13.2, 13.3, 13.4, 13.5, 13.6, 13.7, 13.8, 13.9, 14.0, 14.1, 14.2, 14.3, 14.4, 14.5, 14.6, 14.7, 14.8, 14.9, 15.0, 15.1, 15.2, 15.3, 15.4, 15.5, 15.6, 15.7, 15.8, 15.9, 16.0, 16.1, 16.2, 16.3, 16.4, 16.5, 16.6, 16.7, 16.8, 16.9, 17.0, 17.1, 17.2, 17.3, 17.4, 17.5, 17.6, 17.7, 17.8, 17.9, 18.0, 18.1, 18.2, 18.3, 18.4, 18.5, 18.6, 18.7, 18.8, 18.9, 19.0, 19.1, 19.2, 19.3, 19.4, 19.5, 19.6, 19.7, 19.8, 19.9, 20.0, 20.1, 20.2, 20.3, 20.4, 20.5, 20.6, 20.7, 20.8, 20.9, 21.0, 21.1, 21.2, 21.3, 21.4, 21.5, 21.6, 21.7, 21.8, 21.9, 22.0, 22.1, 22.2, 22.3, 22.4, 22.5, 22.6, 22.7, 22.8, 22.9, 23.0, 23.1, 23.2, 23.3, 23.4, 23.5, 23.6, 23.7, 23.8, 23.9, 24.0, 24.1, 24.2, 24.3, 24.4, 24.5, 24.6, 24.7, 24.8, 24.9, 25.0, 25.1, 25.2, 25.3, 25.4, 25.5, 25.6, 25.7, 25.8, 25.9, 26.0, 26.1, 26.2, 26.3, 26.4, 26.5, 26.6, 26.7, 26.8, 26.9, 27.0, 27.1, 27.2, 27.3, 27.4, 27.5, 27.6, 27.7, 27.8, 27.9, 28.0, 28.1, 28.2, 28.3, 28.4, 28.5, 28.6, 28.7, 28.8, 28.9, 29.0, 29.1, 29.2, 29.3, 29.4, 29.5, 29.6, 29.7, 29.8, 29.9, 30.0, 30.1, 30.2, 30.3, 30.4, 30.5, 30.6, 30.7, 30.8, 30.9, 31.0, 31.1, 31.2, 31.3, 31.4, 31.5, 31.6, 31.7, 31.8, 31.9, 32.0, 32.1, 32.2, 32.3, 32.4, 32.5, 32.6, 32.7, 32.8, 32.9, 33.0, 33.1, 33.2, 33.3, 33.4, 33.5, 33.6, 33.7, 33.8, 33.9, 34.0, 34.1, 34.2, 34.3, 34.4, 34.5, 34.6, 34.7, 34.8, 34.9, 35.0, 35.1, 35.2, 35.3, 35.4, 35.5, 35.6, 35.7, 35.8, 35.9, 36.0, 36.1, 36.2, 36.3, 36.4, 36.5, 36.6, 36.7, 36.8, 36.9, 37.0, 37.1, 37.2, 37.3, 37.4, 37.5, 37.6, 37.7, 37.8, 37.9, 38.0, 38.1, 38.2, 38.3, 38.4, 38.5, 38.6, 38.7, 38.8, 38.9, 39.0, 39.1, 39.2, 39.3, 39.4, 39.5, 39.6, 39.7, 39.8, 39.9, 40.0, 40.1, 40.2, 40.3, 40.4, 40.5, 40.6, 40.7, 40.8, 40.9, 41.0, 41.1, 41.2, 41.3, 41.4, 41.5, 41.6, 41.7, 41.8, 41.9, 42.0, 42.1, 42.2, 42.3, 42.4, 42.5, 42.6, 42.7, 42.8, 42.9, 43.0, 43.1, 43.2, 43.3, 43.4, 43.5, 43.6, 43.7, 43.8, 43.9, 44.0, 44.1, 44.2, 44.3, 44.4, 44.5, 44.6, 44.7, 44.8, 44.9, 45.0, 45.1, 45.2, 45.3, 45.4, 45.5, 45.6, 45.7, 45.8, 45.9, 46.0, 46.1, 46.2, 46.3, 46.4, 46.5, 46.6, 46.7, 46.8, 46.9, 47.0, 47.1, 47.2, 47.3, 47.4, 47.5, 47.6, 47.7, 47.8, 47.9, 48.0, 48.1, 48.2, 48.3, 48.4, 48.5, 48.6, 48.7, 48.8, 48.9, 49.0, 49.1, 49.2, 49.3, 49.4, 49.5, 49.6, 49.7, 49.8, 49.9, 50.0, 50.1, 50.2, 50.3, 50.4, 50.5, 50.6, 50.7, 50.8, 50.9, 51.0, 51.1, 51.2, 51.3, 51.4, 51.5, 51.6, 51.7, 51.8, 51.9, 52.0, 52.1, 52.2, 52.3, 52.4, 52.5, 52.6, 52.7, 52.8, 52.9, 53.0, 53.1, 53.2, 53.3, 53.4, 53.5, 53.6, 53.7, 53.8, 53.9, 54.0, 54.1, 54.2, 54.3, 54.4, 54.5, 54.6, 54.7, 54.8, 54.9, 55.0, 55.1, 55.2, 55.3, 55.4, 55.5, 55.6, 55.7, 55.8, 55.9, 56.0, 56.1, 56.2, 56.3, 56.4, 56.5, 56.6, 56.7, 56.8, 56.9, 57.0, 57.1, 57.2, 57.3, 57.4, 57.5, 57.6, 57.7, 57.8, 57.9, 58.0, 58.1, 58.2, 58.3, 58.4, 58.5, 58.6, 58.7, 58.8, 58.9, 59.0, 59.1, 59.2, 59.3, 59.4, 59.5, 59.6, 59.7, 59.8, 59.9, 60.0, 60.1, 60.2, 60.3, 60.4, 60.5, 60.6, 60.7, 60.8, 60.9, 61.0, 61.1, 61.2, 61.3, 61.4, 61.5, 61.6, 61.7, 61.8, 61.9, 62.0, 62.1, 62.2, 62.3, 62.4, 62.5, 62.6, 62.7, 62.8, 62.9, 63.0, 63.1, 63.2, 63.3, 63.4, 63.5, 63.6, 63.7, 63.8, 63.9, 64.0, 64.1, 64.2, 64.3, 64.4, 64.5, 64.6, 64.7, 64.8, 64.9, 65.0, 65.1, 65.2, 65.3, 65.4, 65.5, 65.6, 65.7, 65.8, 65.9, 66.0, 66.1, 66.2, 66.3, 66.4, 66.5, 66.6, 66.7, 66.8
```



This graph indicates that the graduation rates of indigenous undergraduate students vary greatly, with an average graduation rate of 50%.

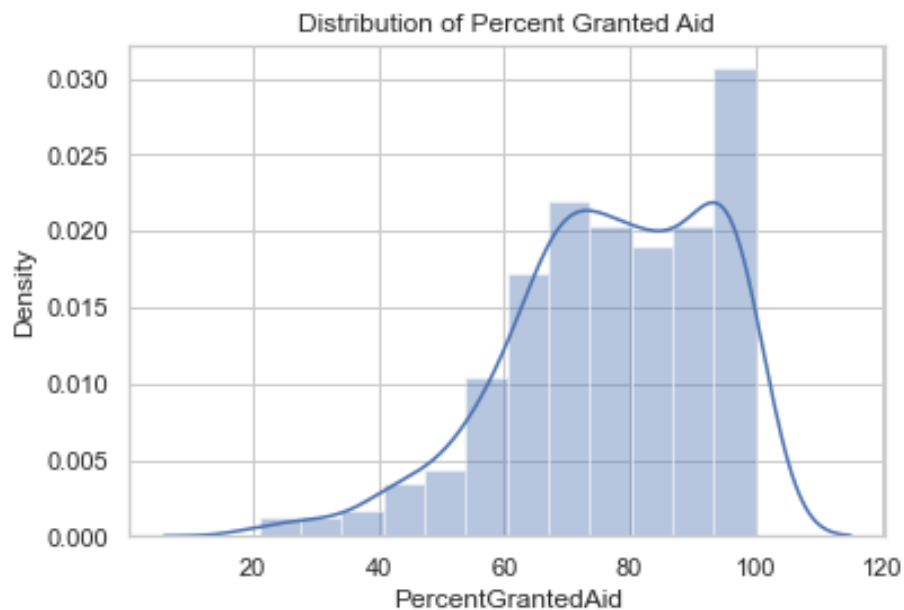
```
fig, ax = plt.subplots()
ax.set_title('Distribution of Pacific Islander Graduation Rates')
sns.distplot(df.PacificIslanderGradRate)
```

<AxesSubplot:title={'center': 'Distribution of Pacific Islander Graduation Rates'}, xlabel='P



The graph indicates that the average graduation rate of pacific islander students is highly varied.

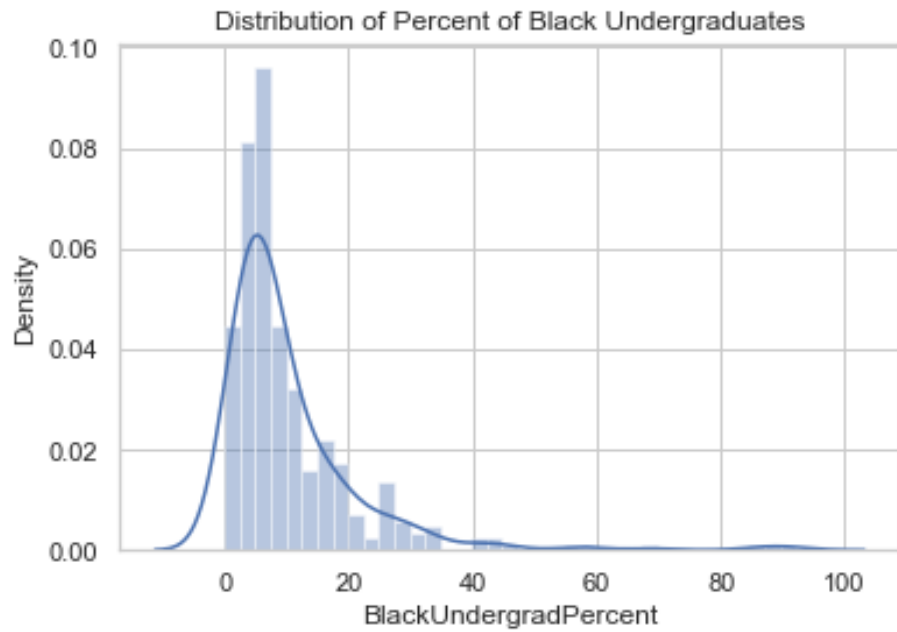
```
fig, ax = plt.subplots()
ax.set_title('Distribution of Percent Granted Aid')
sns.distplot(df.PercentGrantedAid)
<AxesSubplot:title={'center': 'Distribution of Percent Granted Aid'}, xlabel='PercentGrantedAid'
```

This graph shows that the average percent of students granted financial aid is near 90% and is left-skewed due to the long tail on the left, indicating that the mean is less than the median. The bulk of the observations are toward 80% or greater.

```
fig, ax = plt.subplots()
ax.set_title('Distribution of Percent of Black Undergraduates')
sns.distplot(df.BlackUndergradPercent)
```

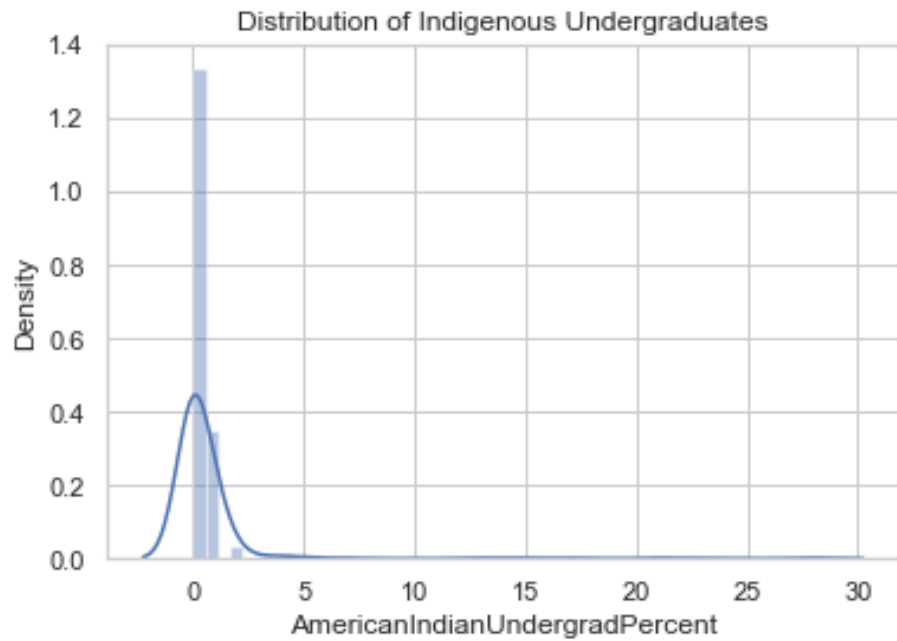
```
<AxesSubplot:title={'center': 'Distribution of Percent of Black Undergraduates'}, xlabel='Bl
```



This graph shows that the average percent of black undergraduate students is approximately near 10%, and the graph is right-skewed, indicating that the majority of Title IV schools have less than 15% of black undergraduate students.

```
fig, ax = plt.subplots()
ax.set_title('Distribution of Indigenous Undergraduates')
sns.distplot(df.AmericanIndianUndergradPercent)

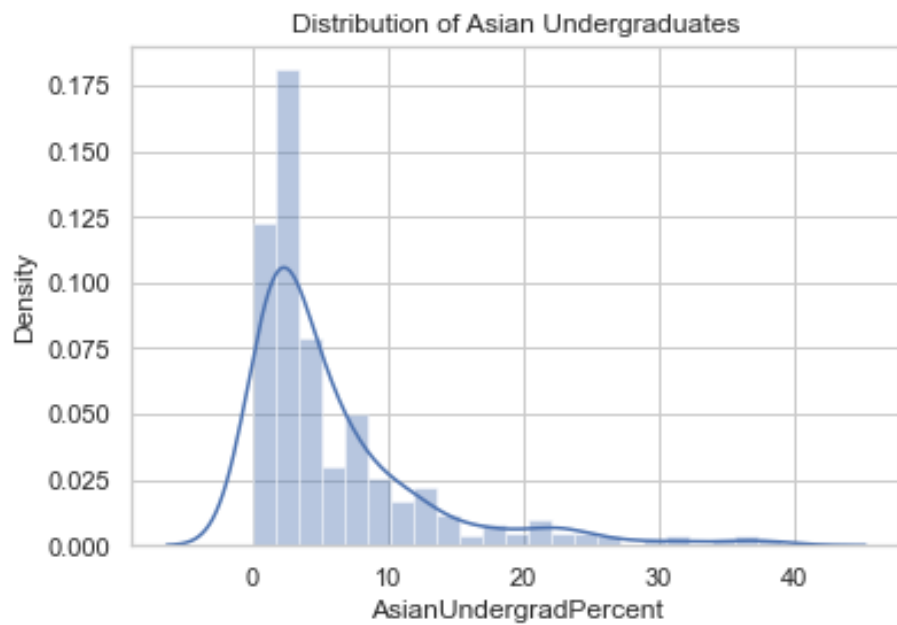
<AxesSubplot:title={'center': 'Distribution of Indigenous Undergraduates'}, xlabel='American
```



This graph indicates the breakdown of American Indian Undergraduates. The graph is right-skewed, indicating that the majority of undergraduate populations in Title IV schools have less than 5% of Indigenous students.

```
fig, ax = plt.subplots()
ax.set_title('Distribution of Asian Undergraduates')
sns.distplot(df.AsianUndergradPercent)

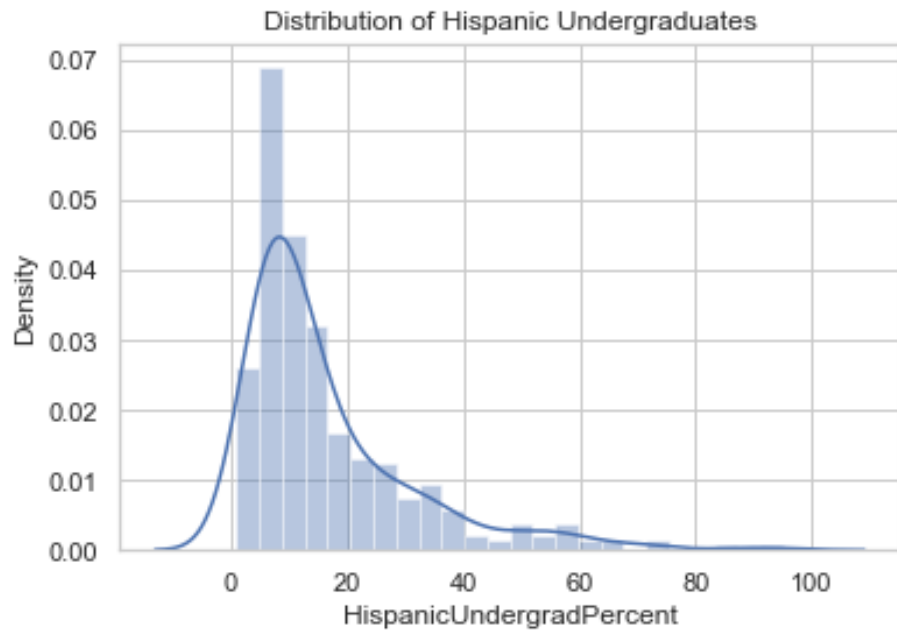
<AxesSubplot:title={'center': 'Distribution of Asian Undergraduates'}, xlabel='AsianUndergradPercent'
```



This graph indicates the breakdown of Asian Undergraduates. The graph is right-skewed, indicating that the majority of undergraduate populations in Title IV schools have less than 20% of Asian students.

```
fig, ax = plt.subplots()
ax.set_title('Distribution of Hispanic Undergraduates')
sns.distplot(df.HispanicUndergradPercent)

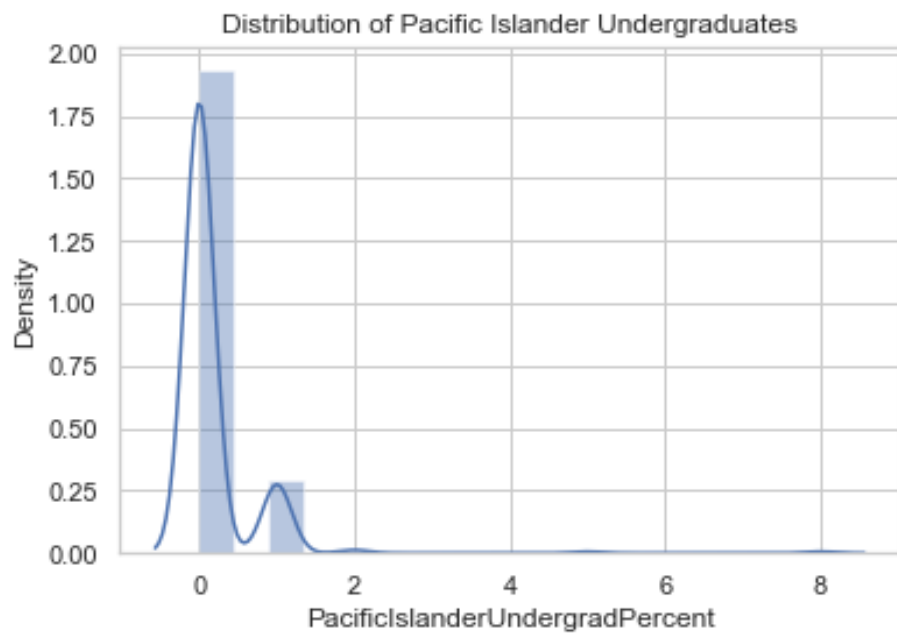
<AxesSubplot:title={'center': 'Distribution of Hispanic Undergraduates'}, xlabel='HispanicUn
```



This graph indicates the breakdown of Hispanic Undergraduates. The graph is right-skewed, indicating that the majority of undergraduate populations in Title IV schools have less than 30% of Hispanic students.

```
fig, ax = plt.subplots()
ax.set_title('Distribution of Pacific Islander Undergraduates')
sns.distplot(df.PacificIslanderUndergradPercent)
```

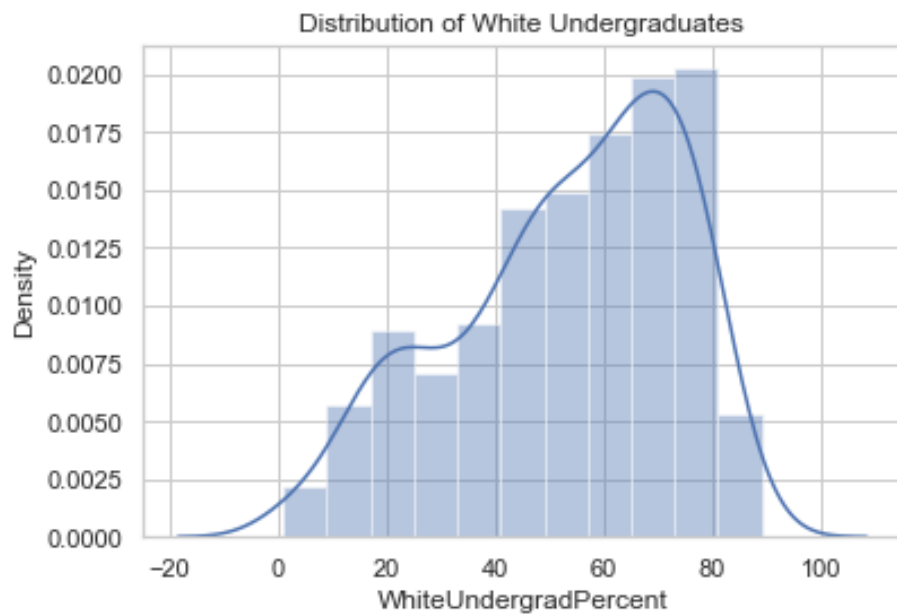
```
<AxesSubplot:title={'center': 'Distribution of Pacific Islander Undergraduates'}, xlabel='Pacific Islander Undergrad Percent'>
```



This graph indicates the breakdown of Pacific Islander Undergraduates. The graph is right-skewed, indicating that the majority of undergraduate populations in Title IV schools have less than 2% of Pacific Islander students.

```
fig, ax = plt.subplots()
ax.set_title('Distribution of White Undergraduates')
sns.distplot(df.WhiteUndergradPercent)

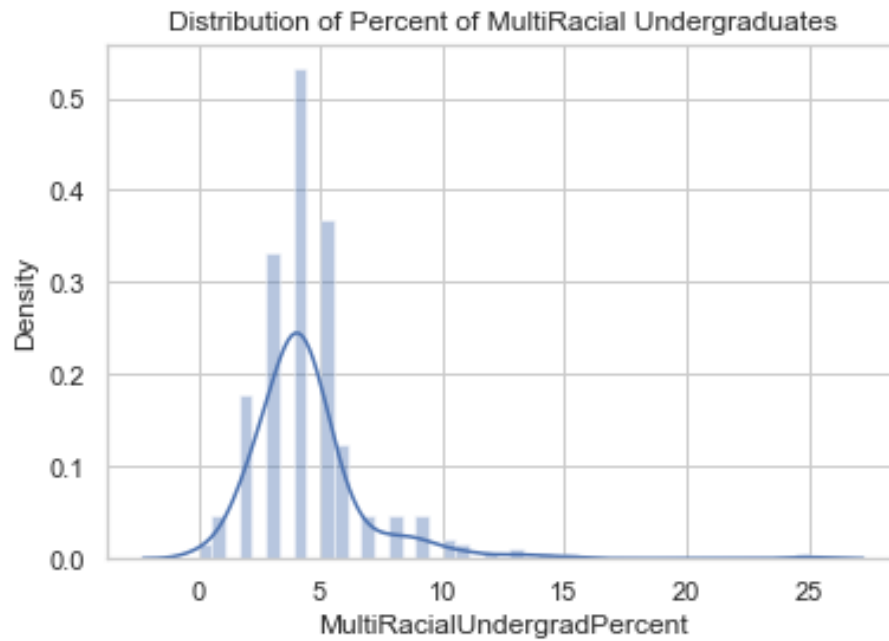
<AxesSubplot:title={'center': 'Distribution of White Undergraduates'}, xlabel='WhiteUndergradPercent'>
```



This graph indicates the breakdown of White Undergraduates. The graph is slightly left-skewed, indicating that the majority of undergraduate populations in Title IV schools have more than 50% of White students.

```
fig, ax = plt.subplots()
ax.set_title('Distribution of Percent of MultiRacial Undergraduates')
sns.distplot(df.MultiRacialUndergradPercent)

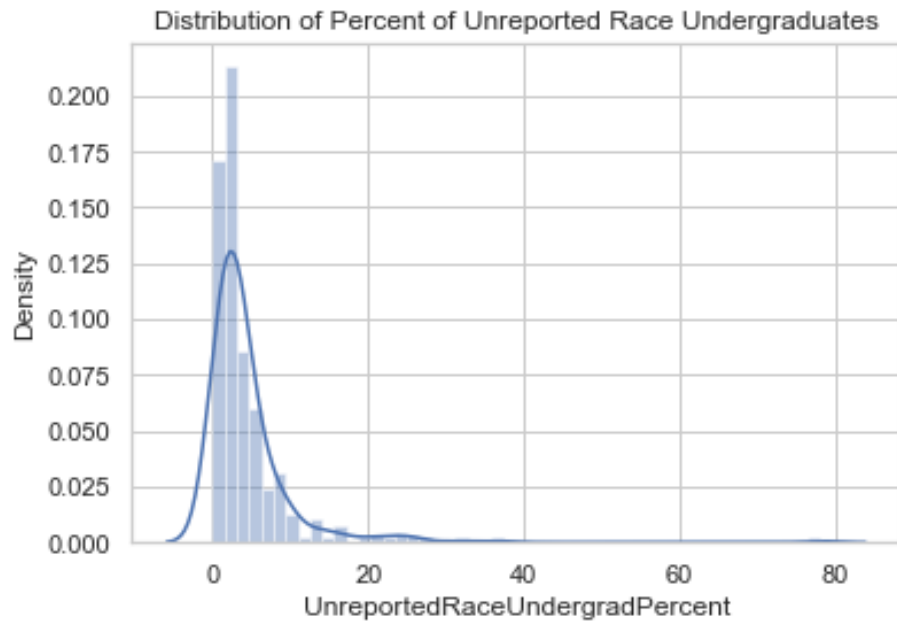
<AxesSubplot:title={'center': 'Distribution of Percent of MultiRacial Undergraduates'}, xlabel=
```



This graph indicates the breakdown of Multi-Racial Undergraduate students. This graph is left-skewed, indicating that the majority of undergraduate populations have between 0-10% of Multi-Racial students.

```
fig, ax = plt.subplots()
ax.set_title('Distribution of Percent of Unreported Race Undergraduates')
sns.distplot(df.UnreportedRaceUndergradPercent)
```

```
<AxesSubplot:title={'center': 'Distribution of Percent of Unreported Race Undergraduates'}, x
```

This graph indicates the breakdown of Unreported Race Undergraduate students. This graph is left-skewed, indicating that the majority of undergraduate populations have between 0-20% of Unreported Race students.

```
subset = df.sample(n=200)
```

```
subset.shape
```

```
(200, 26)
```

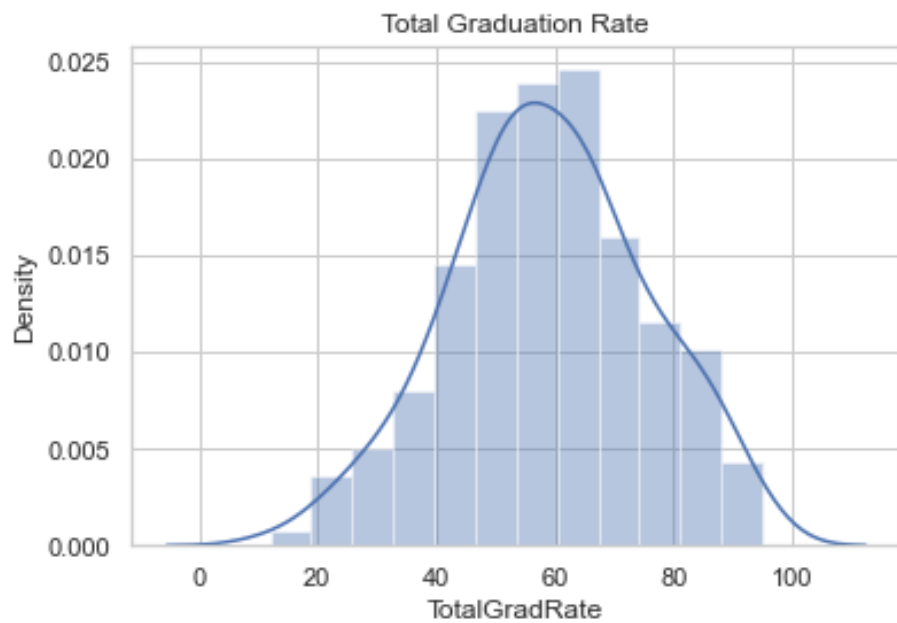
Creating a subset of size (n=200) randomly sampled schools to extract insights from a smaller group of schools. This allows us to perform deeper analysis on variables using a smaller subset of all 5,000 Title IV schools.

```
fig, ax = plt.subplots()
```

```
ax.set_title('Total Graduation Rate')
```

```
sns.distplot(subset.TotalGradRate)
```

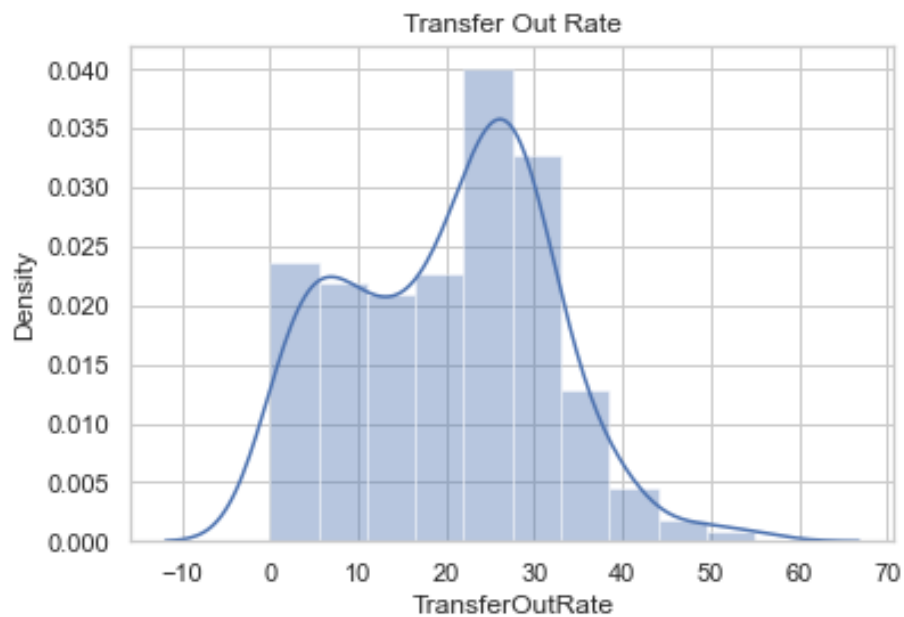
```
<AxesSubplot:title={'center': 'Total Graduation Rate'}, xlabel='TotalGradRate', ylabel='Dens
```



This graph indicates the distribution of total graduation rates. This graph is slightly right-skewed, indicating that the majority of institutions have a total graduation rate of greater than 60%.

```
fig, ax = plt.subplots()
ax.set_title('Transfer Out Rate')
sns.distplot(subset.TransferOutRate)
```

```
<AxesSubplot:title={'center': 'Transfer Out Rate'}, xlabel='TransferOutRate', ylabel='Density'
```



This graph indicates the distribution of transfer rates across Title IV schools. This graph is right-skewed with a few notable outliers. The average transfer out rate is approximately 30%, and the right-skewed nature indicates that the majority of schools have a transfer out rate that is greater than 30%.

```
#Get dataframe where school is HBCU status
```

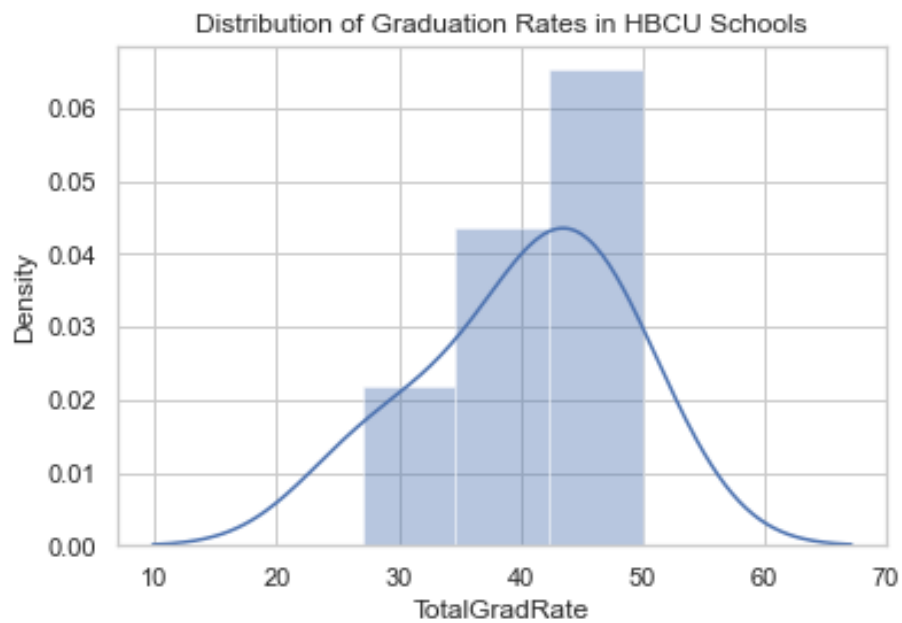
```
HBCU_df = df[df['HBCUStatus'] == 1]
```

```
fig, ax = plt.subplots()
```

```
ax.set_title("Distribution of Graduation Rates in HBCU Schools")
```

```
sns.distplot(HBCU_df.TotalGradRate)
```

```
<AxesSubplot:title={'center': 'Distribution of Graduation Rates in HBCU Schools'}, xlabel='Total Graduation Rate'>
```



The graduation rates among HBCU schools are left-skewed, with an average graduation rate of 45%. The left-skewed nature of the graph indicates that the majority of HBCU schools hold graduation rates greater than 45%.

```
fig, ax = plt.subplots()
ax.set_title("Distribution of Transfer Out Rates in HBCUs")
sns.distplot(HBCU_df.TransferOutRate)
```

```
<AxesSubplot:title={'center': 'Distribution of Transfer Out Rates in HBCUs'}, xlabel='Transf
```



The transfer out rates among HBCU schools are slightly left-skewed, with an average transfer rate of 30%. The left-skewed nature of the graph indicates that the majority of HBCU schools hold transfer rates greater than 30%.

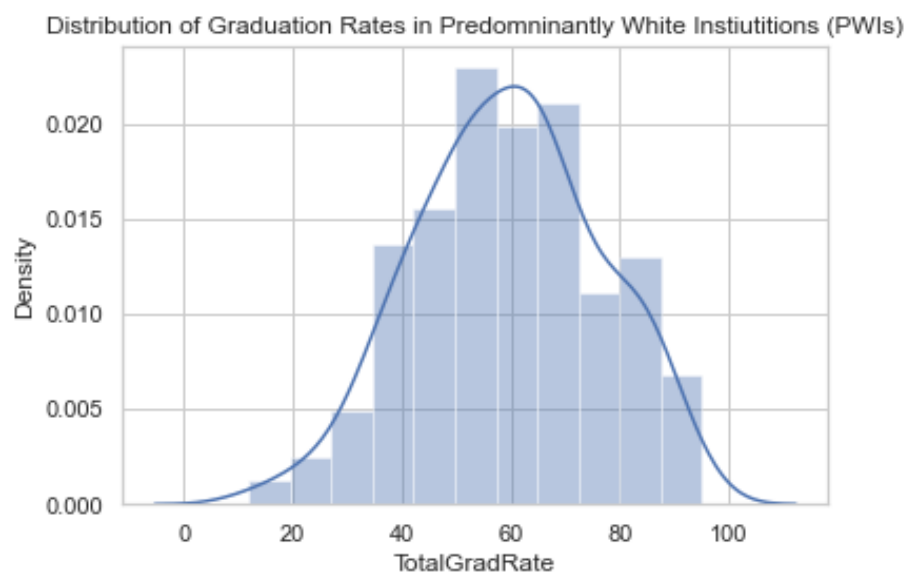
```
PWI_df = df[df['WhiteUndergradPercent'] >= 50]
```

```
fig, ax = plt.subplots()
```

```
ax.set_title("Distribution of Graduation Rates in Predominantly White Institutions (PWIs)")
```

```
sns.distplot(PWI_df.TotalGradRate)
```

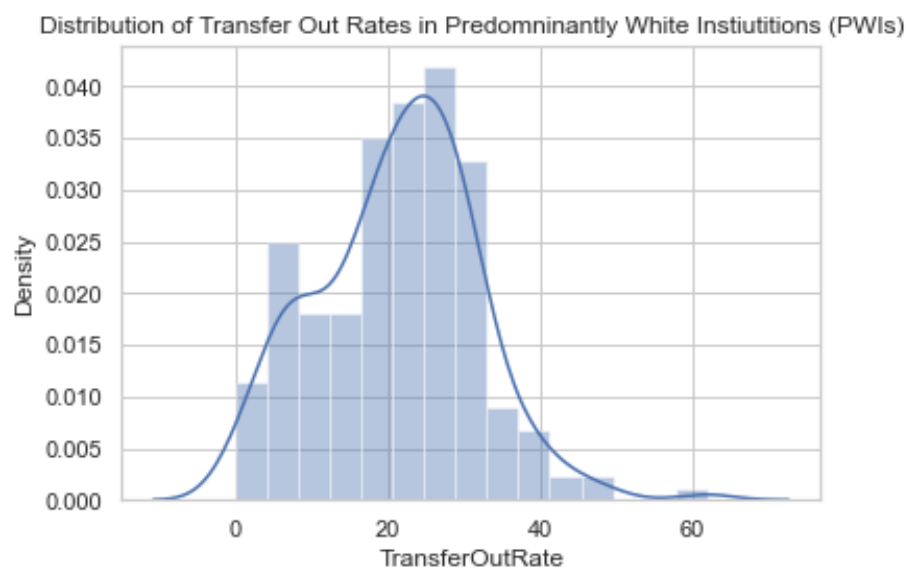
```
<AxesSubplot:title={'center': 'Distribution of Graduation Rates in Predominantly White Institutions'}
```



This graph indicates the distribution of the average graduation rate in institutions where white undergraduates account for greater than 50% of the undergraduate population. The average total graduation rate among PWIs is approximately 60%.

```
fig, ax = plt.subplots()
ax.set_title("Distribution of Transfer Out Rates in Predominantly White Institutions (PWIs)")
sns.distplot(PWI_df.TransferOutRate)
```

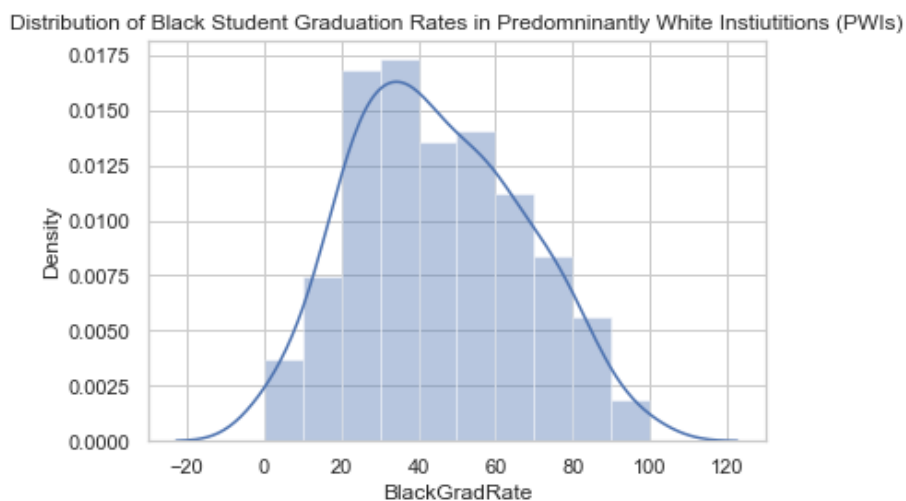
<AxesSubplot:title={'center': 'Distribution of Transfer Out Rates in Predominantly White Ins



This graph indicates the distribution of the transfer out rate in institutions where white undergraduates account for greater than 1/2 of the undergraduate population. The average transfer rate among PWIs is approximately 30%. The right-skewed nature of this graph indicates that the majority of sampled institutions have a transfer rate less than 30%.

```
fig, ax = plt.subplots()
ax.set_title("Distribution of Black Student Graduation Rates in Predominantly White Institutions")
sns.distplot(PWI_df.BlackGradRate)

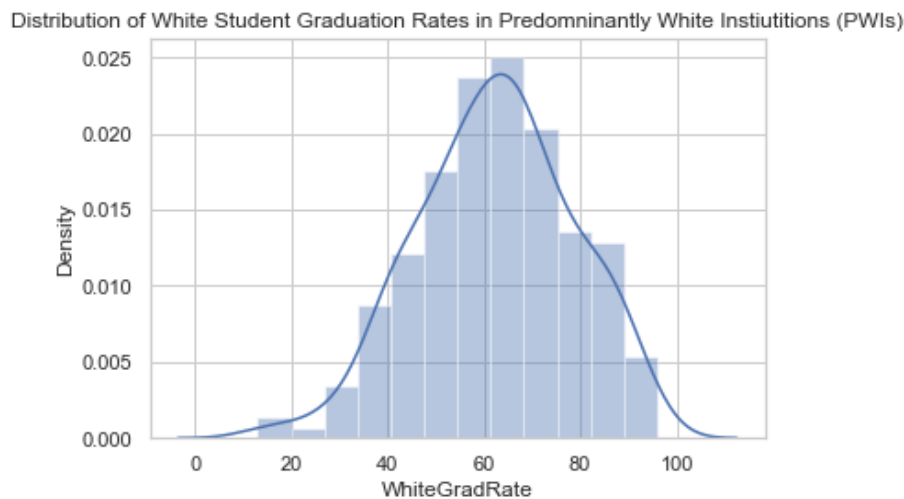
<AxesSubplot:title={'center': 'Distribution of Black Student Graduation Rates in Predominantly White Institutions'}
```



This graph indicates the distribution of black undergraduate graduation rates in institutions where white undergraduates account for greater than 1/2 of the undergraduate population. The average black graduation rate in PWIs is approximately 40%. The left-skewed nature of this graph indicates that the majority of sampled institutions have a black undergraduate graduation rate of greater than 40%.

```
fig, ax = plt.subplots()
ax.set_title("Distribution of White Student Graduation Rates in Predominantly White Institutions")
sns.distplot(PWI_df.WhiteGradRate)

<AxesSubplot:title={'center': 'Distribution of White Student Graduation Rates in Predominantly White Institutions'}
```

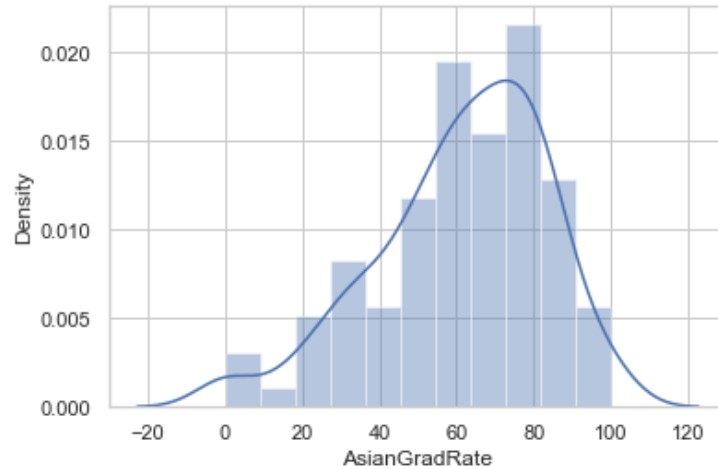


This graph indicates the distribution of the white undergraduate graduation rates in institutions where white undergraduates account for greater than 1/2 of the undergraduate population. The average transfer rate among PWIs is approximately 70%. The left-skewed nature of this graph indicates that the majority of sampled institutions have a white undergraduate graduation rate greater than 70%.

```
fig, ax = plt.subplots()
ax.set_title("Distribution of Asian Student Graduation Rates in Predominantly White Institutions")
sns.distplot(PWI_df.AsianGradRate)

<AxesSubplot:title={'center': 'Distribution of Asian Student Graduation Rates in Predominantly White Institutions'}
```


Distribution of Asian Student Graduation Rates in Predominantly White Institutions (PWIs)

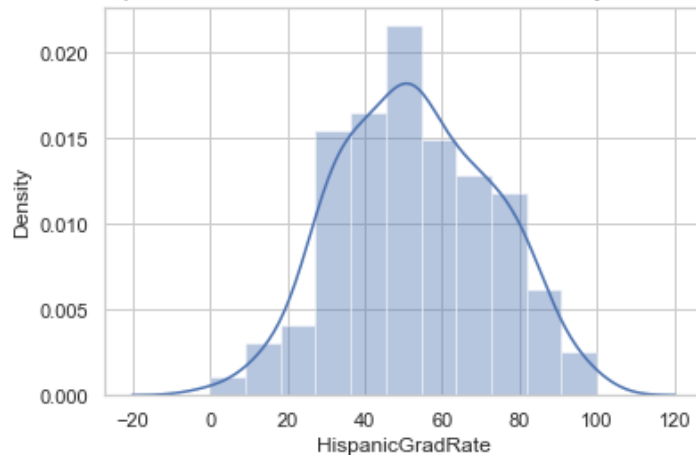


This graph indicates the distribution of the Asian graduation rates in institutions where white undergraduates account for greater than 1/2 of the undergraduate population. The average Asian graduation rate among PWIs is approximately 75%. The left-skewed nature of this graph indicates that the majority of sampled institutions have an asian undergraduate graduation rate of greater than 75%.

```
fig, ax = plt.subplots()
ax.set_title("Distribution of Hispanic Student Graduation Rates in Predominantly White Inst")
sns.distplot(PWI_df.HispanicGradRate)

<AxesSubplot:title={'center': 'Distribution of Hispanic Student Graduation Rates in Predominantly White Institutions (PWIs)'}>
```

Distribution of Hispanic Student Graduation Rates in Predominantly White Institutions (PWIs)

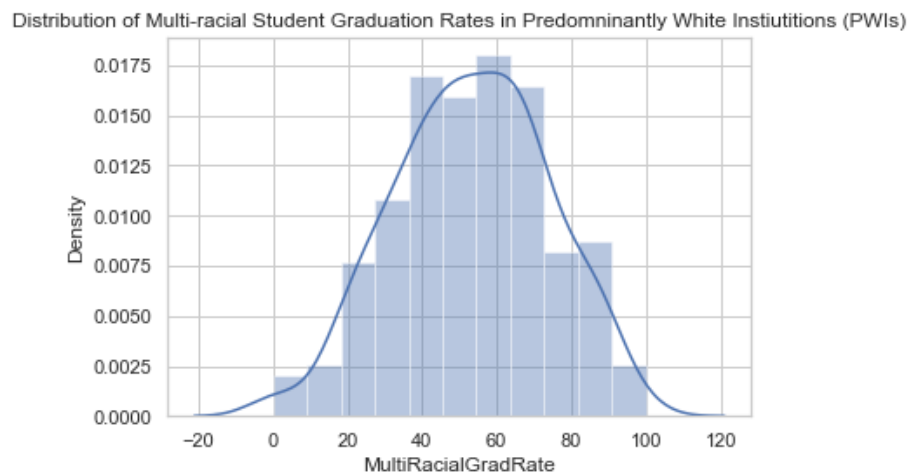


This graph indicates the distribution of the Hispanic undergraduate graduation rates in institutions where white undergraduates account for greater than 1/2

of the undergraduate population. The average Hispanic graduation rate among PWIs is approximately 50%.

```
fig, ax = plt.subplots()
ax.set_title("Distribution of Multi-racial Student Graduation Rates in Predominantly White Institutions (PWIs)")
sns.distplot(PWI_df.MultiRacialGradRate)
```

<AxesSubplot:title={'center': 'Distribution of Multi-racial Student Graduation Rates in Predominantly White Institutions (PWIs)'}>

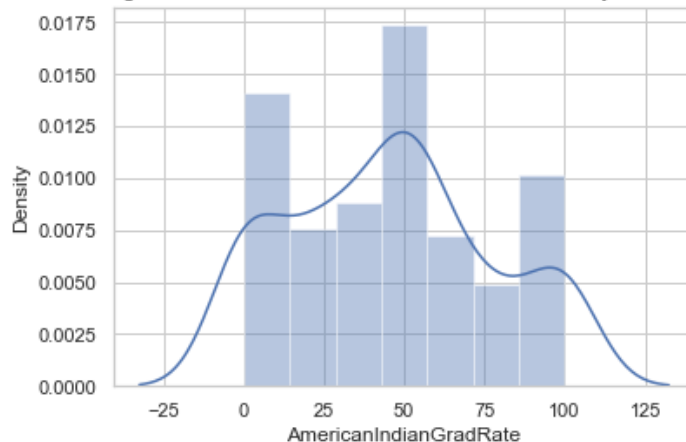


This graph indicates the distribution of the multi-racial undergraduate graduation rates in institutions where white undergraduates account for greater than 1/2 of the undergraduate population. The average multi-racial graduation rate among PWIs is approximately 55%.

```
fig, ax = plt.subplots()
ax.set_title("Distribution of Indigenous Student Graduation Rates in Predominantly White Institutions (PWIs)")
sns.distplot(PWI_df.AmericanIndianGradRate)
```

<AxesSubplot:title={'center': 'Distribution of Indigenous Student Graduation Rates in Predominantly White Institutions (PWIs)'}>

Distribution of Indigenous Student Graduation Rates in Predominantly White Institutions (PWIs)

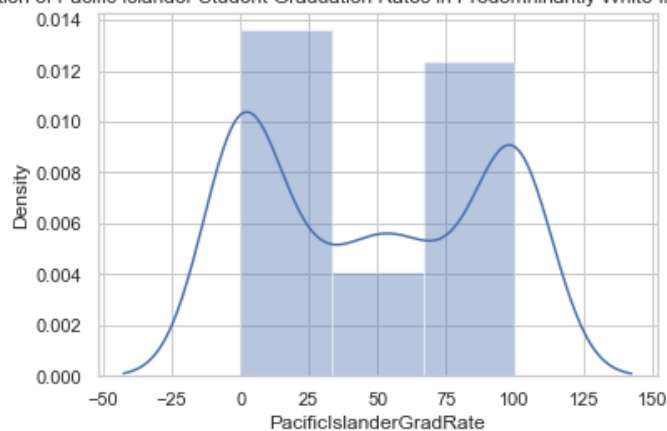


This graph indicates the distribution of the American Indian undergraduate graduation rate in institutions where white undergraduate students account for greater than 1/2 of the undergraduate population. The distribution for the American Indian Graduation Rate varies greatly. The average American Indian Graduation Rate among PWIs is approximately 50%.

```
fig, ax = plt.subplots()
ax.set_title("Distribution of Pacific Islander Student Graduation Rates in Predominantly White Institutions (PWIs)")
sns.distplot(PWI_df.PacificIslanderGradRate)
```

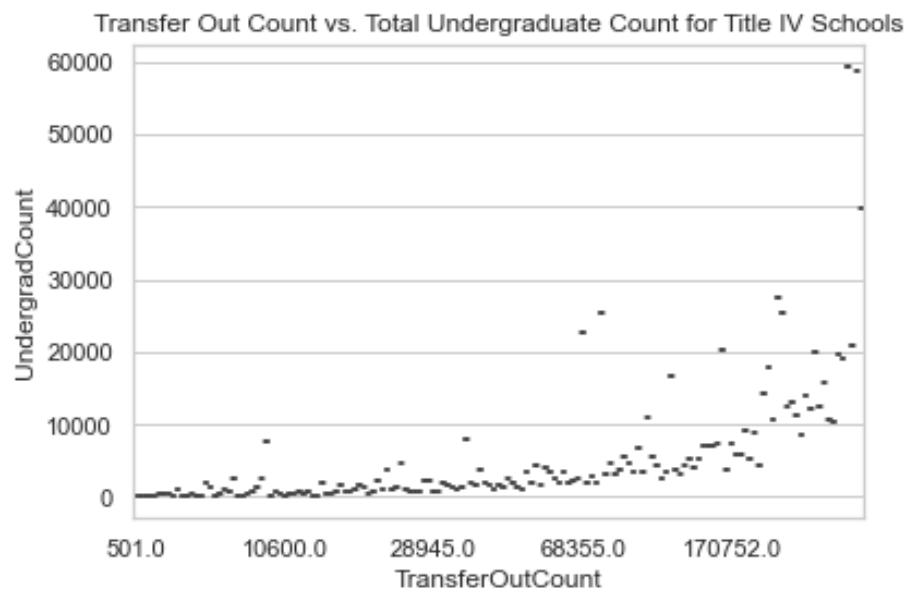
<AxesSubplot:title={'center': 'Distribution of Pacific Islander Student Graduation Rates in Predominantly White Institutions (PWIs)'}>

Distribution of Pacific Islander Student Graduation Rates in Predominantly White Institutions (PWIs)



This graph indicates the distribution of the Pacific Islander Graduation rate in institutions where white undergraduates account for greater than 1/2 of the undergraduate population. The average pacific islander graduation rate greatly varies among PWI institutions.

```
fig, ax = plt.subplots()
ax.set_title('Transfer Out Count vs. Total Undergraduate Count for Title IV Schools')
sns.boxplot(x='TransferOutCount', y='UndergradCount', data=subset)
plt.locator_params(axis='x', nbins=5)
```



This plot shows the distribution of students that transferred out versus the students that stayed in their institution.