Group 3 NYC AirBnB Data Principle Component Analysis  Step-1: Import Packages, Data Exploration & Cleaning  Intent should be to maximize information content while minimizing the number of dimensions or variables in the datset. Most of the variables in the dataframe have non null values except 'last_review', 'reviews_per_month', and 'license'. These columns are believed to have either little information value, or are believed to have other variables that would explain their effects anyway. Thus the aforementioned columns were removed. In addition all observations with zero values in the 'price' column were removed. It is assumed that the 'price' will be the primary variable of interest.
In [1]: import pandas as pd import numpy as np from sklearn.decomposition import PCA import matplotlib.pyplot as plt from sklearn.preprocessing import scale import seaborn as sb import plotly.express as px  In [2]: df = pd.read_csv('Grp3Project_InitialData/listings.csv')
<pre>class 'pandas.core.frame.DataFrame' &gt; RangeIndex: 39881 entries, 0 to 3980  Data columns (total 18 columns):  # Column</pre>
6 latitude 39881 non-null float64 7 longitude 39881 non-null float64 8 room_type 39881 non-null object 9 price 39881 non-null int64 10 minimum_nights 39881 non-null int64 11 number_of_reviews 39881 non-null int64 12 last_review 31519 non-null object 13 reviews per_month 31519 non-null float64 14 calculated_host_listings_count 39881 non-null int64 15 availability_365 39881 non-null int64 16 number_of_reviews_ltm 39881 non-null int64 17 license 5 non-null object
<pre>dtypes: float64(3), int64(8), object(7) memory usage: 5.5+ MB None  In [3]: # remove deemed to be less than useful</pre>
Int64Index: 39851 entries, 0 to 39880  Data columns (total 12 columns):  # Column
8 calculated_host_listings_count 39851 non-null int64 9 availability_365 39851 non-null int64 10 number_of_reviews_ltm 39851 non-null int64 11 lat_lon 39851 non-null float64 dtypes: float64(3), int64(7), object(2) memory usage: 4.0+ MB None  15000
40.8 - 40.7 - 40.5 - 40
-73.8 -74.0 -74.2
-0.546 $-0.548$ $-0.550$ $-0.552$
The pairplot was developed to gain an understanding of the relationship between 'latitude', 'longitude', and 'price'. Based on this visualization, that correlation appears to be weak.
Step 2: Determine & Remove the Outliers  Process was to remove all outliers larger than q3 + (1.5 * IQR), or less than or equal to the minimum. In addition the latitude and longitude columns were removed from the dataset and replaced with the transformation 'lat_lon', which is the ratio latitude to longitude. The result reduced the number of observations in the data from 39881 to 36524.  In [4]:    df_stats = df.groupby('room_type')['price'].describe()   df_stats.reindex()   display(df_stats)    count mean std min 25% 50% 75% max   room_type   com_type
Entire home/apt 22761.0 251.546022 338.044654 10.0 128.00 180.0 275.0 15000.0  Hotel room 172.0 436.470930 282.491141 10.0 258.75 308.0 501.5 1998.0  Private room 16361.0 122.936495 356.373737 10.0 55.00 75.0 110.0 16500.0  Shared room 567.0 119.398564 454.106078 10.0 42.00 66.0 92.0 10000.0  In [5]: # Create Table with outliers df_stats["15%"] - df_stats["15%"] - df_stats["10,"] = df_stats["15%"] - df_stats["10,"] = df_stats["10,"] = df_stats["10,"] + (1.5 * df_stats['10,"]) df_stats['10,"] - (1.5 * df_stats['10,"]) df_stats['10,"] - (1.5 * df_stats['10,"])
display(df_stats)    Count   Mean   std   Min   25%   50%   75%   Max   IQR   upper_outlier   lower_outlier
<pre>In [6]: df1 = df[(df['room_type'] == "Entire home/apt") &amp; (df['price']&gt;=10) &amp; (df['price']&lt;= 496)]     df2 = df[(df['room_type'] == "Private room") &amp; (df['price']&gt;= 10) &amp; (df['price']&lt;= 192)]     df3 = df[(df['room_type'] == "Shared room") &amp; (df['price']&gt;= 10) &amp; (df['price']&lt;= 167)]     df4 = df[(df['room_type'] == "Hotel room") &amp; (df['price']&gt;= 100) &amp; (df['price']&lt;= 865)]     frames = [df1, df2, df3, df4]     df = pd.concat(frames)  display(df.groupby('neighbourhood_group')['price'].describe())     display(df.groupby('room_type')['price'].describe())     sb.boxplot(y = df['price'], x = df['room_type'])     plt.show()</pre>
sb.boxplot(y = df['price'], x = df['neighbourhood_group']) plt.show() display(df.info()) display(df.corr())   recount mean std min 25% 50% 75% max  neighbourhood_group  Bronx 1530. 108.036601 68.56101 10.0 60.0 89.0 135.00 450.0  Brooklyn 14192.0 133.914459 88.425233 10.0 69.0 110.0 175.00 529.0
Manhattan         14/426, 180.040690         104.040690         105.05595         10.0         100.0         25.00         256.0           Queens         595.0         114.189620         79.028227         10.0         60.0         90.0         145.00         495.0           Staten Island         422.0         113.646919         618.27482         33.0         69.0         99.5         140.75         412.0           room_type           Entire home/apt         2117.0         196.250793         95.38230         10.0         125.0         496.0         496.0           Hotel room         155.0         363.490323         163.47128         100.0         248.5         307.0         470.0         826.0
Private room 14751.0 78.448648 34.294310 10.0 52.0 72.0 97.0 192.0  Shared room 501.0 66.179641 32.916921 10.0 40.0 60.0 85.0 165.0  800
500 - <del> </del> <del> </del> <del> </del> <del> </del> 400 - <del> </del> 200 - <del> </del> 100
Entire home/apt Private room Shared room room_type
600 - 500 - 9 400 - 300 - 200 -
Queens Manhattan Brooklyn Bronx Staten Island neighbourhood_group <class 'pandas.core.frame.dataframe'=""> Int64Index: 36524 entries, 2 to 19816 Data columns (total 12 columns):</class>
# Column Non-Null Court bype
11 lat_lon
minimum_nights         -0.114942         0.03942         -0.09853         -0.08823         1.00000         -0.147075         0.124210         -0.054423         -0.234587         0.017041           number_of_reviews         -0.187285         -0.032738         0.040181         -0.00074         -0.147075         1.00000         -0.096784         0.089859         0.629522         0.009497           calculated_host_listings_count         0.041516         0.031671         -0.063658         0.152782         0.124210         -0.096784         1.00000         0.131778         -0.059411         0.00320           number_of_reviews_ltm         -0.084761         -0.02163         0.143466         0.155715         -0.054423         0.089859         0.131778         1.000000         0.160993         -0.048121           number_of_reviews_ltm         -0.084761         -0.042166         0.077105         0.08203         -0.29522         -0.059411         0.160993         1.000000         0.000450           lat_lon         -0.08676         -0.886394         -0.512761         0.082805         0.017041         0.002320         -0.048121         0.000450         0.000450
In [7]: af select = df[['price','lat_lon','neighbourhood_group','room_type', 'minimum_nights', 'number_of_reviews','calculated_host_listings_count','availability_365','number_of_reviews_ltm']]  df_select.info() <class 'pandas.core.frame.dataframe'=""></class>
5 number_of_reviews 36524 non-null int64 6 calculated_host_listings_count 36524 non-null int64 7 availability_365 36524 non-null int64 8 number_of_reviews_ltm 36524 non-null int64 dtypes: float64(1), int64(6), object(2) memory usage: 2.8+ MB  Step 3: Incorporate Categorical Variables into the Data and Scale the Dataset  The purpose of this step was to incorporate room type and neighborhood or NYC burrough into the dataset as "dummy variables" using one hot encoding. Also conducted during this step was a targeted correlation between price and the remaining variables. The strongest positive correlations were room types of 'Entire home/apt', neighborhood groups of 'Manhattan' and
365 day availability. Strongest negative correlation was a room type of "Private room". Finally the dataset was normalized.  In [8]:  # One-hot encode the neighbourhood_group column
display(df.info()) display(df.focolumns[0:]].corr()['price'][:].sort_values(ascending=False).to_frame()) <class 'pandas.core.frame.dataframe'=""> Int64Index: 36524 entries, 2 to 19816  Data columns (total 16 columns):  # Column Non-Null Count Dtype</class>
5 availability_365 36524 non-null int64 6 number_of_reviews_ltm 36524 non-null int64 7 Bronx 36524 non-null uint8 8 Brooklyn 36524 non-null uint8 9 Manhattan 36524 non-null uint8 10 Queens 36524 non-null uint8 11 Staten Island 36524 non-null uint8 12 Entire home/apt 36524 non-null uint8 13 Hotel room 36524 non-null uint8 14 Private room 36524 non-null uint8 15 Shared room 36524 non-null uint8 16 Shared room 36524 non-null uint8 17 Shared room 36524 non-null uint8 18 Shared room 36524 non-null uint8 19 Figure room 36524 non-null uint8 19 Figure room 36524 non-null uint8
memory usage: 2.5 MB None  price  price  1.00000  Entire home/apt 0.585895  Manhattan 0.269626  availability_365 0.155715  calculated_host_listings_count 0.152782
Hotel room       0.144978         lat_lon       0.082805         number_of_reviews_ltm       0.082033         number_of_reviews       -0.000474         Staten Island       -0.037760         Bronx       -0.085095         minimum_nights       -0.088323         Shared room       -0.098770
Brooklyn -0.112220 Queens -0.151669 Private room -0.585485  In [10]: X = scale(df) # Normalize the data display(pd.DataFrame(X).describe().transpose()) # Display the normalize data statistics corr_ = pd.DataFrame(X).corr() #Display the correlation matrix for the normalized data pd.DataFrame(corr_).style.background_gradient(cmap = 'YlOrRd')
count         man         std         min         25%         50%         75%         max           0         36524.0         5.539893e-15         1.00014         -1.415415         -0.746793         -0.283901         0.43913         6.98363           1         36524.0         9.434973e-14         1.00014         -3.78356         -0.643971         0.17789         0.56800         5.630077           2         36524.0         -1.498075e-14         1.00014         -0.587472         -0.556147         0.258299         0.32948         38.537241           3         36524.0         -1.983690e-14         1.00014         -0.68374         -0.263174         -0.228420         7.59138           5         36524.0         7.306705e-15         1.00014         -0.93275         -0.429781         1.02812         1.727913           6         36524.0         7.110498e-15         1.00014         -0.47298         -0.47298         -0.44790         5.3849669
7         36524.0         4.900318e-14         1.000014         -0.209098         -0.209098         -0.209098         -0.209098         4.782457           8         36524.0         2.021351e-15         1.000014         -0.807972         -0.807972         -0.807972         -0.807972         -0.807972         -0.807972         -0.41323         -0.441323         -0.441323         -0.441323         -0.441323         -0.441323         -0.41323         -0.441323         -0.451323         -0.25913           10         36524.0         4.617554e-15         1.000014         -0.108116         -0.108116         -0.108116         -0.108116         -0.108116         -0.108116         -0.108116         -0.854167           12         36524.0         -1.630806e-13         1.000014         -1.17731         -1.17731         0.854167         0.854167           13         36524.0         2.449373e-15         1.000014         -0.065283
14 36240 -7.222863e15 1,00014 -0.82308
4 0.152782 0.002320 0.124210 -0.096784 1.000000 0.13178 -0.059441 -0.046750 -0.101159 0.114416 0.014788 -0.025402 0.014032 0.001746 -0.008779 -0.023518 5 0.155715 -0.048121 -0.054423 0.09859 0.131778 1.000000 0.160993 0.102307 -0.075341 -0.06966 0.105373 0.066505 0.055361 0.044967 -0.064226 0.010751 6 0.082033 0.00460 -0.234587 0.62952 -0.059441 0.160993 1.00000 0.026521 0.00209 0.026521 0.00209 0.026521 0.00209 0.026521 0.00209 0.026521 0.00209 0.026521 0.00209 0.026521 0.00209 0.026521 0.00209 0.026521 0.00000 0.026521 0.00209 0.00209 0.026521 0.00209 0.0026221 0.00209 0
11 -0.037760
The PCA was conducted on the normalized data and indicate the following results:  • The most important variables in the dataset are 'price', 'Manhattan', 'Entire home/apt', and 'Private room'  • The least significant variables seem to be 'Hotel room' and 'Queens'.  • Information content in the data is dispersed. There dont seem to be any significant clusters of variability. As indicated by the scree plot and the measures of the cumulative variability, you dont get above 90% of the variability explained until you get to 11 of the 16 principal components. Once you get to roughly 13 principal components, marginal variability apears to go to zero. As evidenced by the correlation matrices and the heatmaps, most of the variables could be considered to have a significant component loading depending on the dimension in which they are viewed. The practical effect of this seems to be that action taken to reduce the dimensionality of the dataset won't add any significant benefits.  In [11]:  # Select the number of components peal = PCA(n_components=16)
# fit the PCA model pcal.fit(x)  Out[11]: PCA(n_components=16)  In [12]: #The amount of variance that each PC explains var = pcal.explained_variance_ratio_ display(var)  array([1.65195096e-01, 1.34942542e-01, 1.17200904e-01, 8.08398602e-02, 7.70139078e-02, 6.83907993e-02, 6.70856015e-02, 6.42368108e-02, 6.33128426e-02, 5.42576742e-02, 4.61935056e-02, 2.1700494e-02, 2.17283626e-02, 1.3432045e-02, 6.2380475e-33, 5.16641682e-33])
<pre>In [13]: #Cumulative Variance explains</pre>
plt.ylabel('Cumulative Variance Explained') plt.savefig('scree_plot.png',dpi=100,bbox_inches='tight') plt.show()  Scree Plot  100 -
Cumulative Variance Expl
In [15]:  # Select the number of components pcal6 = PCA(n_components=16) pcal6.fit(X) data_pcal6 = pcal6.transform(X)
# Convert the numpy array to pandas DataFrame data_pcal6 = pd.DataFrame(data_pcal6) data_pcal6.columns = ["PC"+str(i) for i in range(1,17)]  # Show the head of the DataFrame display(data_pcal6.head()) display(data_pcal6.corr())  PC1 PC2 PC3 PC4 PC5 PC6 PC7 PC8 PC9 PC10 PC11 PC12 PC13 PC14 PC15 PC16  0 1.985462 -0.577469 0.687516 -1.269858 -2.613957 -0.497312 -0.356578 0.348281 -0.725955 0.259770 -1.033230 1.721738 -1.038442 0.247408 3.330669e-16 7.771561e-16
1 1.649282
PC3         -5.197114e-17         1.836350e-16         1.000000e-10         7.357477e-17         1.097143e-16         -2.04536e-16         -4.49512e-16         3.2321e-16         -5.8999e-17         -5.8999e-17         -5.8999e-17         -4.9360e-17         -7.9360e-17         1.259448e-16         -2.93436e-16         -2.55150e-16         -7.57477e-17         1.00000e-10         -8.831691e-18         -4.67186r-16         -3.279430e-16         -3.279430e-16         -2.48872e-17         -1.063903e-15         -5.80496e-16         -2.75468e-16         -2.48870e-17         5.34649e-16         -0.45838e-16         -0.45368e-16         -0.45368e-16         -0.45368e-16         -2.77548e-16         -2.75468e-16         -2.75468e-16         -2.75468e-16         -2.75468e-16         -2.75468e-16         -2.75468e-16         -0.071309         -0.88569e-17         -0.44870e-16         -2.75468e-16         -2.01748e-16         -2.01748e-16 <th< th=""></th<>
PC11 -1.990741e-16 -4.979036e-16 7.49365e-17 -5.804964e-16 5.086496e-16 6.99940e-16 6.08923e-16 -7.572605e-16 -7.5
All pairwise combinations are cose to zero which is indicative of no multicollinearity among principal components.  In [16]: # Show the loadings for the 1st components peals.components.  Out[16]: array([[ 0.49436566,  0.04118856, -0.01297004,  0.02628435,  0.08448953,  0.0832199,  0.05827946,  -0.05827946,  -0.05822177,  -0.13351556,  0.27544056,  -0.15112606,  0.0024715,  0.05822177,  -0.55213588,  -0.04485018]])  In [17]: corr_16 = pd.DataFrame(pca16.componentstranspose(),  index=df.loc[i, 'price': 'Shared room'].columns,  columns,  columns,  range(].1711)
index=df.loc[; 'price': 'Shared room'] .columns; range(l,17]) display(pd.DataFrame(corr_16, shyle.background_gradient(cmap = 'YlorRd')) ax = sb.heatmap(corr_16, show())  PC1 PC2 PC3 PC6 PC6 PC6 PC7 PC8 PC9
calculated_host_listings_cout         0.08449         -0.09415         -0.14024         0.16965         0.033719         -0.28856         0.03019         -0.15564         0.03179         -0.03564         0.05179         -0.03016         -0.000000         -0.000000         -0.000000         -0.000000         -0.000000         -0.000000         -0.0000000         -0.0000000         -0.0000000
Staten Island 0.002472 0.63468 0.034519 0.20928 0.209281 0.060490 0.77529 -0.252848 -0.049195 0.14261 0.200540 0.042905 -0.024820 0.042905 -0.031886 0.104164 -0.075548 0.091312 0.024684 0.091312 0.024684 -0.274174 0.008817 -0.147154 0.016972 -0.098540 0.045075 0.103288 0.046521 -0.267217 0.091210 -0.001003 0.423029 0.5536174 0.056074
minimum_nights - number_of_reviews - calculated_host_listings_count - availability_365 - number_of_reviews_ltm - Bronx - Brooklyn - Manhattan - Manhattan -  - 0.75  - 0.75  - 0.50  - 0.50  - 0.25  - 0.25
Queens -