

DM_Assignment_2_b

August 16, 2021

0.1 Choose an appropriate application and perform partitional clustering using K Means Algorithm

About the Data

Context

Statistics for a large number of US Colleges from the 1995 issue of US News and World Report.

Content

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

Source

This dataset was taken from the StatLib library which is maintained at Carnegie Mellon University.

The dataset was used in the ASA Statistical Graphics Section's 1995 Data Analysis Exposition.

```
[1]: # basic imports

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
[3]: # creating a dataframe

df=pd.read_csv("/content/drive/MyDrive/AI-ML/DM/College.csv",index_col=0)
```

```
[4]: df.head()
```

```
[4]:
```

	Private	Apps	...	Expend	Grad.Rate
Abilene Christian University	Yes	1660	...	7041	60
Adelphi University	Yes	2186	...	10527	56
Adrian College	Yes	1428	...	8735	54
Agnes Scott College	Yes	417	...	19016	59
Alaska Pacific University	Yes	193	...	10922	15

[5 rows x 18 columns]

```
[5]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 777 entries, Abilene Christian University to York College of Pennsylvania
Data columns (total 18 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Private                777 non-null   object
1   Apps                   777 non-null   int64
2   Accept                 777 non-null   int64
3   Enroll                 777 non-null   int64
4   Top10perc              777 non-null   int64
5   Top25perc              777 non-null   int64
6   F.Undergrad            777 non-null   int64
7   P.Undergrad            777 non-null   int64
8   Outstate               777 non-null   int64
9   Room.Board             777 non-null   int64
10  Books                  777 non-null   int64
11  Personal               777 non-null   int64
12  PhD                    777 non-null   int64
13  Terminal               777 non-null   int64
14  S.F.Ratio              777 non-null   float64
15  perc.alumni            777 non-null   int64
16  Expend                 777 non-null   int64
```

```
17 Grad.Rate    777 non-null    int64
dtypes: float64(1), int64(16), object(1)
memory usage: 135.3+ KB
```

```
[6]: df.describe()
```

```
[6]:
```

	Apps	Accept	...	Expend	Grad.Rate
count	777.000000	777.000000	...	777.000000	777.00000
mean	3001.638353	2018.804376	...	9660.171171	65.46332
std	3870.201484	2451.113971	...	5221.768440	17.17771
min	81.000000	72.000000	...	3186.000000	10.00000
25%	776.000000	604.000000	...	6751.000000	53.00000
50%	1558.000000	1110.000000	...	8377.000000	65.00000
75%	3624.000000	2424.000000	...	10830.000000	78.00000
max	48094.000000	26330.000000	...	56233.000000	118.00000

```
[8 rows x 17 columns]
```

```
[9]: # visualizations

# graduation rate vs room & board costs for private and public

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

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning:
Pass the following variables as keyword args: x, y. From version 0.12, the only
valid positional argument will be `data`, and passing other arguments without an
explicit keyword will result in an error or misinterpretation.
FutureWarning
```

```
[9]: <seaborn.axisgrid.FacetGrid at 0x7f9146129d50>
```



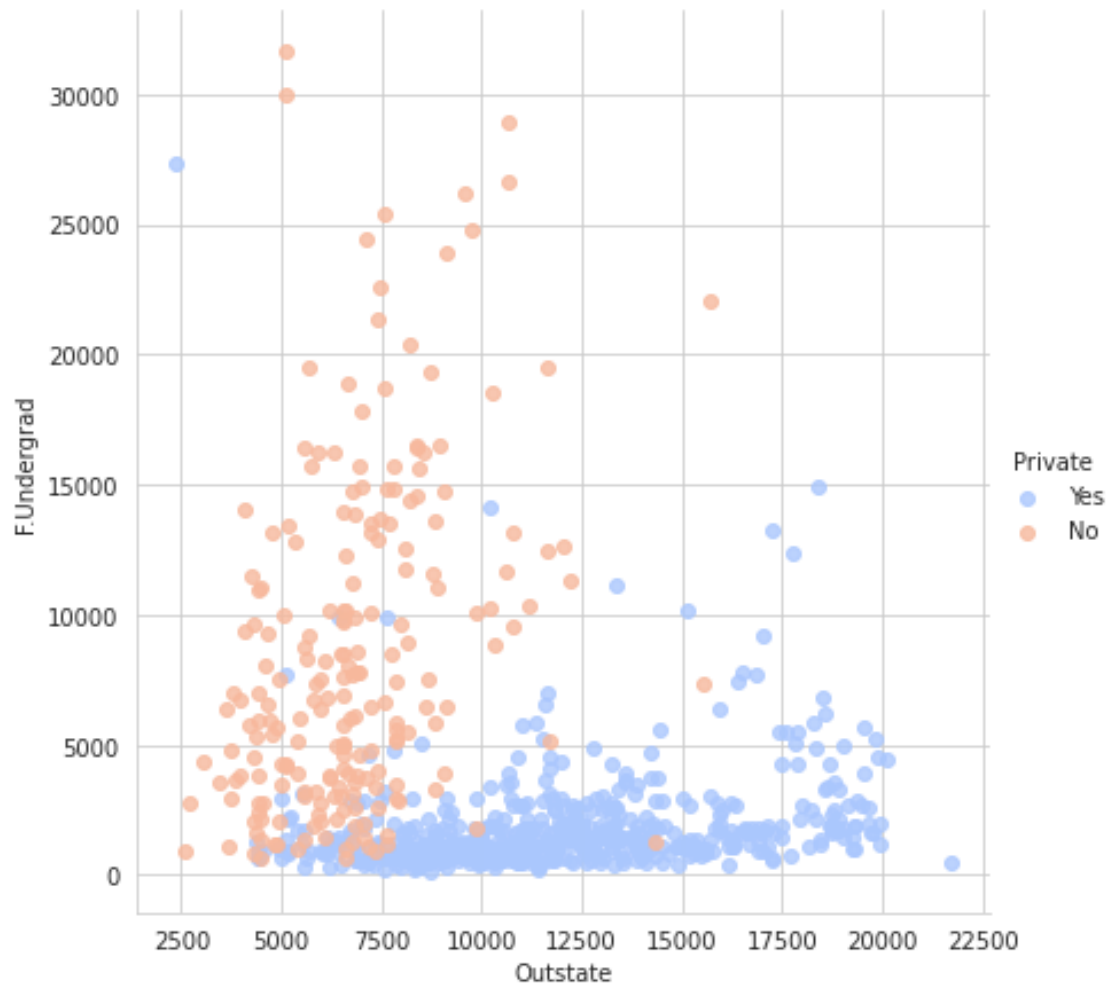
```
[10]: # Number of fulltime undergraduates vs out of state tuition

sns.set_style('whitegrid')
sns.lmplot('Outstate', 'F.Undergrad', data=df, hue='Private',
          palette='coolwarm', height=6, aspect=1, fit_reg=False)
```

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

FutureWarning

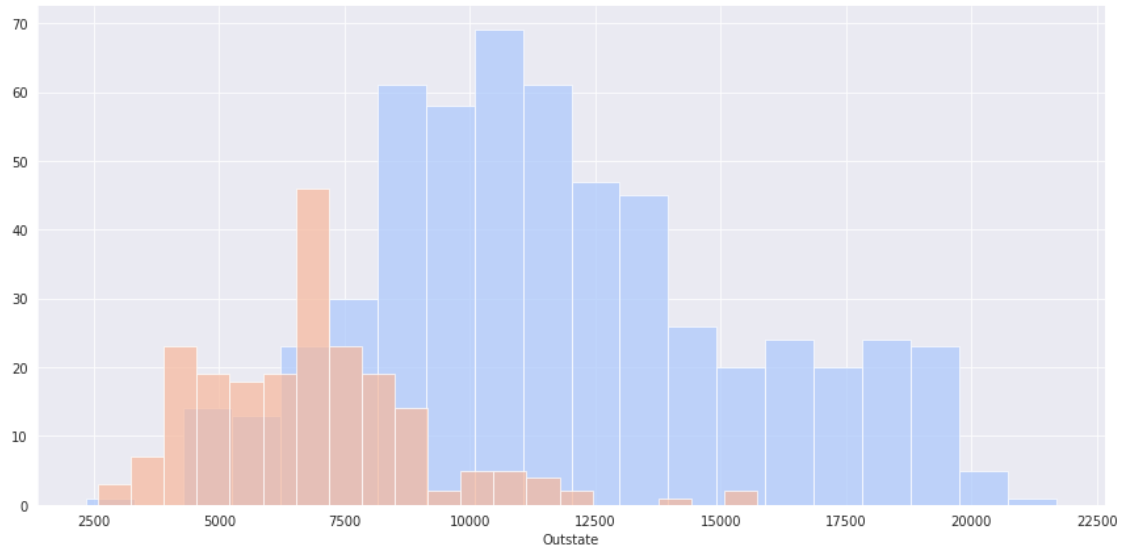
```
[10]: <seaborn.axisgrid.FacetGrid at 0x7f9145ff94d0>
```



```
[11]: # stacked histogram of outstate with colors for private or not
```

```
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)
```

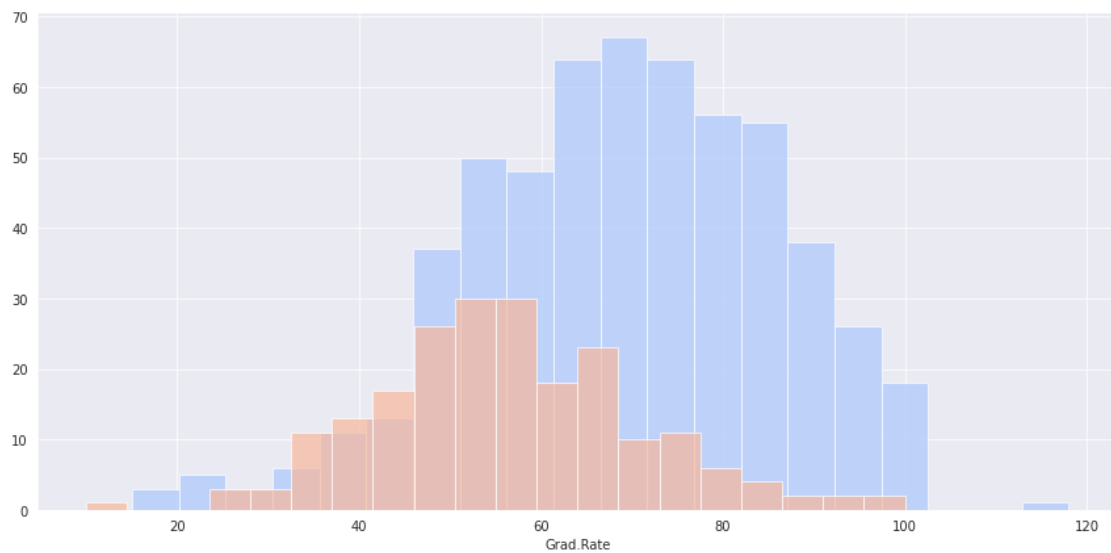
```
/usr/local/lib/python3.7/dist-packages/seaborn/axisgrid.py:316: UserWarning: The
`size` parameter has been renamed to `height`; please update your code.
  warnings.warn(msg, UserWarning)
```



[12]: *# same above for grad rate*

```
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)
```

/usr/local/lib/python3.7/dist-packages/seaborn/axisgrid.py:316: UserWarning: The `size` parameter has been renamed to `height`; please update your code.
warnings.warn(msg, UserWarning)



```
[13]: # grad rate of cazenovia college is greater than 100 which doesnt make any sense

# setting the grad rate of this college to 100.
df['Grad.Rate']['Cazenovia College'] = 100
```

/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:1:

SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

"""Entry point for launching an IPython kernel.

0.1.1 Kmeans clustering part

```
[18]: from sklearn.cluster import KMeans

# only 2 clusters selected
# private or not

kmeans=KMeans(n_clusters=2)

# here dropping the private column
# we need to predict or cluster the data into two clusters
# private or not

kmeans.fit(df.drop('Private',axis=1))
```

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

```
[19]: # cluster vectors
```

```
kmeans.cluster_centers_
```

```
[19]: 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]])
```

```
kmeans.labels_
```

[illegible]

0.1.2 Evaluation

since we are clustering for private or not and we already know the real labels. we can do evaluations. But this may not be possible in real life scenarios

```
# creating a new column on the existing dataframe to encode
# the string
# 1 - private
```



```

# 0 - not private

def converter(cluster):
    if cluster=='Yes':
        return 1
    else:
        return 0

df['Cluster'] = df['Private'].apply(converter)

# creates new column "Cluster" where if private --> 1
# not private --> 0

df.head()

```

```

[20]:

```

	Private	Apps	Accept	...	Expend	Grad.Rate
Cluster						
Abilene Christian University	Yes	1660	1232	...	7041	60
1						
Adelphi University	Yes	2186	1924	...	10527	56
1						
Adrian College	Yes	1428	1097	...	8735	54
1						
Agnes Scott College	Yes	417	349	...	19016	59
1						
Alaska Pacific University	Yes	193	146	...	10922	15
1						

[5 rows x 19 columns]

```

[21]: # Evaluating using confusion matrix

from sklearn.metrics import confusion_matrix, classification_report

print(confusion_matrix(df['Cluster'], kmeans.labels_))
print(classification_report(df['Cluster'], kmeans.labels_))

```

```

[[ 74 138]
 [ 34 531]]

```

	precision	recall	f1-score	support
0	0.69	0.35	0.46	212
1	0.79	0.94	0.86	565
accuracy			0.78	777
macro avg	0.74	0.64	0.66	777
weighted avg	0.76	0.78	0.75	777

