K Means Clustering Project For this project we will attempt to use KMeans Clustering to cluster Universities into to 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. df = pd.read\_csv('19 K Means Clustering Project.csv',index\_col=0) Check the head of the data In [3]: df.head() Private Apps Accept Enroll Top10perc Top25perc F.Undergrad P.Undergrad Outstate Room.Board Books Personal PhD Terminal S.F.Ratio perc.alumni Expend Grad.Rate Out[3]: Yes 1660 1232 721 23 52 2885 537 7440 3300 450 2200 70 78 18.1 7041 60 Christian 12 University Adelphi Yes 2186 1924 512 16 29 2683 1227 12280 6450 750 1500 29 30 12.2 10527 16 56 University Adrian 1428 1097 336 22 50 1036 11250 53 12.9 8735 Yes 99 3750 400 1165 66 30 54 College **Agnes Scott** 417 137 60 89 12960 92 7.7 19016 Yes 349 510 63 5450 450 875 97 59 College **Alaska Pacific** 193 146 55 16 44 249 869 7560 4120 800 1500 76 72 11.9 2 10922 15 Yes University Check the info() and describe() methods on the data. In [4]: 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 1 777 non-null Apps int64 2 Accept 777 non-null int64 Enroll 777 non-null int64 777 non-null Top10perc int64 777 non-null Top25perc int64 F.Undergrad 777 non-null int64 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 PhD 777 non-null 12 int64 Terminal 777 non-null int64 13 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: 115.3+ KB In [5]: df.describe() Accept Enroll Top10perc Top25perc F.Undergrad P.Undergrad Outstate Room.Board **Books** Personal PhD Terminal S.F.Ratio Out[5]: **Apps** 777.000000 777.000000 777.000000 777.000000 777.000000 777.000000 777.000000 777.000000 777.000000 777.000000 777.000000 777.000000 777.0 count 777.000000 777.000000 1340.642214 3001.638353 855.298584 10440.669241 549.380952 2018.804376 779.972973 27.558559 55.796654 3699.907336 4357.526384 72.660232 79.702703 14.089704 22.7 mean 16.328155 14.722359 3.958349 std 3870.201484 2451.113971 929.176190 17.640364 19.804778 4850.420531 1522.431887 4023.016484 1096.696416 165.105360 677.071454 12.3 81.000000 72.000000 35.000000 1.000000 9.000000 139.000000 1.000000 2340.000000 1780.000000 96.000000 250.000000 8.000000 24.000000 2.500000 0.0 min 242.000000 992.000000 3597.000000 850.000000 **25**% 776.000000 604.000000 15.000000 41.000000 95.000000 7320.000000 470.000000 62.000000 71.000000 11.500000 13.0 1200.000000 1558.000000 54.000000 353.000000 9990.000000 75.000000 82.000000 13.600000 **50**% 1110.000000 434.000000 23.000000 1707.000000 4200.000000 500.000000 21.0 35.000000 967.000000 3624.000000 2424.000000 902.000000 69.000000 4005.000000 12925.000000 5050.000000 600.000000 1700.000000 85.000000 92.000000 16.500000 31.0 75% max 48094.000000 26330.000000 6392.000000 96.000000 100.000000 31643.000000 21836.000000 21700.000000 8124.000000 2340.000000 6800.000000 103.000000 100.000000**EDA** It's time to create some data visualizations! Create a scatterplot of Grad.Rate versus Room.Board where the points are colored by the Private column. 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) C:\Users\Yasin\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, th e only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation. warnings.warn( C:\Users\Yasin\anaconda3\lib\site-packages\seaborn\regression.py:580: UserWarning: The `size` parameter has been renamed to `height`; please update your code. warnings.warn(msg, UserWarning) Out[6]: <seaborn.axisgrid.FacetGrid at 0x2054fa38190> 100 Private 60 40 2000 3000 4000 5000 7000 8000 6000 Room.Board 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) C:\Users\Yasin\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, th e only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation. warnings.warn( C:\Users\Yasin\anaconda3\lib\site-packages\seaborn\regression.py:580: UserWarning: The `size` parameter has been renamed to `height`; please update your code. warnings.warn(msg, UserWarning) Out[7]: <seaborn.axisgrid.FacetGrid at 0x20574831be0> 30000 25000 20000 Private 15000 10000 Create a stacked histogram showing Out of State Tuition based on the Private column. Try doing this using sns.FacetGrid. 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) C:\Users\Yasin\anaconda3\lib\site-packages\seaborn\axisgrid.py:316: UserWarning: The `size` parameter has been renamed to `height`; please update your code. warnings.warn(msg, UserWarning) 60 50 20 0 2500 5000 7500 10000 12500 15000 17500 20000 22500 Outstate 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) C:\Users\Yasin\anaconda3\lib\site-packages\seaborn\axisgrid.py:316: UserWarning: The `size` parameter has been renamed to `height`; please update your code. warnings.warn(msg, UserWarning) 60 50 40 20 Grad.Rate Notice how there seems to be a private school with a graduation rate of higher than 100%. What is the name of that school? In [10]: df[df['Grad.Rate'] > 100] Private Apps Accept Enroll Top10perc Top25perc F.Undergrad P.Undergrad Outstate Room.Board Books Personal PhD Terminal S.F.Ratio perc.alumni Expend Grad.Rate Out[10]: Cazenovia Yes 3847 3433 35 1010 12 9384 4840 600 500 14.3 20 7697 118 College Set that school's graduation rate to 100 so it makes sense. You may get a warning not an error) when doing this operation, so use dataframe operations or just re-do the histogram visualization to make sure it actually went through. In [11]: df['Grad.Rate']['Cazenovia College'] = 100 <ipython-input-11-bc95ac68ab2d>: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 df['Grad.Rate']['Cazenovia College'] = 100 In [12]: df[df['Grad.Rate'] > 100] Private Apps Accept Enroll Top10perc Top25perc F.Undergrad P.Undergrad Outstate Room.Board Books Personal PhD Terminal S.F.Ratio perc.alumni Expend Grad.Rate Out[12]: In [13]: 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) C:\Users\Yasin\anaconda3\lib\site-packages\seaborn\axisgrid.py:316: UserWarning: The `size` parameter has been renamed to `height`; please update your code. warnings.warn(msg, UserWarning) 60 20 Grad.Rate K Means Cluster Creation 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(n\_clusters=2) What are the cluster center vectors? kmeans.cluster\_centers\_ Out[17]: array([[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], [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]]) **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) df.head() Out[20]: Private Apps Accept Enroll Top10perc Top25perc F.Undergrad P.Undergrad Outstate Room.Board Books Personal PhD Terminal S.F.Ratio perc.alumni Expend Grad.Rate Cluste Abilene Christian 23 52 2885 7440 450 2200 70 78 18.1 7041 60 Yes 1660 1232 721 537 3300 12 University Adelphi Yes 2186 1924 512 16 29 2683 1227 12280 6450 750 1500 29 30 12.2 16 10527 56 University Adrian 1428 1097 336 22 50 1036 99 11250 3750 400 1165 53 12.9 30 8735 54 College **Agnes** 510 12960 450 97 19016 Scott Yes 417 349 137 60 89 63 5450 875 92 7.7 37 59 College Alaska **Pacific** Yes 193 146 55 16 249 869 7560 4120 800 1500 76 10922 15 University Create a confusion matrix and classification report to see how well the Kmeans clustering worked without being given any labels. In [21]: kmeans.labels\_ 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, Θ, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, Θ, Θ, Θ, Θ, Θ, 0, 1, 0, 0, 0, Ο, Θ, Θ, Ο, Õ, 0, 0, 0, 1, Θ, Θ, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, Θ, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, Θ, 0, 1, 1, 0, 0, 0, 0, Θ, 0, 0, 0, Θ, Θ, Ο, Ο, Θ, Ο, Θ, 0, 0, 0, 0, 0, 0, 0, Θ, Θ, Θ, Ο, Ο, Ο, Θ, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, Θ, Θ, 0, 0, 0, Θ, Θ, 0, 0, 0, 0, 0, Θ, 0, 0, 0, Θ, Θ, Θ, Θ, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, Θ, Θ, Θ, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, Θ, Ο, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 1, 1, Θ, 1, 1,  $0,\ 1,\ 0,\ 1,\ 0,\ 0,\ 0,\ 0,\ 1,\ 1,\ 1,\ 1,\ 0,\ 0,\ 0,\ 0,\ 1,$ 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0]) from sklearn.metrics import confusion\_matrix,classification\_report print(confusion\_matrix(df['Cluster'], kmeans.labels\_)) print(classification\_report(df['Cluster'], kmeans.labels\_)) [[138 74] [531 34]] precision recall f1-score support 0 0.21 0.65 0.31 212 0.31 0.06 0.10 accuracy 0.22 777 macro avg 0.26 0.36 0.21 777 777 weighted avg 0.29 0.22 0.16