

```
#needed libraries
import pandas as pd
import requests as r
from scipy import stats
import numpy as np
import matplotlib.pyplot as plt

#reading flat files into DataFrame object
salaries_region = pd.read_csv("datasets/salaries-by-region.csv")
print(salaries_region.head())
```

Milestone 2

Perform at least 5 data transformation and/or cleansing steps to your flat file data. For example:

1. Replace Headers
2. Format data into a more readable format
3. Identify outliers and bad data
4. Find duplicates
5. Fix casing or inconsistent values
6. Conduct Fuzzy Matching

Reading in flat files

```
In [98]: salaries_by_region.csv
#Reading Flat files into DataFrame object
salaries_region = pd.read_csv("datasets/salaries-by-region.csv")
print(salaries_region.head())

0      School Name      Region \
1  California Institute of Technology (CIT)  California
2      Harvey Mudd College  California
3  University of California, Berkeley  California
4      Occidental College  California

Starting Median Salary Mid-Career Median Salary \
0      $70,400.00      $129,000.00
1      $75,500.00      $123,000.00
2      $71,800.00      $122,000.00
3      $59,900.00      $112,000.00
4      $51,900.00      $105,000.00

Mid-Career 10th Percentile Salary Mid-Career 25th Percentile Salary \
0      $68,400.00      $93,100.00
1      $74,000.00      $104,000.00
2      NaN      $96,000.00
3      $59,500.00      $81,000.00
4      NaN      $54,800.00

Mid-Career 75th Percentile Salary Mid-Career 90th Percentile Salary \
0      $184,000.00      $229,200.00
1      $161,000.00      NaN
2      $180,000.00      NaN
3      $149,000.00      $201,000.00
4      $157,000.00      NaN
```

```
In [98]: salaries-by-college-type.csv
salaries_type = pd.read_csv("datasets/salaries-by-college-type.csv")
print(salaries_type.head())

0      School Name      School Type \
1  Massachusetts Institute of Technology (MIT)  Engineering
2  California Institute of Technology (CIT)  Engineering
3  University of California, Berkeley  Engineering
4  Polytechnic University of New York, Brooklyn  Engineering

Starting Median Salary Mid-Career Median Salary \
0      $72,200.00      $126,000.00
1      $71,800.00      $122,000.00
2      $71,800.00      $122,000.00
3      $62,400.00      $114,000.00
4      $52,200.00      $114,000.00

Mid-Career 10th Percentile Salary Mid-Career 25th Percentile Salary \
0      $76,800.00      $104,000.00
1      NaN      $104,000.00
2      NaN      $96,000.00
3      $66,800.00      $137,000.00
4      NaN      $80,200.00

Mid-Career 75th Percentile Salary Mid-Career 90th Percentile Salary \
0      $168,000.00      $220,000.00
1      $161,000.00      NaN
2      $180,000.00      NaN
3      $145,000.00      $190,000.00
4      $142,000.00      NaN
```

Merging 2 Salaries Flat Files

```
In [98]: salaries_merged = pd.merge(salaries_region, salaries_type, on='School Name')
salaries_merged.head()

Out[98]:
```

	School Name	Region	Starting Median Salary_x	Mid-Career Median Salary_x	Mid-Career 10th Percentile Salary_x	Mid-Career 25th Percentile Salary_x	Mid-Career 75th Percentile Salary_x	Mid-Career 90th Percentile Salary_x	School Type	Starting Median Salary_y	Mid-Career Median Salary_y	Mid-Career 10th Percentile Salary_y
0	Institute of Technology (CIT)	California	\$75,500.00	\$123,000.00	NaN	\$104,000.00	\$161,000.00	NaN	Engineering	\$75,500.00	\$123,000.00	NaN
1	Harvey Mudd College	California	\$71,800.00	\$122,000.00	NaN	\$96,000.00	\$180,000.00	NaN	Engineering	\$71,800.00	\$122,000.00	NaN
2	University of California, Berkeley	California	\$59,900.00	\$112,000.00	\$59,500.00	\$81,000.00	\$149,000.00	\$201,000.00	State	\$59,900.00	\$112,000.00	\$59,500.00
3	Occidental College	California	\$51,900.00	\$105,000.00	NaN	\$54,800.00	\$157,000.00	NaN	Liberal Arts	\$51,900.00	\$105,000.00	NaN
4	Cal Poly San Luis Obispo	California	\$57,200.00	\$101,000.00	\$55,000.00	\$74,700.00	\$133,000.00	\$178,000.00	State	\$57,200.00	\$101,000.00	\$55,000.00

Removing Duplicate and Unnecessary Columns

For the percentile columns, we have values for the 10th, 25th, 75th, and 90th percentiles. Given that we can find the interquartile range for a dataset by using Q3 (75th) and Q1 (25th), I don't think we will need the 10th and 90th percentile fields. The IQR will give us a good idea of statistical dispersion, and is commonly used as a robust measure of scale.

I am also going to drop the 10th and 90th percentile columns from the first DataFrame.

```
In [98]: #removing second set of salary columns as they contain same data from first dataframe
salaries_merged = salaries_merged.drop(['Mid-Career 10th Percentile Salary_x', 'Mid-Career 90th Percentile Salary_x'], axis=1)

In [98]: print(salaries_merged.head())

0      School Name      Region \
1  California Institute of Technology (CIT)  California
2      Harvey Mudd College  California
3  University of California, Berkeley  California
4      Cal Poly San Luis Obispo  California

Starting Median Salary_x Mid-Career Median Salary_x \
0      $75,500.00      $123,000.00
1      $71,800.00      $122,000.00
2      $59,900.00      $112,000.00
3      $51,900.00      $105,000.00
4      $57,200.00      $101,000.00

Mid-Career 25th Percentile Salary_x Mid-Career 75th Percentile Salary_x \
0      $104,000.00      $161,000.00
1      $96,000.00      $180,000.00
2      $81,000.00      $149,000.00
3      $54,800.00      $137,000.00
4      $74,700.00      $133,000.00

School Type
0      Engineering
1      Engineering
2      State
3      Liberal Arts
4      State
```

Rearranging columns in Merged DataFrame

```
In [98]: #retrieve columns of the merged dataframe in list form
cols = salaries_merged.columns.tolist()

Out[98]:
```

```
['School Name',
 'Region',
 'Starting Median Salary_x',
 'Mid-Career Median Salary_x',
 'Mid-Career 25th Percentile Salary_x',
 'Mid-Career 75th Percentile Salary_x',
 'School Type']
```

```
In [98]: #assign cols list in the order wanted
cols = ['School Name', 'Region', 'School Type', 'Starting Median Salary_x', 'Mid-Career Median Salary_x', 'Mid-Career 25th Percentile Salary_x', 'Mid-Career 75th Percentile Salary_x']
salaries_merged = salaries_merged[cols]

In [98]: print(salaries_merged.head())

0      School Name      Region      School Type \
1  California Institute of Technology (CIT)  California  Engineering
2      Harvey Mudd College  California  Engineering
3  University of California, Berkeley  California  State
4      Occidental College  California  Liberal Arts
5      Cal Poly San Luis Obispo  California  State

Starting Median Salary_x Mid-Career Median Salary_x \
0      $75,500.00      $123,000.00
1      $71,800.00      $122,000.00
2      $59,900.00      $112,000.00
3      $51,900.00      $105,000.00
4      $57,200.00      $101,000.00

Mid-Career 25th Percentile Salary_x Mid-Career 75th Percentile Salary_x \
0      $104,000.00      $161,000.00
1      $96,000.00      $180,000.00
2      $81,000.00      $149,000.00
3      $54,800.00      $137,000.00
4      $74,700.00      $133,000.00
```

Renaming Existing Columns

```
In [98]: salaries_merged = salaries_merged.rename(columns={'School Name':'school_name', 'School Type':'school_type', 'Starting Median Salary_x':'starting_median_salary', 'Mid-Career Median Salary_x':'midCareer_median_salary', 'Mid-Career 25th Percentile Salary_x':'midCareer_25th_salary', 'Mid-Career 75th Percentile Salary_x':'midCareer_75th_salary'})

In [98]: print(salaries_merged.head())

0      school_name      Region      school_type \
1  California Institute of Technology (CIT)  California  Engineering
2      Harvey Mudd College  California  Engineering
3  University of California, Berkeley  California  State
4      Occidental College  California  Liberal Arts
5      Cal Poly San Luis Obispo  California  State

starting_median_salary midCareer_median_salary midCareer_25th_salary \
0      $75,500.00      $123,000.00      $104,000.00
1      $71,800.00      $122,000.00      $96,000.00
2      $59,900.00      $112,000.00      $81,000.00
3      $51,900.00      $105,000.00      $54,800.00
4      $57,200.00      $101,000.00      $74,700.00

midCareer_75th_salary
0      $161,000.00
1      $180,000.00
2      $149,000.00
3      $137,000.00
4      $133,000.00
```

Checking for Duplicates

```
In [91]: #can check for duplicates using pandas duplicated function on each column

#loop through each column in dataframe
print("Checking for duplicates in DataFrame\n")
for col in salaries_merged.columns:
    print(col + ": " + str(any(salaries_merged[col].duplicated())))

Checking for duplicates in DataFrame

school_name: True
Region: True
school_type: True
starting_median_salary: True
midCareer_median_salary: True
midCareer_25th_salary: True
midCareer_75th_salary: True

There are duplicates in the 'school_name' column which is not ideal as this is technically the key for our DataFrame and each record for each School Name should be unique.
```

I am going to investigate the duplicates in this column.

```
In [91]: #looking at the duplicate values in the school_name column
salaries_merged[salaries_merged.duplicated('school_name', keep=False)]

#20 schools show up twice in the dataframe

Out[91]:
```

	school_name	Region	school_type	starting_median_salary	midCareer_median_salary	midCareer_25th_salary	midCareer_75th_salary
11	University of California, Santa Barbara (UCSB)	California	Party	\$50,500.00	\$95,000.00	\$71,200.00	\$129,000.00
12	University of California, Santa Barbara (UCSB)	California	State	\$50,500.00	\$95,000.00	\$71,200.00	\$129,000.00
34	Arizona State University (ASU)	Western	Party	\$47,400.00	\$84,100.00	\$60,700.00	\$114,000.00
35	Arizona State University (ASU)	Western	State	\$47,400.00	\$84,100.00	\$60,700.00	\$114,000.00
68	University of Illinois at Urbana-Champaign (UIUC)	Midwestern	Party	\$52,900.00	\$96,100.00	\$68,900.00	\$132,000.00
69	University of Illinois at Urbana-Champaign (UIUC)	Midwestern	State	\$52,900.00	\$96,100.00	\$68,900.00	\$132,000.00
80	Indiana University (IU)	Midwestern	Party	\$46,300.00	\$84,000.00	\$60,400.00	\$119,000.00
81	Indiana University (IU)	Midwestern	State	\$46,300.00	\$84,000.00	\$60,400.00	\$119,000.00
82	University of Iowa (UI)	Midwestern	Party	\$44,700.00	\$83,900.00	\$61,100.00	\$116,000.00
83	University of Iowa (UI)	Midwestern	State	\$44,700.00	\$83,900.00	\$61,100.00	\$116,000.00
106	Ohio State University (OSU)	Midwestern	Party	\$42,200.00	\$73,400.00	\$52,800.00	\$106,000.00
107	Ohio State University (OSU)	Midwestern	State	\$42,200.00	\$73,400.00	\$52,800.00	\$106,000.00
136	University of Maryland, College Park	Southern	Party	\$52,000.00	\$95,000.00	\$68,300.00	\$126,000.00
137	University of Maryland, College Park	Southern	State	\$52,000.00	\$95,000.00	\$68,300.00	\$126,000.00
139	University of Texas (UT) - Austin	Southern	Party	\$49,700.00	\$93,900.00	\$67,400.00	\$129,000.00
140	University of Texas (UT) - Austin	Southern	State	\$49,700.00	\$93,900.00	\$67,400.00	\$129,000.00
141	University of Florida (UF)	Southern	Party	\$47,100.00	\$87,900.00	\$62,900.00	\$120,000.00
142	University of Florida (UF)	Southern	State	\$47,100.00	\$87,900.00	\$62,900.00	\$120,000.00
143	Louisiana State University (LSU)	Southern	Party	\$46,900.00	\$87,800.00	\$61,300.00	\$120,000.00
144	Louisiana State University (LSU)	Southern	State	\$46,900.00	\$87,800.00	\$61,300.00	\$120,000.00
147	University of Georgia (UGA)	Southern	Party	\$44,100.00	\$86,000.00	\$57,800.00	\$118,000.00
148	University of Georgia (UGA)	Southern	State	\$44,100.00	\$86,000.00	\$57,800.00	\$118,000.00
151	Randolph-Macon College	Southern	Party	\$42,600.00	\$83,600.00	\$54,100.00	\$123,000.00
152	Randolph-Macon College	Southern	Liberal Arts	\$42,600.00	\$83,600.00	\$54,100.00	\$123,000.00
158	University of Alabama, Tuscaloosa	Southern	Party	\$41,300.00	\$81,400.00	\$56,500.00	\$117,000.00
159	University of Alabama, Tuscaloosa	Southern	State	\$41,300.00	\$81,400.00	\$56,500.00	\$117,000.00
164	University of Mississippi	Southern	Party	\$41,400.00	\$79,700.00	\$53,500.00	\$108,000.00
165	University of Mississippi	Southern	State	\$41,400.00	\$79,700.00	\$53,500.00	\$108,000.00
169	West Virginia University (WVU)	Southern	Party	\$43,100.00	\$78,100.00	\$55,700.00	\$106,000.00
170	West Virginia University (WVU)	Southern	State	\$43,100.00	\$78,100.00	\$55,700.00	\$106,000.00
174	University of Tennessee	Southern	Party	\$45,800.00	\$74,600.00	\$53,200.00	\$106,000.00
175	University of Tennessee	Southern	State	\$45,800.00	\$74,600.00	\$53,200.00	\$106,000.00
179	Florida State University (FSU)	Southern	Party	\$42,100.00	\$73,000.00	\$52,800.00	\$107,000.00
180	Florida State University (FSU)	Southern	State	\$42,100.00	\$73,000.00	\$52,800.00	\$107,000.00
231	University of New York at Albany (SUNY) at Albany	Northeastern	Party	\$44,500.00	\$92,200.00	\$63,100.00	\$135,000.00
232	University of New York at Albany (SUNY) at Albany	Northeastern	State	\$44,500.00	\$92,200.00	\$63,100.00	\$135,000.00
239	Pennsylvania State University (PSU)	Northeastern	Party	\$49,900.00	\$85,700.00	\$62,000.00	\$117,000.00
240	Pennsylvania State University (PSU)	Northeastern	State	\$49,900.00	\$85,700.00	\$62,000.00	\$117,000.00
253	University of New Hampshire (UNH)	Northeastern	Party	\$41,800.00	\$78,300.00	\$56,400.00	\$114,000.00
254	University of New Hampshire (UNH)	Northeastern	State	\$41,800.00	\$78,300.00	\$56,400.00	\$114,000.00

Removing Duplicates

```
In [91]: #sort dataframe by key, school_name
salaries_merged = salaries_merged.sort_values("school_name")

In [91]: #reset new dataframe and drop duplicates
#drop = 'first' --> keep the first occurrence of the duplicate and remove the other one, leave one unique record
salaries_unique = salaries_merged.drop_duplicates(subset="school_name", keep="first")

In [91]: orig_size = salaries_merged.shape
new_size = salaries_unique.shape

In [91]: print("Original size before removing duplicates: " + str(orig_size))

Original size before removing duplicates: (248, 7)

In [91]: print("New size after removing duplicates: " + str(new_size))

New size after removing duplicates: (248, 7)
```

Identifying Missing Values

```
In [91]: #loop through columns in dataframe
#check for any NaN values
for col in salaries_unique.columns:
    print(col + ": " + str(salaries_unique[col].isnull().values.any()))

#no missing values!

school_name: False
Region: False
school_type: False
starting_median_salary: False
midCareer_median_salary: False
midCareer_25th_salary: False
midCareer_75th_salary: False
```

Describing Data & Checking field types

```
In [91]: #describing our dataframe with duplicates removed
salaries_unique.describe()

Out[91]:
```

	school_name	Region	school_type	starting_median_salary	midCareer_median_salary	midCareer_25th_salary	midCareer_75th_salary	
count	248	248	248	248	248	248	248	
unique	248	5	5	145	166	178	110	
top	Western	Michigan (WVNU)	Northeastern	State	\$42,600.00	\$72,100.00	\$54,100.00	\$122,000.00
freq	1	67	164	6	5	5	7	

```
In [91]: salaries_unique.dtypes

Out[91]:
```

```
school_name      object
Region            object
school_type       object
starting_median_salary      object
midCareer_median_salary     object
midCareer_25th_salary       object
midCareer_75th_salary       object
dtype: object
```

All of the columns in our dataframe are of the Object type. We will need to cast the salary fields to be of type 'numeric', so we can properly work with them and apply functions/transformations.

Casting Variables

```
In [91]: #removing special characters from the salary fields so they can be cast to numeric
for col in ['starting_median_salary', 'midCareer_median_salary', 'midCareer_25th_salary', 'midCareer_75th_salary']:
    salaries_unique[col] = salaries_unique[col].str.replace(' ','')
    salaries_unique[col] = salaries_unique[col].str.replace('$','')
    salaries_unique[col] = salaries_unique[col].str.replace('\n','')

Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

Python-input-921-62552876a08c13: FutureWarning: The default value of regex will change from True to False in a future version.
Python-input-921-62552876a08c1
```



```
In [938]: print(df.shape)
Out[938]: (569, 9)
In [939]: print(df.head())
```

Rank	School	Student to Faculty Ratio	Graduation Rate	Retention Rate	Acceptance Rate	Enrollment Rate	Institutional Aid Rate	Default Rate
0	1	Harvard University	7 to 1	0.98	0.98	0.06	0.44	0.98
1	2	Yale University	6 to 1	0.97	0.99	0.07	0.52	0.97
2	3	University of Pennsylvania	6 to 1	0.95	0.98	0.10	0.07	0.94
3	4	Johns Hopkins University	10 to 1	0.94	0.97	0.14	0.05	0.94
4	5	Cornell University	9 to 1	0.93	0.97	0.15	0.08	0.93
Retention Rate	Acceptance Rate	Enrollment Rate	Institutional Aid Rate	Default Rate				
0	0.98	0.98	0.06	0.44	1	0.98	0.06	0.44
1	0.97	0.99	0.07	0.52	2	0.97	0.07	0.52
2	0.95	0.98	0.10	0.07	3	0.95	0.10	0.07
3	0.94	0.97	0.14	0.05	4	0.94	0.14	0.05
4	0.93	0.97	0.15	0.08	5	0.93	0.15	0.08
Default Rate								
0	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
1	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
2	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
3	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
4	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A

Rank	School	Student to Faculty Ratio	Graduation Rate	Retention Rate	Acceptance Rate	Enrollment Rate	Institutional Aid Rate	Default Rate
566	567	Touro University Worldwide	13 to 1	0.98	0.98	0.06	0.44	0.98
569	568	Unitek College	16 to 1	0.97	0.99	0.07	0.52	0.97
568	569	University of Western States	16 to 1	0.94	0.97	0.14	0.05	0.94
569	570	Virginia Baptist College	5 to 1	0.94	0.97	0.14	0.05	0.94
570	571	West Virginia Junior College-Morgantown	23 to 1	0.93	0.97	0.15	0.08	0.93
Graduation Rate	Retention Rate	Acceptance Rate	Enrollment Rate	Institutional Aid Rate	Default Rate			
566	0.98	0.98	0.06	0.44	N/A	1	0.98	0.06
567	0.97	0.99	0.07	0.52	N/A	2	0.97	0.07
568	0.95	0.98	0.10	0.07	N/A	3	0.95	0.10
569	0.94	0.97	0.14	0.05	N/A	4	0.94	0.14
570	0.93	0.97	0.15	0.08	N/A	5	0.93	0.15
Institutional Aid Rate	Default Rate							
566	0.98	0.98	0.06	0.44	N/A	1	0.98	0.06
567	0.97	0.99	0.07	0.52	N/A	2	0.97	0.07
568	0.95	0.98	0.10	0.07	N/A	3	0.95	0.10
569	0.94	0.97	0.14	0.05	N/A	4	0.94	0.14
570	0.93	0.97	0.15	0.08	N/A	5	0.93	0.15

Renaming Column Names

```
In [941]: #renaming the columns with spaces so they are easier to call/access
#ratio, rate columns
df = df.rename(columns={'Student to Faculty Ratio':'Stud_Fac_Ratio','Graduation Rate':'Grad_Rate','Retention Rate':'Reten_Rate'})
```

```
In [942]: #looking at the first of data, new columns are shown
df.iloc[0]
```

Rank	School	Stud_Fac_Ratio	Grad_Rate	Reten_Rate	Accept_Rate	Enroll_Rate	Inst_Aid_Rate	Default_Rate
0	1	Harvard University	7 to 1	0.98	0.98	0.06	0.44	0.98
1	2	Yale University	6 to 1	0.97	0.99	0.07	0.52	0.97
2	3	University of Pennsylvania	6 to 1	0.95	0.98	0.10	0.07	0.94
3	4	Johns Hopkins University	10 to 1	0.94	0.97	0.14	0.05	0.94
4	5	Cornell University	9 to 1	0.93	0.97	0.15	0.08	0.93

Finding Duplicates

```
In [943]: #using the DataFrame method duplicated to determine whether each row is a duplicate
df.duplicated()
```

```
Out[943]: 0 False
1 False
2 False
3 False
4 False
566 False
567 False
568 False
569 False
570 False
length: 571, dtype: bool
```

```
In [944]: #count the number of duplicates
df.duplicated().sum()
```

```
Out[944]: 0
```

```
In [945]: print("No duplicates were found! All unique rows were loaded into the Data Frame")
No duplicates were found! All unique rows were loaded into the Data Frame
```

Finding Missing Data

```
In [946]: #loop through columns in dataframe
#check for any NaN values
for col in df.columns:
    print(col + " : " + str(df[col].isnull().values.any()))
```

```
Rank: False
School: False
Stud_Fac_Ratio: False
Grad_Rate: False
Reten_Rate: False
Accept_Rate: False
Enroll_Rate: False
Inst_Aid_Rate: False
Default_Rate: False
Name: 0, dtype: object
```

```
In [947]: print("No missing data was found, but I can see NaN's in the dataset. Going to look into those!")
No missing data was found, but I can see NaN's in the dataset. Going to look into those!
```

Replacing N/A values with np.NaN

```
In [948]: #replacing N/A values with np.NaN so they are recognized as missing values
na = df.replace('N/A',np.NaN)
```

```
In [949]: #checking for missing data again now that NaN values are in place
for col in df.columns:
    print(col + " : " + str(df[col].isnull().values.any()))
```

```
Rank: False
School: False
Stud_Fac_Ratio: False
Grad_Rate: True
Reten_Rate: True
Accept_Rate: True
Enroll_Rate: True
Inst_Aid_Rate: False
Default_Rate: True
Name: 0, dtype: object
```

```
In [950]: print("Multiple columns have missing data: Graduation Rate, Retention Rate, Acceptance Rate, Enrollment Rate, Default Rate")
Multiple columns have missing data: Graduation Rate, Retention Rate, Acceptance Rate, Enrollment Rate, and Default Rate
```

```
In [951]: #getting counts of missing values in data frame
df.isnull().sum()
```

```
Out[951]: Rank      0
School      0
Stud_Fac_Ratio      0
Grad_Rate      6
Reten_Rate      4
Accept_Rate      29
Enroll_Rate      29
Inst_Aid_Rate      0
Default_Rate      291
dtype: int64
```

Checking & Changing Data Types

```
In [952]: #using dtypes function to find data types
df.dtypes
```

```
Out[952]: Rank      object
School      object
Stud_Fac_Ratio      object
Grad_Rate      float64
Reten_Rate      object
Accept_Rate      object
Enroll_Rate      object
Inst_Aid_Rate      object
Default_Rate      object
dtype: object
```

The rate columns are marked as type 'object'. I am going to cast them to be numeric, so we can utilize the numerical values as they represent rates in percentages.

```
In [953]: def percent_to_float(col):
df[col] = df[col].str.strip('%').astype('float')/100.0
```

```
In [954]: for col in ['Grad_Rate','Reten_Rate','Accept_Rate','Enroll_Rate','Inst_Aid_Rate','Default_Rate']:
    percent_to_float(col)
```

```
In [955]: df.dtypes
```

```
Rank      object
School      object
Stud_Fac_Ratio      object
Grad_Rate      float64
Reten_Rate      float64
Accept_Rate      float64
Enroll_Rate      float64
Inst_Aid_Rate      float64
Default_Rate      float64
dtype: object
```

```
In [956]: df.head()
```

Rank	School	Stud_Fac_Ratio	Grad_Rate	Reten_Rate	Accept_Rate	Enroll_Rate	Inst_Aid_Rate	Default_Rate
0	1	Harvard University	7 to 1	0.98	0.98	0.06	0.44	NaN
1	2	Yale University	6 to 1	0.97	0.99	0.07	0.52	NaN
2	3	University of Pennsylvania	6 to 1	0.95	0.98	0.10	0.07	0.54
3	4	Johns Hopkins University	10 to 1	0.94	0.97	0.14	0.05	0.94
4	5	Cornell University	9 to 1	0.93	0.97	0.15	0.08	0.55

Now that the rate values are classified as being of type 'float', we can replace any missing values with aggregated numerical values so as to best fit with the type of the columns.

Filling Missing Values

```
In [957]: #fill the missing Default Rate data with median value
#round median values to two decimal places for percentage purposes
def fill_na_median(data, inplace=True):
    return data.fillna(data.median(),2), inplace=inplace
```

Default Rate has the most missing values, 291 of them. This count is over half of the total size (number of rows) of the Data Frame. I am going to replace these with the median value for Default Rate, so we can handle any outliers as well.

```
In [958]: #median value for Default Rate
print("Median Default Rate: " + str(df['Default_Rate'].median()))
```

```
Median Default Rate: 0.86
```

```
In [959]: fill_na_median(df['Default_Rate'])
```

```
In [960]: df['Default_Rate']
```

```
Out[960]: 0    0.06
1    0.06
2    0.06
3    0.06
4    0.06
566    0.04
567    0.06
568    0.06
569    0.06
570    0.06
Name: Default_Rate, Length: 571, dtype: float64
```

For the other rate values, I will perform the same missing data activity by replacing the missing values with the median values for the corresponding columns.

```
In [961]: #median values
print("Median Graduation Rate: " + str(df['Grad_Rate'].median())) #0.66
print("Median Retention Rate: " + str(df['Reten_Rate'].median())) #0.83
print("Median Acceptance Rate: " + str(df['Accept_Rate'].median())) #0.65
print("Median Enrollment Rate: " + str(df['Enroll_Rate'].median())) #0.18
```

```
Median Graduation Rate: 0.66
Median Retention Rate: 0.83
Median Acceptance Rate: 0.645
Median Enrollment Rate: 0.18
```

```
In [962]: #fill in rest of missing data with median values
fill_na_median(df['Grad_Rate'])
fill_na_median(df['Reten_Rate'])
fill_na_median(df['Accept_Rate'])
fill_na_median(df['Enroll_Rate'])
```

```
In [963]: #checking to make sure we handled all missing data
#getting counts of missing values in data frame
df.isnull().sum()
```

```
Out[963]: Rank      0
School      0
Stud_Fac_Ratio      0
Grad_Rate      0
Reten_Rate      0
Accept_Rate      0
Enroll_Rate      0
Inst_Aid_Rate      0
Default_Rate      0
dtype: int64
```

Number of missing values in each column is now 0! We handled the missing data for the rate columns by replacing any missing data with the median value for the columns. The median will allow us to get a value which accounts for outliers, rather than taking the mean.

Detecting and Filtering Outliers

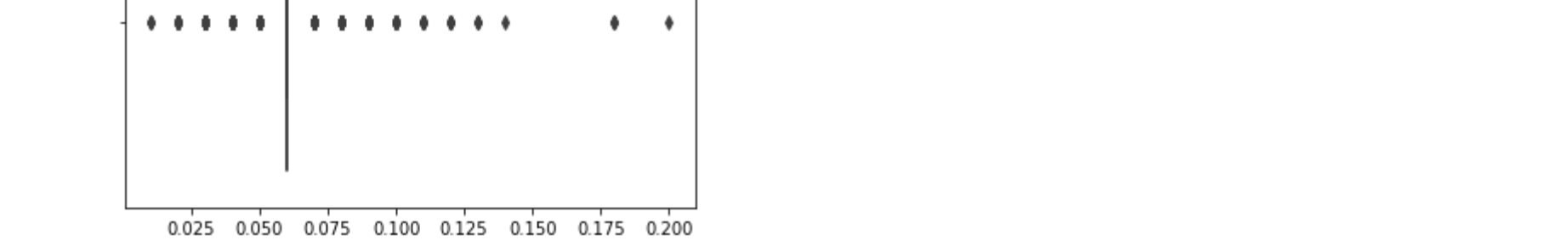
```
In [964]: #describing the data
#gives us an idea of the distribution
df.describe()
```

	Grad_Rate	Reten_Rate	Accept_Rate	Enroll_Rate	Inst_Aid_Rate	Default_Rate
count	571.000000	571.000000	571.000000	571.000000	571.000000	571.000000
mean	0.676533	0.827968	0.628546	0.201156	0.719772	0.060893
std	0.130232	0.085167	0.172070	0.120318	0.179739	0.021236
min	0.350000	0.250000	0.060000	0.040000	0.060000	0.010000
25%	0.560000	0.780000	0.530000	0.120000	0.580000	0.060000
50%	0.660000	0.830000	0.650000	0.180000	0.740000	0.060000
75%	0.740000	0.880000	0.750000	0.240000	0.870000	0.060000
max	1.000000	1.000000	1.000000	1.000000	1.000000	0.200000

Boxplots of Numerical Fields

```
In [965]: #boxplot of Graduation Rate
sns.boxplot(x=df['Grad_Rate'])
```

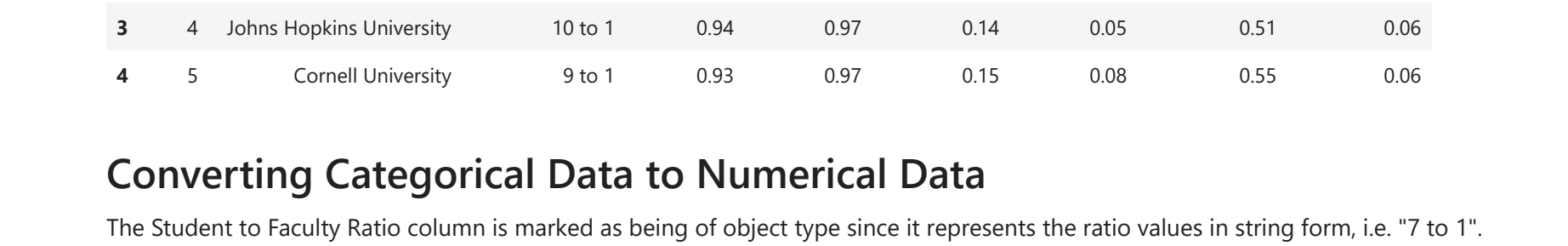
```
Out[965]: <AxesSubplot:xlabel='Grad_Rate'>
```



The Graduation Rate is relatively normally distributed. The median value (represented by the vertical line) is pretty much in the middle of the box, which represents the interquartile range for the column. I am pleased with the distribution, and there are no apparent outliers that need to be removed.

```
Retention Rate
#boxplot of Retention Rate
sns.boxplot(x=df['Reten_Rate'])
```

```
Out[966]: <AxesSubplot:xlabel='Reten_Rate'>
```



The outlier values for Retention Rate fall below ~0.65 or 65%. The minimum is 0.25 or a 25% retention rate, which I think is low for an university but it makes sense from the perspective of university.

This low retention rate can show how there are universities which experience difficulty in keeping students after freshman year, since I'm sure that is a blocker that many schools experience.

In terms of giving us information in regards to how one's university influences their post-grad salary, I think that knowing if the university they attended experiences trouble with retaining students, then it could give more influence to the student for graduating and still staying at the school.

```
Acceptance Rate
#boxplot of Acceptance Rate
sns.boxplot(x=df['Accept_Rate'])
```

```
Out[967]: <AxesSubplot:xlabel='Accept_Rate'>
```



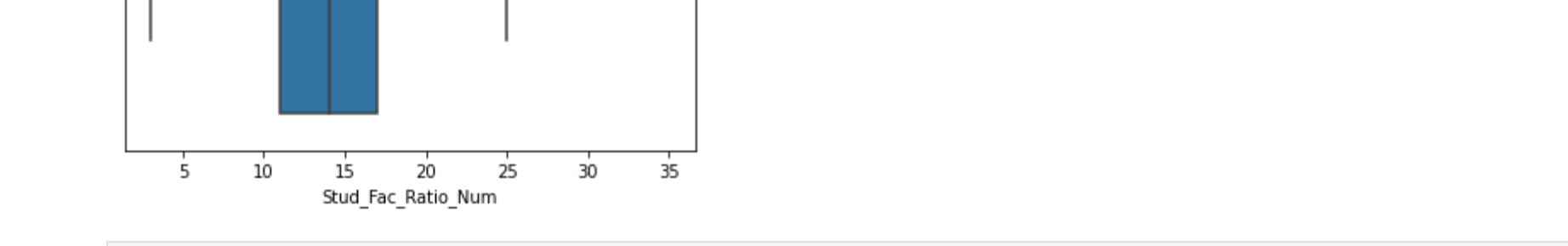
There are acceptance rate outlier values, which fall below about 0.2 or 20% acceptance rate.

Values falling below this are not surprising as many Ivy League schools have a low acceptance rate, since they are prestigious and more difficult to attend/get into.

Therefore, I will not be removing these outliers that fall below ~0.2 as they give representation to the schools that are more difficult to get accepted into given their prestige.

```
Enrollment Rates
#boxplot of Enrollment Rate
sns.boxplot(x=df['Enroll_Rate'])
```

```
Out[968]: <AxesSubplot:xlabel='Enroll_Rate'>
```



Enrollment rate represents the percentage of 18-to-24-year-olds enrolled as undergraduate or graduate students in 2- or 4- year institutions.

The outlier values for this field come above ~0.4 or 40% enrollment rate. There are universities with a higher amount of students enrolled in their courses/programs, which is good for that university!

However, for the outlier value at basically 1.0 or 100%, that is difficult for me to fathom that an university has 100% of its original students enrolled for school attendance. It would be very impressive of them, but I am viewing this university with this enrollment rate as an outlier. I am going to remove it.

```
In [969]: #finding row with enrollment rate greater than 0.8
df[df.Enroll_Rate > 0.8]
```

Rank	School	Stud_Fac_Ratio	Grad_Rate	Reten_Rate	Accept_Rate	Enroll_Rate	Inst_Aid_Rate	Default_Rate
539	540	Luther Rice University & Seminary	23 to 1	1.00	1.00	1.00	0.65	0.06
541	542	Midwives College of Utah	5 to 1	0.66	1.0	1.0	1.0	0.30

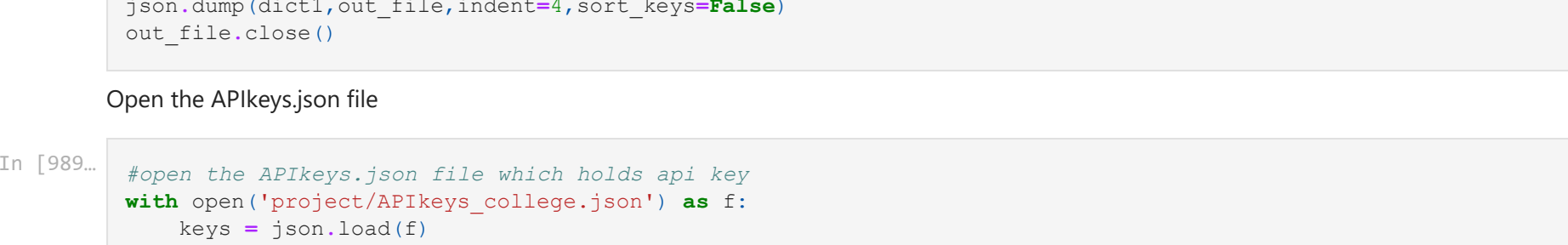
```
In [970]: #removing these two rows with 1.0 Enroll_Rate
df = df[df.Enroll_Rate <= 0.8]
```

```
In [971]: print(df.shape)
print("Now shape: 569 rows and 9 columns")
(569, 9)
Now shape: 569 rows and 9 columns
```

Institutional Aid Rate

```
In [972]: #boxplot of Institutional Aid Rate
sns.boxplot(x=df['Inst_Aid_Rate'])
```

```
Out[972]: <AxesSubplot:xlabel='Inst_Aid_Rate'>
```



Institutional Aid Rate: what percentage of college students get financial aid?

There is a designated outlier in this distribution that is less than about 0.15 or 15%. It looks like it is almost at 0.0. Let's look into the rows.

```
Out[973]: #finding row with institutional aid rate less than 0.1
df[df.Inst_Aid_Rate < 0.15]
```

Rank	School	Stud_Fac_Ratio	Grad_Rate	Reten_Rate	Accept_Rate	Enroll_Rate	Inst_Aid_Rate	Default_Rate
557	558	Mountain State College	17 to 1	0.92	0.83	0.65	0.18	0.18

The designated row actually has an Institutional Aid Rate of 6%, which means that 6% of their college students receive financial aid. This percentage is small, but it could designate that the students don't need financial aid or the university is not equipped/prepared to give it out.

In terms of our business problem of identifying salaries based on university attendance, this could designate that the students from this university are either not in need or they are using other ways to receive aid while in school.

I am going to keep this row in; its other values are representative of data that we want to look into as well.

```
Default Rate
#boxplot of Default Rate
sns.boxplot(x=df['Default_Rate'])
```

```
Out[974]: <AxesSubplot:xlabel='Default_Rate'>
```


Default rate is the percentage of all outstanding loans that a lender has written off as unpaid after a prolonged period of missed payments. In terms of how this relates to an university, it would present the percentage of graduates that default on their student loans in the first 12 months of repayment.

I think this rate value is definitely indicative of how students are doing post-graduation in terms of their incomes/finances, because paying off loans is a huge financial responsibility for students and young adults.

Since the majority of the values in this Default Rate column were replaced with the median value of 0.06, the distribution is not very well distributed and shows any values outside of the median as outliers.

In order to keep variety and distribution for the column, I will not remove any of the outliers as we need to represent the values that do not fall into the median bucket that was used for replacement.

```
In [975]: df.head()
```

Rank	School	Stud_Fac_Ratio	Grad_Rate	Reten_Rate	Accept_Rate	Enroll_Rate	Inst_Aid_Rate	Default_Rate
0	1	Harvard University	7 to 1	0.98	0.98	0.06	0.44	0.06
1	2	Yale University	6 to 1	0.97	0.99	0.07	0.52	0.06
2	3	University of Pennsylvania	6 to 1	0.95	0.98	0.10	0.07	0.54
3	4	Johns Hopkins University	10 to 1	0.94	0.97	0.14	0.05	0.94
4	5	Cornell University	9 to 1	0.93	0.97	0.15	0.08	0.55

Converting Categorical Data to Numerical Data

The Student to Faculty Ratio column is marked as being of object type since it represents the ratio values in string form, i.e. "7 to 1".

I am going to create a new column which represents the ratios in numerical form by dividing the first number in the ratio by the second number.

```
In [976]: #function to retrieve the numbers in the ratio field from around "to"
#divides the two numbers that come from the split by eachother
def convertRatio(x):
    a,b = x.split("to")
    return int(a)/int
```



```

Out[93]. # Dictionary all the desired fields
year = "latest"
#renaming columns and pulling wanted ones for df
fields =

# School Category
'School_Name': 'school.name',
'School_ID': 'id',
'School_State': 'school.state',
'School_Ownership': 'school.ownership',
'Full_Time_Faculty': 'school.ft_faculty.rate',
'Faculty_avg_sal_monthly': 'school.faculty_salary',

# Student Category
'Student_Enroll_Size': year + ".student.size",
'Student_Enroll_All': year + ".student.enrollment.all",
'percent_male_student': year + ".student.demographics.men",
'percent_female_student': year + ".student.demographics.women",
'4_yr_retention': year + ".student.retention.rate.four_year.full_time",
#Cost Category
'Attendance_Cost_per_Year': year + ".cost.attendance.academic.year",

# Completion Category
'1500_Completion_Rate_4yr': year + ".completion.completion.rate_4yr_1500a",

# Admissions Category
'Admission_Rate': year + ".admissions.admission.rate.overall",
'DAT_Avg_Overall': year + ".admissions.sat_scores.average.overall",
'DAT_75th_Percentile_Math': year + ".admissions.sat_scores.75th.percentile.math",
'DAT_75th_Percentile_Reading': year + ".admissions.sat_scores.75th.percentile.critical_reading",
'DAT_75th_Percentile_Writing': year + ".admissions.sat_scores.75th.percentile.writing",

# Earnings Category
# 6 Years after Enrollment:
'Mean_Earnings_6yrs': year + ".earnings.6_yrs_after_entry.working_not_enrolled.mean_earnings",
'Mean_male_earnings_6yrs': year + ".earnings.6_yrs_after_entry.mean_earnings_male_students",
'Mean_female_earnings_6yrs': year + ".earnings.6_yrs_after_entry.mean_earnings_female_students",
'DAT_earnings_6yrs': year + ".earnings.6_yrs_after_entry.earning_not_enrolled.std_dev",
'Percent_above_25k_6yrs': year + ".earnings.6_yrs_after_entry.percent_greater_than_25000",
'Low_income_6yrs': year + ".earnings.6_yrs_after_entry.working_not_enrolled.income.lowest_tercile",
'medium_income_6yrs': year + ".earnings.6_yrs_after_entry.working_not_enrolled.income.middle_tercile",
'high_income_6yrs': year + ".earnings.6_yrs_after_entry.working_not_enrolled.income.highest_tercile",
'Low_mean_earn_6yrs': year + ".earnings.6_yrs_after_entry.mean_earnings_lowest_tercile",
'med_mean_earn_6yrs': year + ".earnings.6_yrs_after_entry.mean_earnings_middle_tercile",
'high_mean_earn_6yrs': year + ".earnings.6_yrs_after_entry.mean_earnings_highest_tercile",

# 10 Years after Enrollment:
'Mean_Earnings_10yrs': year + ".earnings.10_yrs_after_entry.working_not_enrolled.mean_earnings",
'Mean_male_earn_10yrs': year + ".earnings.10_yrs_after_entry.mean_earnings_male_students",
'Mean_female_earn_10yrs': year + ".earnings.10_yrs_after_entry.mean_earnings_female_students",
'std_earn_10yrs': year + ".earnings.10_yrs_after_entry.working_not_enrolled.std_dev",
'Percent_above_25k_10yrs': year + ".earnings.10_yrs_after_entry.percent_greater_than_25000",
'Low_income_10yrs': year + ".earnings.10_yrs_after_entry.working_not_enrolled.income.lowest_tercile",
'medium_income_10yrs': year + ".earnings.10_yrs_after_entry.working_not_enrolled.income.middle_tercile",
'high_income_10yrs': year + ".earnings.10_yrs_after_entry.working_not_enrolled.income.highest_tercile",
'Low_mean_earn_10yrs': year + ".earnings.10_yrs_after_entry.mean_earnings_lowest_tercile",
'med_mean_earn_10yrs': year + ".earnings.10_yrs_after_entry.mean_earnings_middle_tercile",
'high_mean_earn_10yrs': year + ".earnings.10_yrs_after_entry.mean_earnings_highest_tercile"
}

# Appending all the fields values to construct the field url
fields_url = ""
for key, val in fields.items():
    fields_url = fields_url + val + ","

# To remove the extra "," at the end of fields url
fields_url = fields_url[:-1]

Out[93]. {'school.name,id,school.state,school.ownership,school.ft_faculty.rate,school.faculty_salary,latest.student.size,latest.student.enrollment.all,latest.student.demographics.men,latest.student.demographics.women,latest.student.retention.rate.four_year.full_time,latest.cost.attendance.academic.year,latest.completion.completion.rate_4yr_1500a,latest.admissions.admission.rate.overall,latest.admissions.sat_scores.average.overall,latest.admissions.sat_scores.75th.percentile.math,latest.admissions.sat_scores.75th.percentile.critical_reading,latest.admissions.sat_scores.75th.percentile.writing,latest.earnings.6_yrs_after_entry.earning_not_enrolled.mean_earnings,latest.earnings.6_yrs_after_entry.mean_earnings_male_students,latest.earnings.6_yrs_after_entry.mean_earnings_female_students,latest.earnings.6_yrs_after_entry.earning_not_enrolled.std_dev,latest.earnings.6_yrs_after_entry.percent_greater_than_25000,latest.earnings.6_yrs_after_entry.working_not_enrolled.income.lowest_tercile,latest.earnings.6_yrs_after_entry.working_not_enrolled.income.middle_tercile,latest.earnings.6_yrs_after_entry.working_not_enrolled.income.highest_tercile,latest.earnings.6_yrs_after_entry.mean_earnings_lowest_tercile,latest.earnings.6_yrs_after_entry.mean_earnings_middle_tercile,latest.earnings.6_yrs_after_entry.mean_earnings_highest_tercile,latest.earnings.10_yrs_after_entry.working_not_enrolled.mean_earnings,latest.earnings.10_yrs_after_entry.mean_earnings_male_students,latest.earnings.10_yrs_after_entry.mean_earnings_female_students,latest.earnings.10_yrs_after_entry.earning_not_enrolled.std_dev,latest.earnings.10_yrs_after_entry.percent_greater_than_25000,latest.earnings.10_yrs_after_entry.working_not_enrolled.income.lowest_tercile,latest.earnings.10_yrs_after_entry.working_not_enrolled.income.middle_tercile,latest.earnings.10_yrs_after_entry.working_not_enrolled.income.highest_tercile,latest.earnings.10_yrs_after_entry.mean_earnings_lowest_tercile,latest.earnings.10_yrs_after_entry.mean_earnings_middle_tercile,latest.earnings.10_yrs_after_entry.mean_earnings_highest_tercile'}

In[94]. # Getting number of records returned to set the max page number
base_url = f"base_url/{fields_url}(parameters.url){append}"
response = requests.get(query_url,json={})
#finding max number of pages that will need to be looked through for retrieving data
max_page_num = response["metadata"]["total"]//100 + 1
max_page_num

Out[94]. 21

Putting API Data into a DataFrame

In[95]. # Constructing the dataframe from the API request response

#initializing variables
school_df = {}
#page = 100

#looping through each page in the dataset and retrieving 100 records from each page
for page_num in range(0,max_page_num):
    query_url = f"base_url/{fields_url}(parameters.url){append}"
    max GET request to the API --> return JSON object
    response = requests.get(query_url,json={})

#retrieving records from the JSON object
for x in range(len(response["results"])):
    result_row = {}

    for key, val in fields.items():
        try:
            result_row[key] = response["results"][x][val]

```

```

except KeyError:
    print(f"Key '{key}' not found")

#append rows to the array which will become the DataFrame
school_df.append(result_row)

#create DataFrame from array variable data
school_df = pd.DataFrame(school_df)

```

In [996]:

school_df.head()

Out[996]:

	School_Name	School_ID	School_State	School_Ownership	Full_time_Faculty_Rate	Faculty_avg_sal_monthly	Stud_Enroll_Size	Stud_Enroll_Ratio
0	Alabama A & M University	100654	AL	1	0.7110	7709.0	5271.0	None
1	University of Alabama at Birmingham	100663	AL	1	0.7754	11049.0	13328.0	None
2	University of Alabama in Huntsville	100706	AL	1	0.6434	9688.0	7785.0	None
3	Alabama State University	100724	AL	1	0.6501	7221.0	3750.0	None

```

4 The University of Alabama      AL      1      0.7604      10291.0      31900.0
5 rows x 40 columns

In [997]: #checking shape of school_df
          print(school_df.shape)

          (2006, 40)

In [998]: print("2006 rows and 40 columns")

          2006 rows and 40 columns

```

```
In [100]: print(school_df.dtypes)
```

school_name	object
school_id	int64
school_state	object
school_ownership	int64
full_time_faculty_rate	float64
faculty_avg_sal_monthly	float64
stud_enroll_size	float64
stud_enroll_all	object
percent_male_stud	float64
percent_fem_stud	float64
4_yr_retention	float64
attendance_cost_per_year	float64
150%_completion_rate_4yr	float64
admission_rate	float64
avg_sas_scoreall	float64

sat_75th_percentile_math	float64
sat_75th_percentile_reading	float64
sat_75th_percentile_writing	float64
mean_earnings_6yrs	float64
mean_male_earnings_6yrs	float64
mean_fem_earnings_6yrs	float64
std_earnings_6yrs	float64
percent_above_25k_6yrs	float64
low_income_6yrs	float64
medium_income_6yrs	float64
high_income_6yrs	float64
low_mean_earn_6yrs	float64
med_mean_earn_6yrs	float64
high_mean_earn_6yrs	float64
mean_earnings_10yrs	float64
mean_male_earn_10yrs	float64
mean_fem_earn_10yrs	float64
std_earn_10yrs	float64

```

percent_above_25k_10yrs      float64
low_income_10yrs             float64
medium_income_10yrs          float64
high_income_10yrs            float64
low_mean_earn_10yrs          float64
med_mean_earn_10yrs          float64
high_mean_earn_10yrs         float64
dtype: object

```

I notice that the 'School Ownership' is defined as an int64 column. I am going to look into it, as I suspect that this is a categorical column, and it would be helpful to know what the numerical values correspond to in terms of classes.

```

In [108]: school_df['school_ownership'].unique()

Out[108]: array([1, 2, 3], dtype=int64)

```

There are three distinct numerical values for the 'School Ownership' column: 1, 2, 3. This column corresponds to the control of the institution in terms of public vs. private. I am going to create a new column which aligns the numerical values with their categories/string names.

- 1 = Public
- 2 = Private Non-Profit
- 3 = Private For-Profit

```
In [98]: # Updating School Ownership: 1: "Public", 2: "Private NonProfit", 3: "Private ForProfit"
school_df['loc[school_df["school_ownership"] == 1, "school_ownership_cat"] = "Public"
school_df['loc[school_df["school_ownership"] == 2, "school_ownership_cat"] = "Private NonProfit"
school_df['loc[school_df["school_ownership"] == 3, "school_ownership_cat"] = "Private ForProfit"]

In [100]: school_df['school_ownership_cat']
```

```
Out[100]: 0      Public
          1      Public
          2      Public
          3      Public
          4      Public
          ...
2001    Private ForProfit
2002    Private NonProfit
2003    Private NonProfit
2004    Private NonProfit
2005    Private ForProfit
Name: school_ownership_cat, Length: 2006, dtype: object

In [108]: school_df['school_ownership_cat'].unique()

Out[108]: array(['Public', 'Private NonProfit', 'Private ForProfit'], dtype=object)
```

Step #5: Finding Duplicates

```
In [100]: #using the DataFrame method duplicated to determine whether each row is a duplicate
school_df.duplicated()
```

Out[100]:	0	False
	1	False
	2	False
	3	False
	4	False
	...	
	2001	False
	2002	False
	2003	False
	2004	False

```

2005      False
Length: 2006, dtype: bool

In [100]: #using the DataFrame method duplicated to determine whether each row is a duplicate
#how many duplicate rows are there?
school_df.duplicated().sum()

Out[100]:
0

No duplicates were found in the DataFrame containing the College Scorecard API data!

Step #6: Finding Missing Data

In [100]: #loop through columns in dataframe
#check for any NaN values
for col in school_df.columns:

```

```
print(col + " " + str(school_df[col].isnull().values.any()))

school_name: False
school_id: False
school_state: False
school_ownership: False
full_time_faculty_rate: True
faculty_avg_sal_monthly: True
stud_enroll_rate: True
stud_enroll_all: True
percent_male_stud: True
percent_fem_stud: True
4_yr_retention: True
attendance_cost_per_year: True
150k_completion_rate_4yr: True
admission_rate: True
sat_avg_overall: True
sat_25th_percentile_math: True
```

```
sat_75th_percentile_reading: True
sat_75th_percentile_writing: True
mean_earnings_6yrs: True
mean_male_earning_6yrs: True
mean_fem_earning_6yrs: True
std_earning_6yrs: True
percent_above_25k_6yrs: True
low_income_6yrs: True
medium_income_6yrs: True
high_income_6yrs: True
low_mean_earn_6yrs: True
med_mean_earn_6yrs: True
high_mean_earn_6yrs: True
mean_earnings_10yrs: True
mean_male_earn_10yrs: True
mean_fem_earn_10yrs: True
std_earn_10yrs: True
percent_above_25k_10yrs: True
```

```

low_income_10yrs: True
medium_income_10yrs: True
high_income_10yrs: True
low_mean_earn_10yrs: True
med_mean_earn_10yrs: True
high_mean_earn_10yrs: True
school_ownership_cat: False

In [100]: #finding number of missing values in columns with missing data
          #getting counts of missing values in data frame
          school_df.isnull().sum()

Out[100]: school_name      0
          school_id      0
          school_state    0
          school_ownership 0
          full_time_faculty_rate    157
```

faculty_avg_sal_monthly	61
stud_enroll_size	2
stud_enroll_all	2006
percent_male_stud	2
percent_fem_stud	2
4_yr_retention	162
attendance_cost_per_year	188
1500_completion_rate_4yr	164
admission_rate	361
sat_avg_overall	770
sat_75th_percentile_math	825
sat_75th_percentile_reading	825
sat_75th_percentile_writing	1294
mean_earnings_6yrs	239
mean_male_earnings_6yrs	405
mean_fem_earnings_6yrs	405
std_earnings_6yrs	239
percent_above_25k_6yrs	256

low_income_6yrs	287
medium_income_6yrs	333
high_income_6yrs	341
low_mean_earn_6yrs	420
med_mean_earn_6yrs	491
high_mean_earn_6yrs	422
mean_earnings_10yrs	260
mean_male_earn_10yrs	433
mean_fem_earn_10yrs	433
std_earn_10yrs	260
percent_above_25k_10years	274
low_income_10yrs	305
medium_income_10yrs	346
high_income_10yrs	353
low_mean_earn_10yrs	439
med_mean_earn_10yrs	482
high_mean_earn_10yrs	440
school_ownership_cat	0

```
dtype: int64
```

Given that the column 'Student Enrollment All' has 2066 missing values, which is also equal to the number of rows in the dataset, I am going to drop this column since we won't know how to fill it and it is not providing us with any information/data.

Drop 'Student Enrollment All'

```
In [100]: school_df = school_df.drop(['stud_enroll_all'],axis=1)
```

```
In [101]: school_df.shape
```

```
Out[101]: (2066, 40)
```

```
In [101]: print("One less column in the DataFrame ... Student Enrollment All has been removed")
```

One less column in the DataFrame ... Student Enrollment All has been removed

Drop 'SAT 75th Percentile Critical Writing'

This field has 1,294 missing values. There are 2006 total rows in the DataFrame, which means that about 65% of the entire column is missing.

Since we already have columns for representing other SAT percentiles and also the overall average, I think it is okay to drop the specific column for writing, since it will not tell us much about the scores for the students at the universities if it mostly contains a replacement value anyway.

```
In [101]: school_df = school_df.drop(['sat_75th_percentile_writing'],axis=1)

In [101]: school_df.shape
```

```
Out[10]: (2006, 39)
```

Step #7: Filling Numerical Missing Values

From looking over the numerical (float) values currently in the DataFrame, most of them pertain to some 'mean' value that has already been calculated for the universities. Therefore, for filling the missing values in these columns, I'm going to take the mean of all of the values and use this as the filler. This will help align with the measurement of the values already in the dataset.

```
In [10]: #fill the missing Default Rate data with mean value  
#round mean values to two decimal places for percentage purposes  
def fill_na_mean(data, inplace=True):  
    return data.fillna(round(data.mean(),2), inplace=inplace)
```

```
In [10]:
```

```
for col in school_df[school_df.columns[school_df.isnull().any()]].columns:
    print("Col: " + col)
    fill_na_mean(school_df(col))

Col: full_time_faculty_rate
Col: faculty_avg_sal_monthly
Col: stud_enroll_size
Col: percent_male_stud
Col: percent_fem_stud
Col: 4_yr_retention
Col: attendance_cost_per_year
Col: 150k_completion_rate_4yr
Col: admission_rate
Col: sat_avg_overall
Col: sat_75th_percentile_math
Col: sat_75th_percentile_reading
Col: mean_earnings_4yrs
Col: mean_male_earning_4yrs
```

```
Col: mean_fem_earn_6yrs
Col: std_earn_6yrs
Col: percent_above_25k_6yrs
Col: low_income_6yrs
Col: medium_income_6yrs
Col: high_income_6yrs
Col: low_mean_earn_6yrs
Col: med_mean_earn_6yrs
Col: high_mean_earn_6yrs
Col: mean_earnings_10yrs
Col: mean_male_earn_10yrs
Col: mean_fem_earn_10yrs
Col: std_earn_10yrs
Col: percent_above_25k_10years
Col: low_income_10yrs
Col: medium_income_10yrs
Col: high_income_10yrs
Col: low_mean_earn_10yrs
```

```

Col: med_mean_earn_10yrs
Col: high_mean_earn_10yrs

In [101]: #check for any missing values again --> verify that they were handled accordingly
          school_df.isnull().sum()

Out[101]: school_name      0
          school_id      0
          school_state    0
          school_ownership 0
          full_time_faculty_rate 0
          faculty_avg_sal_monthly 0
          stud_enroll_size 0
          percent_male_stud 0
          percent_fem_stud 0
          4_yr_retention 0
          attendance_cost_per_year 0

```

150th_completion_rate_4yrs	0
admission_rate	0
sat_avg_overall	0
sat_75th_percentile_math	0
sat_75th_percentile_reading	0
mean_earnings_6yrs	0
mean_male_earning_6yrs	0
mean_fem_earning_6yrs	0
std_earning_6yrs	0
percent_above_25k_6yrs	0
low_income_6yrs	0
medium_income_6yrs	0
high_income_6yrs	0
low_mean_earn_6yrs	0
med_mean_earn_6yrs	0
high_mean_earn_6yrs	0
mean_earnings_10yrs	0
mean_male_earn_10yrs	0

```
mean_fem_earn_10yrs      0
std_earn_10yrs           0
percent_above_25k_10years 0
low_income_10yrs         0
medium_income_10yrs      0
high_income_10yrs        0
low_mean_earn_10yrs      0
med_mean_earn_10yrs      0
high_mean_earn_10yrs     0
school_ownership_cat      0
dtype: int64
```

There is no more missing data in school_dfl

Step #8: Detecting and Filtering Outliers

```
In [101]: #describing the data
```

school_df.describe()								
	school_id	school_ownership	full_time_faculty_rate	faculty_avg_sal_monthly	stud_enroll_size	percent_male_stud	percent_fem_stud	
count	2006.000000	2006.000000	2006.000000	2006.000000	2006.000000	2006.000000	2006.000000	
mean	22184.706879	1.779661	0.642602	7700.923290	4493.455593	0.438067	0.524541	
std	12553.384822	0.572264	0.284430	2651.030474	7926.812405	0.186746	0.189626	
min	106564.000000	1.000000	0.000000	527.000000	0.000000	0.000000	0.000000	
5%	15628.250000	1.000000	0.465175	6115.250000	636.250000	0.357600	0.496000	
50%	196683.000000	2.000000	0.640000	7520.500000	1683.000000	0.429550	0.568250	
75%	228881.750000	2.000000	0.867075	9058.750000	4539.500000	0.500725	0.618500	
max	494807.000000	3.000000	1.000000	20988.000000	98630.000000	1.000000	1.000000	

```
8 rows x 36 columns
```

```
In [101]: #plotting histograms of numerical columns
school_df.hist(bins=30, figsize=(15, 10))

Out[101]: array([[<AxesSubplot:title='center': 'school_id'>,
<AxesSubplot:title='center': 'school_ownership'>,
<AxesSubplot:title='center': 'full_time_faculty_rate'>,
<AxesSubplot:title='center': 'faculty_avg_sal_monthly'>,
<AxesSubplot:title='center': 'student_enroll_size'>,
<AxesSubplot:title='center': 'percent_male_student'>],
[<AxesSubplot:title='center': 'percent_graduated'>,
<AxesSubplot:title='center': 'student_graduated'>,
<AxesSubplot:title='center': 'avg_year_retention'>,
<AxesSubplot:title='center': 'attendance_cost_per_year'>,
<AxesSubplot:title='center': '150k_completion_rate_4yr'>],
```

```
<AxeSubPlot.title>'center': 'admission_rate',>
<AxeSubPlot.title>'center': 'sat_avg_overall',>
<AxeSubPlot.title>'center': 'sat_75th_percentile_math',>
<AxeSubPlot.title>'center': 'sat_75th_percentile_reading',>
<AxeSubPlot.title>'center': 'mean_earnings_9yrs',>
<AxeSubPlot.title>'center': 'mean_male_earnings_9yrs',>
<AxeSubPlot.title>'center': 'mean_fem_earnings_9yrs',>
<AxeSubPlot.title>'center': 'std_earnings_9yrs',>
<AxeSubPlot.title>'center': 'percent_above_25k_9yrs',>
<AxeSubPlot.title>'center': 'low_income_9yrs',>
<AxeSubPlot.title>'center': 'medium_income_9yrs',>
<AxeSubPlot.title>'center': 'high_income_9yrs',>
<AxeSubPlot.title>'center': 'low_mean_earn_9yrs',>
<AxeSubPlot.title>'center': 'med_mean_earn_9yrs',>
<AxeSubPlot.title>'center': 'high_mean_earn_9yrs',>
<AxeSubPlot.title>'center': 'mean_earnings_10yrs',>
<AxeSubPlot.title>'center': 'mean_male_earn_10yrs',>
<AxeSubPlot.title>'center': 'mean_fem_earn_10yrs',>
```

```

<AxesSubplot:title='center':std_earn_10yrs'>,<br>
<AxesSubplot:title='center':percent_above_20k_10yrs'>,<br>
<AxesSubplot:title='center':low_income_10yrs'>,<br>
<AxesSubplot:title='center':medium_income_10yrs'>,<br>
<AxesSubplot:title='center':high_income_10yrs'>,<br>
<AxesSubplot:title='center':low_mean_earn_10yrs'>,<br>
<AxesSubplot:title='center':med_mean_earn_10yrs'>,<br>
<AxesSubplot:title='center':high_mean_earn_10yrs'>]],<br>
dtype=object)

```

Looking at Outliers in percent_male_stud and percent_female_stud

```
In [10]: #boxplot of percent_male_stud

#Percentage of Male Students
sns.boxplot(smschool_df['percent_male_stud'])

Out[10]: <AxesSubplot:label='percent_male_stud'>
```

```
In [102]: #boxplot of percent_female_stud

#Percentage of Female Students
sns.boxplot(saschon_df["percent_fem_stud"])

Out[102]: <AxesSubplot: xlabel='percent_fem_stud'>
```

From analyzing these boxplots of the fields 'percent_male_stud' and 'percent_fem_stud', there are outliers at the extremes of both percentage scales. For the percentage of male students, I notice distinct outliers above ~.92. For the percentage of female students, I notice distinct outliers less than ~0.05.

Therefore, universities with almost all male students are outliers in this dataset, and universities with barely any female datasets are outliers in this dataset. These correspond with eachother!



There are single-sex universities, and I am sure they boast some advantage to receiving this type of education in terms of gender in classes. I don't want to remove any of these from the dataset, because it could give insights into how one's peer demographics influence their

education and "success" in college.

Closer Look into Distribution of 'Attendance_Cost_per_Year'

```
In [102]: sns.boxplot[school_df['attendance_cost_per_year']]
```

Out[102]: <AxesSubplot: xlabel='attendance_cost_per_year'>



A histogram showing the distribution of 'attendance_cost_per_year'. The x-axis is labeled 'attendance_cost_per_year' and ranges from 10000 to 80000. The y-axis represents frequency. The distribution is highly right-skewed, with a very high frequency at zero (the spike) and a long tail extending to the right. The bars are blue.

No distinct outliers for this field! I wanted to take a closer look since there was a spike in the histogram which seemed out of place for the distribution of the rest of the values. I think this spike is due to the replacement of missing values at the mean, which would lead to such a central focus in the spread of the data

Closer Look into Post- 10 Year Earnings

```
#boxplot for the lowest tercile of the mean earnings 10 years post-grad
#contains the lowest third of the population
```

```
sns.boxplot(x=school_id, y='low_mean_earn_10yrs')

Out[102]:
```

```

25000    50000    75000    100000    125000    150000    175000
low_mean_earn_10yrs

In [102]: #boxplot for the middle/medium tertile of the mean earnings 10 years post-grad
#middle third of the earnings population
sns.boxplot(x=school_dfl['med_mean_earn_10yrs'])

Out[102]: <AxesSubplot:label='med_mean_earn_10yrs'>

```

```
In [102]: #boxplot for the highest tercile of the mean earnings 10 years post-grad
#highest third of the earnings population
sns.boxplot(x=school_dfl["high_mean_earn_10yrs"])

Out[102]: <AxesSubplot: xlabel='high_mean_earn_10yrs'>
```

Terciles represent either of the two points that divide an ordered distribution into three parts, each containing a third of the population. In the case for the mean earnings 10-years post graduation, the distribution was divided into three parts: lowest, middle and highest mean

earnings.

There is an overlap in the means for the terciles and also their interquartile ranges, which represents that the groups are not very different from each other.

In terms of outliers, there are quite a lot of them in all three of the terciles/variables. Most are above about \$70,000. It is difficult with earnings, and this is also what I'm trying to investigate, because people and careers make such different levels of money. It is highly dependent on many factors: industry, career role, location, level in a company, etc. Therefore, I don't want to remove any of the outliers, since these could be key in showing how one's university/education influences the range of earnings/salary that one can make post-grad.

Milestone 5

Now that you have cleaned and transformed your 3 datasets, you need to load them into a database. You can choose what kind of database (SQLite or MySQL, Postgre SQL are all free options). You will want to load each dataset into SQLite as an individual table and

then you must join the datasets together in Python into 1 dataset.

Once all the data is merged together in your database, create 5 visualizations that demonstrate the data you have cleaned. You should have at least 2 visualizations that have data from more than one source (meaning, if you have 3 tables, you must have visualizations that span across 2 of the tables - you are also welcome to use your consolidated dataset that you created in the previous step, if you do that, you have met this requirement).

For the visualization portion of the project, you are welcome to use a python library like Matplotlib, Seaborn, or an R package ggPlot2, Plotly, or Tableau/PowerBI.

PowerBI is a free tool that could be used – Tableau only has a free web author. If your use Tableau/PowerBI you need to submit a PDF with your assignment vs the Tableau/PowerBI file.

Clearly label each visualization. Submit your code for merging and storing the data, with your code for the visualizations along with a 250-500 word summary of what you learned and how it helped in completing the project. In your write up, make sure to address the ethical

implications of cleansing data and your project topic. You can submit a [Jupyter Notebook](#) or a PDF of your code. If you submit a *py* file you need to also include a PDF or attachment of your results.

Column Names in Each Dataset

`salaries_unique`

```
In [102]: salaries_unique.columns

Out[102]: Index(['school_name', 'Region', 'school_type', 'starting_median_salary',
            'midCareer_median_salary', 'midCareer_25th_salary',
            'midCareer_75th_salary'],
            dtype='object')
```

df

```
In [102]: df.columns

Out[102]: Index(['Rank', 'School', 'Stud_Fac_Ratio', 'Stud_Fac_Ratio_Num', 'Grad_Rate',
        'Access_Rate', 'Accept_Rate', 'Enroll_Rate', 'Inst_Aid_Rate',
        'Default_Rate',
        dtype='object'])

school_df

In [102]: school_df.columns

Out[102]: Index(['school_name', 'school_id', 'school_state', 'school_ownership',
        'full_time_faculty_rate', 'faculty_avg_sal_monthly', 'stud_enroll_size',
        'percent_male_stud', 'percent_fom_stud', '4_yr_completion',
        'attendance_cost_per_year', '150% completion_rate_4yr',
        'admission_rate', 'sat_avg_overall', 'sat_75th_percentile_math',
        'sat_75th_percentile_reading', 'mean_senior_gpa'])
```

```

'mean_male_earning_6yrs', 'mean_fem_earning_6yrs', 'std_earning_6yrs',
'percent_above_25k_6yrs', 'low_income_6yrs', 'medium_income_6yrs',
'high_income_6yrs', 'low_mean_earn_6yrs', 'med_mean_earn_6yrs',
'high_mean_earn_6yrs', 'mean_earnings_10yrs', 'mean_male_earn_10yrs',
'mean_fem_earn_10yrs', 'std_earn_10yrs', 'percent_above_25k_10yrs',
'low_income_10yrs', 'medium_income_10yrs', 'high_income_10yrs',
'low_mean_earn_10yrs', 'med_mean_earn_10yrs', 'high_mean_earn_10yrs',
'school_ownership_cat'],
dtype='object')

```

Shapes of Each Dataset

```

In [102]: salaries_unique.shape

Out[102]: (248, 7)

```


In [103]:

```
#import sqlite3 library
import sqlite3 as sqli
```

salaries_unique DataFrame --> salaries-by-region.csv

```
#open a connection to a new database for all data
conn = sqli.connect('univ_sal.db')
```

```
#create a new table in the univ_data database for CSV data
#salaries_unique.to_sql('grad_salaries',conn)
```

```
#testing table's connection
grad_salaries = school_sqli('SELECT * FROM grad_salaries',conn)
```

```
grad_salaries.head()
```

	index	school_name	Region	school_type	starting_median_salary	midCareer_median_salary	midCareer_25th_salary	midCareer_75th_salary	
	0	214	Amherst College	Northeastern	Liberal Arts	54500	107000	84900	16200
	1	189	Appalachian State University	Southern	State	40400	69100	50400	9080
	2	34	Arizona State University	Western	Party	47400	84100	60700	11400
	3	194	Arkansas State University	Southern	State	38700	63300	45300	8390
	4	149	Auburn University	Southern	State	45400	84700	62700	10900

```
df --> 'https://oedb.org/rankings/acceptance-rate/'

# create new table in database for website data
# df_dropped.to_sql('univ_rates',conn)

school_df --> API data "https://api.data.gov/ed/collegescorecard/v1/schools?"
```

df --> "https://oeb.org/rankings/acceptance-rate/"

```
#create new table in database for website data
#df_dropped.to_sql('univ_rates',conn)
```

school_df --> API data "https://api.data.gov/ed/collegescorecard/v1/schools?"

```
#create new table in database for API data
#school_df.to_sql('schools',conn)
```

Merging 3 Tables into One Dataset

```
#initiate cursor for executing SQL
cur = conn.cursor()
```

```
#takes the schools dataset (API data) and joins all its rows to the rates table
#each university should have one corresponding rank and rates information
#join universities and their ranks to the information on salaries as only certain ones were included in the CSV
Four.execute("""INSERT INTO grad_salaries
FROM grad_salaries
LEFT JOIN schools ON schools.school_name = grad_salaries.school_name
LEFT JOIN univ_rates ON univ_rates.School = grad_salaries.school_name
LEFT JOIN grad_salaries ON univ_rates.School = grad_salaries.school_name
""")
```

```
cur.execute("""SELECT *
FROM grad_salaries
LEFT JOIN schools ON schools.school_name = grad_salaries.school_name
LEFT JOIN univ_rates ON univ_rates.School = schools.school_name
""")
```

<sqlite3.Cursor at 0x24901d550a>

```
data = pd.DataFrame(cur.fetchall()) #converts SQL query results into dataframe format
data.columns = [x[0] for x in cur.description] #labels the columns of the dataframe
data
```

245	44	Whitman College	Western	Liberal Arts	43500	80100	64800	111
246	220	Williams College	Northeastern	Liberal Arts	51700	102000	76400	143
247	96	Wittenberg University	Midwestern	Liberal Arts	39200	78200	54100	131
248	206	Worcester Polytechnic Institute	Northeastern	Engineering	61000	114000	91200	137
249	201	Yale University	Northeastern	Ivy League	59100	126000	80600	198

250 rows x 58 columns

```
#converting 'Rank' to be of float type
data['Rank'] = data['Rank'].astype(float, errors = 'raise')
```

Visualizations

```
import matplotlib.pyplot as plt
import seaborn as sns
```

Handling Column Names in Merged Dataset

```
#renaming 'school_name' columns to be different
cols = []
count = 1
for column in data.columns:
    if column == 'school_name':
        cols.append(f'school_name_{count}')
        count+=1
        continue
    cols.append(column)
data.columns = cols
```

```
data.school_name_1
```

0	Amherst College
1	Appalachian State University
2	Arizona State University
3	Arkansas State University
4	Auburn University
...	...
245	Whitman College
246	Williams College
247	Wittenberg University
248	Worcester Polytechnic Institute
249	Yale University

Name: school_name_1, Length: 250, dtype: object

250 rows x 58 columns

```
#converting 'Rank' to be of float type
data['Rank'] = data['Rank'].astype(float, errors = 'raise')
```

Visualizations

```
import matplotlib.pyplot as plt
import seaborn as sns
```

Handling Column Names in Merged Dataset

```
#renaming 'school_name' columns to be different
cols = []
count = 1
```

```
for column in data.columns:
    if column == 'school_name':
        cols.append(f'school_name_{count}')
        count+=1
    continue
cols.append(column)
data.columns = cols
```

data.school_name_1

```
0
1
2
3
4
...
245
246
247
248
249
Name: school_name_1, Length: 250, dtype: object
```

```
data.head()
```

```
Visualization #1: Scatter Plot Between Graduation Rate and Starting Median Salary
```

```
x = data['Grad_Rate']
y = data['starting_median_salary']

plt.scatter(x,y)
plt.xlabel('Graduation Rate')
plt.ylabel('Starting Median Salary')
plt.title('Starting Median Salary vs. University Graduation Rate')
plt.show()
```

5 rows x 58 columns

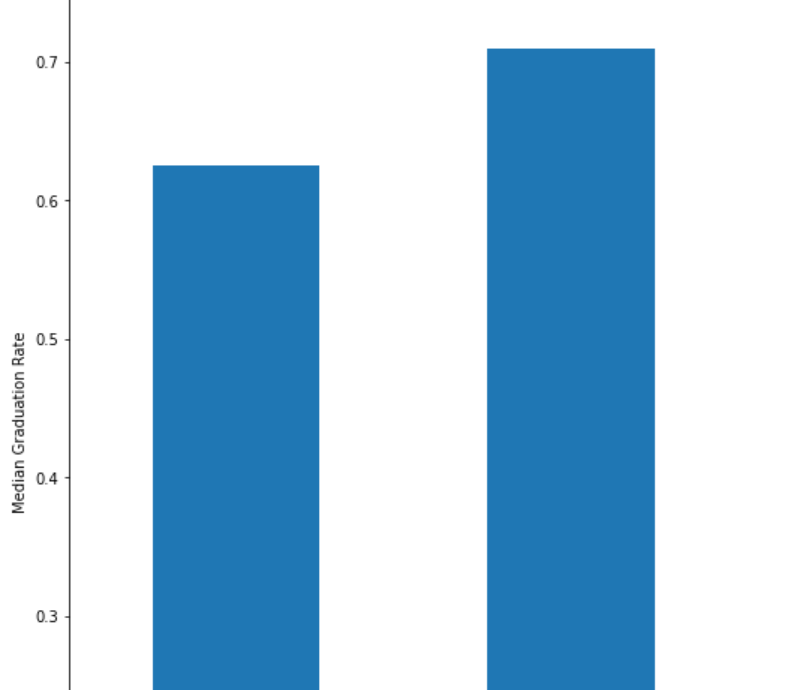
```
data.columns
```

```
Index(['index', 'school_name_1', 'Region', 'school_type',
'starting_median_salary', 'midCareer_median_salary',
'midCareer_25th_salary', 'midCareer_75th_salary', 'index',
'school_name_2', 'school_id', 'school_state', 'school_ownership',
'full_time_faculty_rate', 'faculty_avg_sal_monthly', 'stud_enroll_size',
'percent_male_stud', 'percent_fem_stud', '4_yr_retention',
'attendance_cost_per_year', '150% completion_rate_4yr',
'admission_rate', 'sat_avg_overall', 'sat_75th_percentile_math',
'sat_75th_percentile_reading', 'mean_earnings_6yrs',
'mean_male_earning_6yrs', 'mean_fem_earning_6yrs', 'std_earning_6yrs',
'percent_above_20k_6yrs', 'low_income_6yrs', 'medium_income_6yrs',
'high_income_6yrs', 'low_mean_earn_6yrs', 'med_mean_earn_6yrs',
'high_mean_earn_6yrs', 'mean_earnings_10yrs', 'mean_male_earn_10yrs',
'mean_fem_earn_10yrs', 'std_earn_10yrs', 'percent_above_20k_10years',
'low_income_10yrs', 'medium_income_10yrs', 'high_income_10yrs',
'low_mean_earn_10yrs', 'med_mean_earn_10yrs', 'high_mean_earn_10yrs',
'school_ownership_cat', 'index', 'Rank', 'School', 'Stud_Fac_Ratio_Num',
'Grad_Rate', 'Reten_Rate', 'Accept_Rate', 'Enroll_Rate',
'Inst_Aid_Rate', 'Default_Rate'],
dtype='object')
```

Visualization #1: Scatter Plot Between Graduation Rate and Starting Median Salary

```
x = data['Grad_Rate']
y = data['starting_median_salary']
```

```
plt.scatter(x,y)
plt.xlabel('Graduation Rate')
plt.ylabel('Starting Median Salary')
plt.title('Starting Median Salary vs. University Graduation Rate')
plt.show()
```



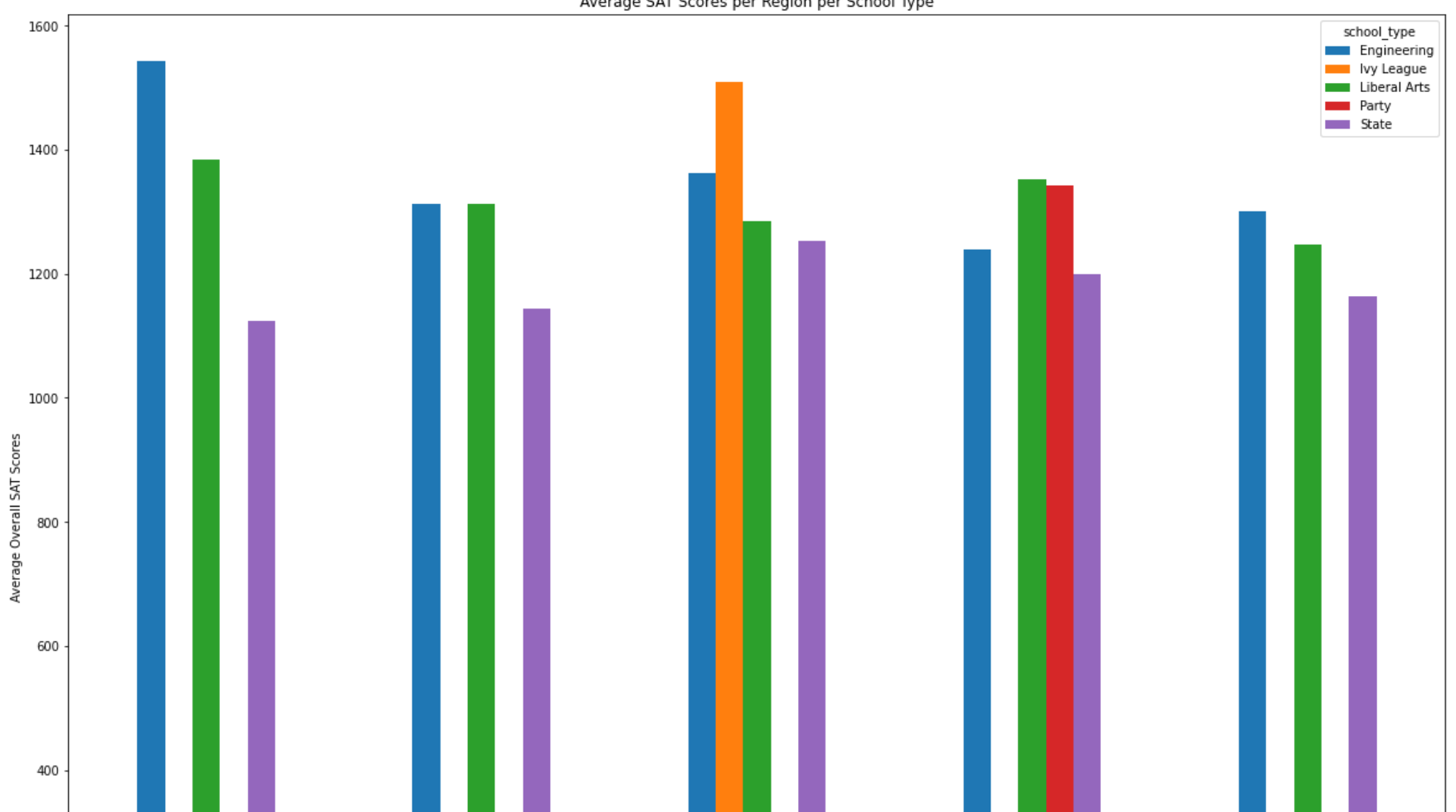
Visualization #2 - Median Graduation Rates Per Region Bar Plot

```
#number of universities per region --> minimizing effect of sample size
data.groupby(['Region'])['Grad_Rate'].median()
```

```
Region
California    0.625
Midwestern    0.710
Northeastern  0.840
Southern      0.660
Western       0.640
Name: Grad_Rate, dtype: float64
```

```
#side-by-side bar chart grouped on region and school type
#Values shown are average SAT Overall
data.groupby(['Region','school_type'])['sat_avg_overall'].mean().unstack().plot(kind='bar', stacked=False,figsize=(20,15))
plt.ylabel('Average Overall SAT Scores')
plt.xlabel('Region')
plt.title('Median Graduation Rate by Region of U.S. Universities')
```

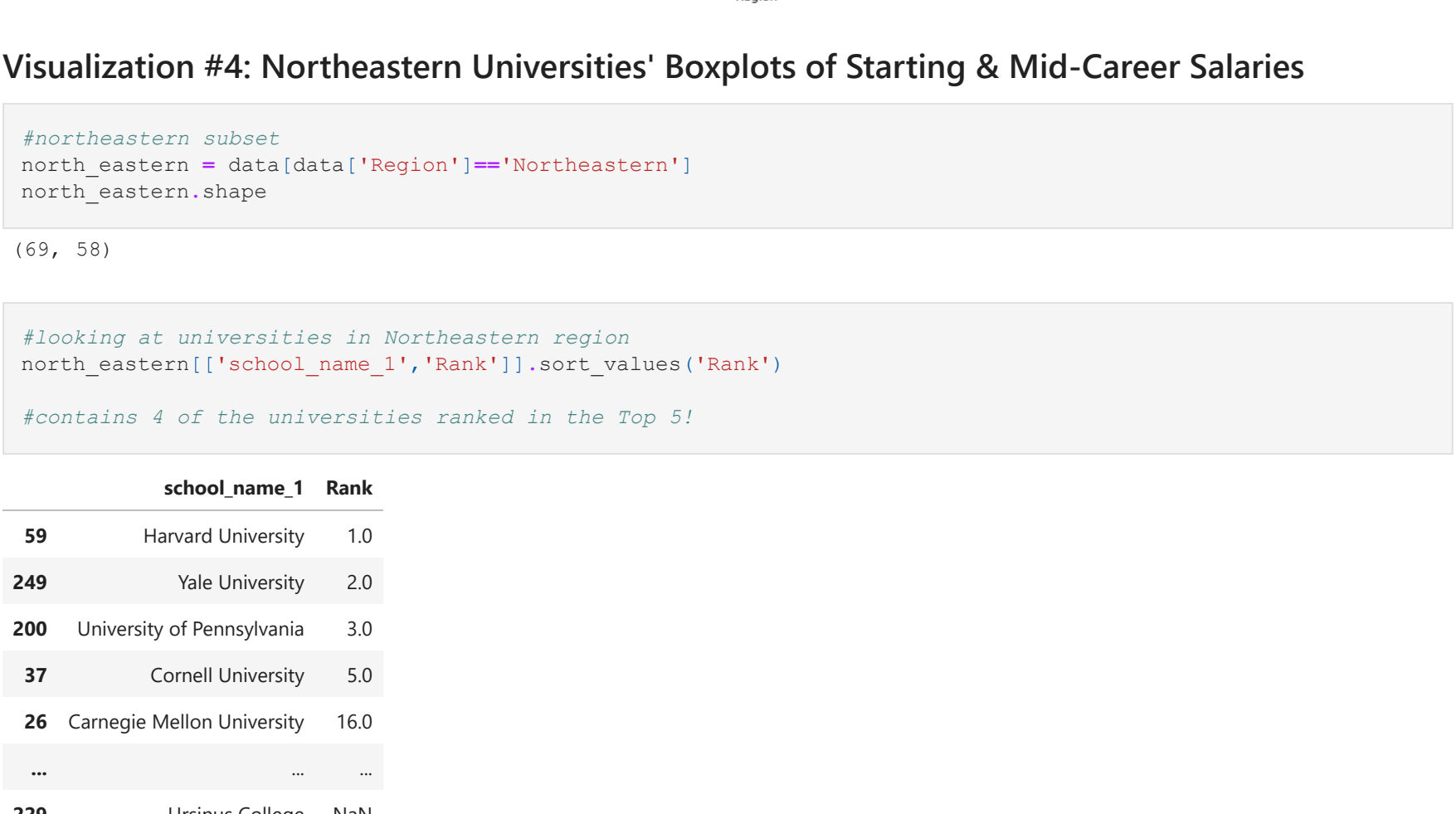
Text (0.5, 1.0, 'Median Graduation Rate by Region of U.S. Universities')



Visualization #3: Average Overall SAT Scores per Region per School Type

```
#side-by-side bar chart grouped on region and school type
#Values shown are average SAT Overall
data.groupby(['Region','school_type'])['sat_avg_overall'].mean().unstack().plot(kind='bar', stacked=False,figsize=(20,15))
plt.ylabel('Average Overall SAT Scores')
plt.xlabel('Region')
plt.title('Average SAT Scores per Region per School Type')
```

Text (0.5, 1.0, 'Average SAT Scores per Region per School Type')



Visualization #4: Northeastern Universities' Boxplots of Starting & Mid-Career Salaries

```
#northeastern subset
north_eastern = data[data['Region']=='Northeastern']
north_eastern.shape
```

```
(69, 58)
```

```
#looking at universities in Northeastern region
north_eastern[['school_name_1','Rank']].sort_values('Rank')
#contains 4 of the universities ranked in the Top 5!
```

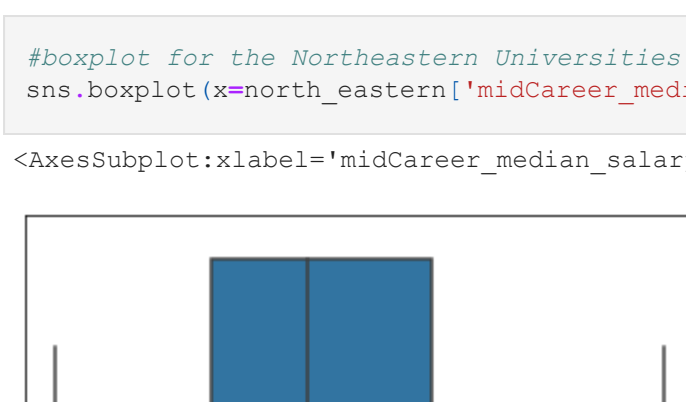
```
school_name_1 Rank
59 Harvard University 1.0
249 Yale University 2.0
200 University of Pennsylvania 3.0
37 Cornell University 5.0
26 Carnegie Mellon University 16.0
```

```
... ..
229 Ursinus College NaN
232 Vassar College NaN
238 Wellesley College NaN
240 Wesleyan University NaN
246 Williams College NaN
```

69 rows x 2 columns

```
#boxplot for the Northeastern Universities' students' starting median salaries
sns.boxplot(x=north_eastern['starting_median_salary'])
```

<AxesSubplot: xlabel='starting_median_salary'>



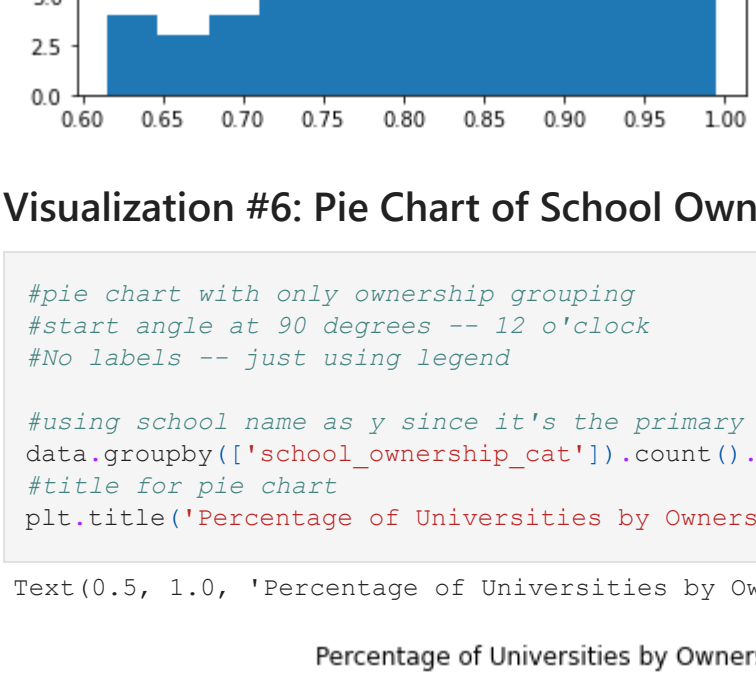
```
#boxplot for the Northeastern Universities' students' mid-career median salaries
sns.boxplot(x=north_eastern['midCareer_median_salary'])
```

<AxesSubplot: xlabel='midCareer_median_salary'>



Visualization #5: Histogram of 4-Year Retention Rates

```
x = data['4_yr_retention']
plt.hist(x,bins=12)
plt.show()
```



Visualization #6: Pie Chart of School Ownership Categories

```
#pie chart with only ownership grouping
#start angle at 90 degrees -- 12 o'clock
#No labels -- just using legend
#using school name as y since it's the primary key and it holds a value for each survey respondent
data.groupby(['school_ownership_cat']).count().plot(kind='pie',y=school_name_1,figsize=(20,10), autopct='%1.1f%%')
plt.title('Percentage of Universities by Ownership Type')
```

Text (0.5, 1.0, 'Percentage of Universities by Ownership Type')



In []: