
K Means Clustering Project -

For this project we will attempt to use KMeans Clustering to cluster Universities into two groups, Private and Public.

It is **very important to note, we actually have the labels for this data set, but we will NOT use them for the KMeans clustering algorithm, since that is an unsupervised learning algorithm.**

When using the Kmeans algorithm under normal circumstances, it is because you don't have labels. In this case we will use the labels to try to get an idea of how well the algorithm performed, but you won't usually do this for Kmeans, so the classification report and confusion matrix at the end of this project, don't truly make sense in a real world setting!.

The Data

We will use a data frame with 777 observations on the following 18 variables.

- Private A factor with levels No and Yes indicating private or public university
- Apps Number of applications received
- Accept Number of applications accepted
- Enroll Number of new students enrolled
- Top10perc Pct. new students from top 10% of H.S. class
- Top25perc Pct. new students from top 25% of H.S. class
- F.Undergrad Number of fulltime undergraduates
- P.Undergrad Number of parttime undergraduates
- Outstate Out-of-state tuition
- Room.Board Room and board costs
- Books Estimated book costs
- Personal Estimated personal spending
- PhD Pct. of faculty with Ph.D.'s
- Terminal Pct. of faculty with terminal degree
- S.F.Ratio Student/faculty ratio
- perc.alumni Pct. alumni who donate
- Expend Instructional expenditure per student
- Grad.Rate Graduation rate

Import Libraries

Import the libraries you usually use for data analysis.

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

Get the Data

Read in the `College_Data` file using `read_csv`. Figure out how to set the first column as the index.

```
In [2]: df = pd.read_csv('College_Data', index_col=0)
```

Check the head of the data

```
In [3]: df.head()
```

Out[3]:

	Private	Apps	Accept	Enroll	Top10perc	Top25perc	F.Undergrad	P.Undergrad	Outsta
Abilene Christian University	Yes	1660	1232	721	23	52	2885	537	744
Adelphi University	Yes	2186	1924	512	16	29	2683	1227	1228
Adrian College	Yes	1428	1097	336	22	50	1036	99	1128
Agnes Scott College	Yes	417	349	137	60	89	510	63	1296
Alaska Pacific University	Yes	193	146	55	16	44	249	869	756

Check the `info()` and `describe()` methods on the data.


```
In [6]: sns.set_style('whitegrid')
sns.lmplot('Room.Board', 'Grad.Rate', data=df, hue='Private',
          palette='coolwarm', size=6, aspect=1, fit_reg=False)
```

```
Out[6]: <seaborn.axisgrid.FacetGrid at 0xa81b550>
```



Create a scatterplot of F.Undergrad versus Outstate where the points are colored by the Private column.

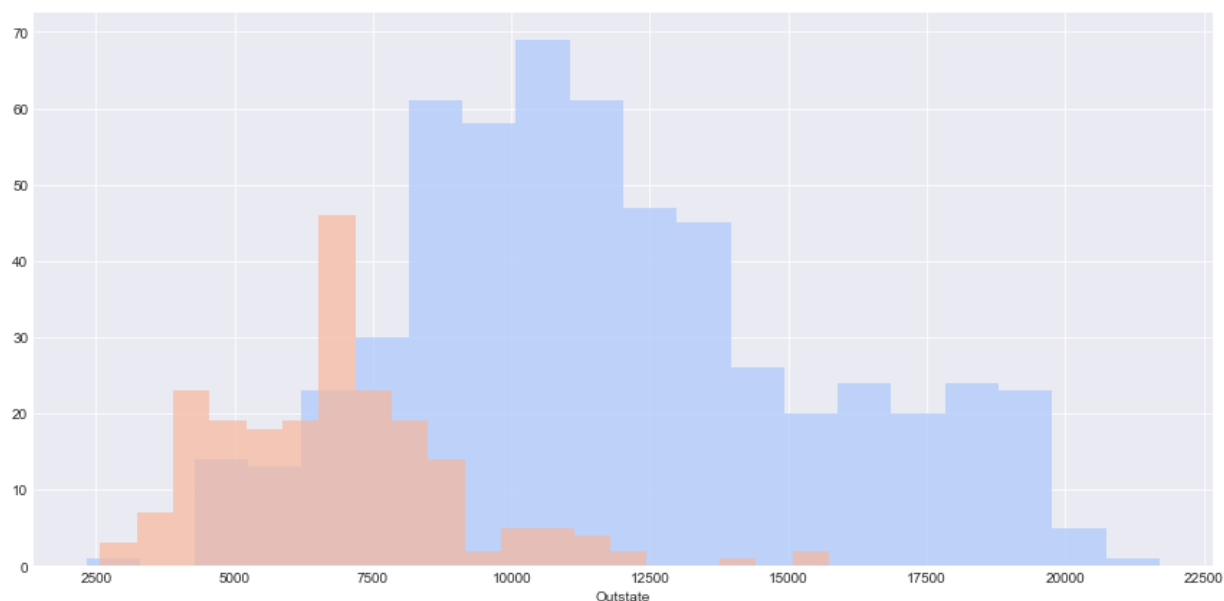
```
In [7]: sns.set_style('whitegrid')
sns.lmplot('Outstate', 'F.Undergrad', data=df, hue='Private',
          palette='coolwarm', size=6, aspect=1, fit_reg=False)
```

```
Out[7]: <seaborn.axisgrid.FacetGrid at 0xa9fe4e0>
```



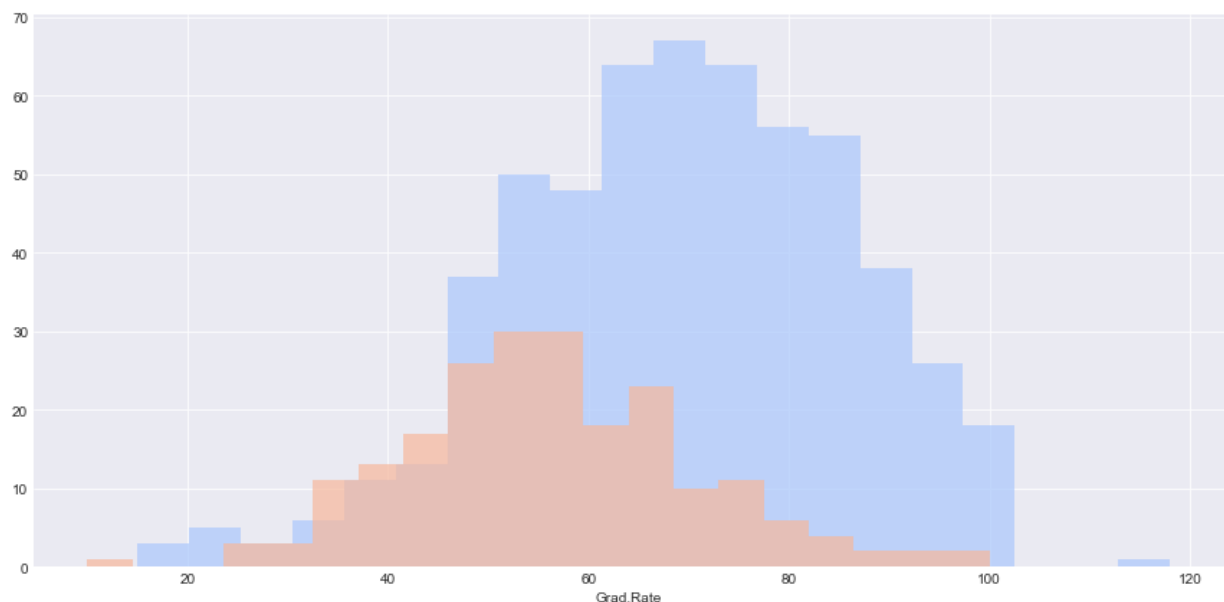
Create a stacked histogram showing Out of State Tuition based on the Private column. Try doing this using [sns.FacetGrid](https://stanford.edu/~mwaskom/software/seaborn/generated/seaborn.FacetGrid.html) (<https://stanford.edu/~mwaskom/software/seaborn/generated/seaborn.FacetGrid.html>). If that is too tricky, see if you can do it just by using two instances of `pandas.plot(kind='hist')`.

```
In [8]: sns.set_style('darkgrid')
g = sns.FacetGrid(df,hue="Private",palette='coolwarm',size=6,aspect=2)
g = g.map(plt.hist,'Outstate',bins=20,alpha=0.7)
```



Create a similar histogram for the Grad.Rate column.

```
In [9]: sns.set_style('darkgrid')
g = sns.FacetGrid(df,hue="Private",palette='coolwarm',size=6,aspect=2)
g = g.map(plt.hist,'Grad.Rate',bins=20,alpha=0.7)
```



Notice how there seems to be a private school with a graduation rate of higher than 100%. What is the name of that school?

Now it is time to create the Cluster labels!

Import KMeans from SciKit Learn.

```
In [14]: from sklearn.cluster import KMeans
```

Create an instance of a K Means model with 2 clusters.

```
In [15]: kmeans = KMeans(n_clusters=2)
```

Fit the model to all the data except for the Private label.

```
In [16]: kmeans.fit(df.drop('Private',axis=1))
```

```
Out[16]: KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
               n_clusters=2, n_init=10, n_jobs=1, precompute_distances='auto',
               random_state=None, tol=0.0001, verbose=0)
```

What are the cluster center vectors?

```
In [17]: kmeans.cluster_centers_
```

```
Out[17]: array([[1.03631389e+04, 6.55089815e+03, 2.56972222e+03, 4.14907407e+01,
                7.02037037e+01, 1.30619352e+04, 2.46486111e+03, 1.07191759e+04,
                4.64347222e+03, 5.95212963e+02, 1.71420370e+03, 8.63981481e+01,
                9.13333333e+01, 1.40277778e+01, 2.00740741e+01, 1.41705000e+04,
                6.75925926e+01],
               [1.81323468e+03, 1.28716592e+03, 4.91044843e+02, 2.53094170e+01,
                5.34708520e+01, 2.18854858e+03, 5.95458894e+02, 1.03957085e+04,
                4.31136472e+03, 5.41982063e+02, 1.28033632e+03, 7.04424514e+01,
                7.78251121e+01, 1.40997010e+01, 2.31748879e+01, 8.93204634e+03,
                6.50926756e+01]])
```

Evaluation

There is no perfect way to evaluate clustering if you don't have the labels, however since this is just an exercise, we do have the labels, so we take advantage of this to evaluate our clusters, keep in mind, you usually won't have this luxury in the real world.

Create a new column for df called 'Cluster', which is a 1 for a Private school, and a 0 for a public school.

```
In [18]: def converter(cluster):
          if cluster=='Yes':
              return 1
          else:
              return 0
```

```
In [19]: df['Cluster'] = df['Private'].apply(converter)
```