

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

from sklearn.preprocessing import KBinsDiscretizer

from sklearn import model_selection

from sklearn.linear_model import LogisticRegression

from sklearn.tree import DecisionTreeClassifier

from sklearn.neighbors import KNeighborsClassifier

from sklearn.svm import SVC

from sklearn.metrics import accuracy_score

from sklearn.metrics import confusion_matrix

from sklearn.metrics import classification_report

from sklearn.model_selection import cross_val_score, cross_val_predict, cross_validate

from sklearn import metrics

from sklearn.model_selection import train_test_split

from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier

warnings.filterwarnings("ignore")

In [2]: listings = pd.read_csv("D:\Trent\Notes\Data Mining\Assignments\Assignment 2\seattle\listin
gs.csv")
print(listings.head(3))

0      10  241832  https://www.airbnb.com/rooms/241832  2016010802432  2016-01-04
1      10  953595  https://www.airbnb.com/rooms/953595  2016010802432  2016-01-04
2      10  5308979  https://www.airbnb.com/rooms/5308979  2016010802432  2016-01-04

0      name \
1  Bright & Airy Queen Anne Apartment
2  New Modern House-Amazing water view

0      summary \
1  Chemically sensitive? We've removed the irrita...
2  New modern house built in 2013. Spectacular s...

0      space \
1  Beautiful, hypoallergenic apartment in an extr...
2  Our house is modern, light and fresh with a wa...

0      description experiences_offered \
1  Chemically sensitive? We've removed the irrita...
2  New modern house built in 2013. Spectacular s...

0      neighborhood_overview ... review_scores_value \
1  Queen Anne is a wonderful, truly functional vi...
2  Upper Queen Anne is a charming neighborhood f...

0      requires_license license jurisdiction_names instant_bookable \
1      f      NaN      WASHINGTON      f
2      NaN      NaN      WASHINGTON      f

0      cancellation_policy require_guest_profile_picture \
1      moderate      t
2      strict      f

0      require_guest_phone_verification calculated_host_listings_count \
1      t      6
2      f      2

0      reviews_per_month \
1      4.07
2      1.15
[3 rows x 92 columns]
```

```
In [3]: print(listings.info())

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3818 entries, 0 to 3817
Data columns (total 92 columns):
id      3818 non-null int64
listing_url      3818 non-null object
scrape_id      3818 non-null int64
last_scraped      3818 non-null object
name      3818 non-null object
summary      3641 non-null object
space      3249 non-null object
description      3818 non-null object
experiences_offered      3818 non-null object
neighborhood_overview      2786 non-null object
notes      2232 non-null object
transit      2884 non-null object
thumbnail_url      3498 non-null object
picture_url      3818 non-null object
xl_picture_url      3498 non-null object
host_id      3818 non-null int64
host_url      3818 non-null object
host_name      3816 non-null object
host_since      3816 non-null object
host_location      3816 non-null object
host_about      2509 non-null object
host_response_time      3295 non-null object
host_response_rate      3816 non-null object
host_is_superhost      3816 non-null object
host_thumbnail_url      3816 non-null object
host_picture_url      3816 non-null object
host_neighborhood      3818 non-null object
host_listings_count      3818 non-null float64
host_total_listings_count      3816 non-null float64
host_verifications      3818 non-null object
host_has_profile_pic      3818 non-null object
host_identity_verified      3816 non-null object
street      3818 non-null object
neighborhood      3402 non-null object
neighborhood_cleansed      3818 non-null object
neighborhood_group_cleansed      3818 non-null object
city      3818 non-null object
state      3818 non-null object
zipcode      3811 non-null object
market      3818 non-null object
smart_location      3818 non-null object
country_code      3818 non-null object
sway      3818 non-null object
latitude      3818 non-null float64
longitude      3818 non-null float64
is_location_exact      3818 non-null object
property_type      3817 non-null object
room_type      3818 non-null object
accommodates      3818 non-null int64
bathrooms      3802 non-null float64
bedrooms      3812 non-null float64
beds      3817 non-null float64
bed_type      3818 non-null object
amenities      3818 non-null object
square_feet      97 non-null float64
price      3818 non-null object
weekly_price      2099 non-null object
monthly_price      1557 non-null object
security_deposit      1866 non-null object
cleaning_fee      2788 non-null object
guests_included      3818 non-null int64
extra_people      3818 non-null object
minimum_nights      3818 non-null int64
maximum_nights      3818 non-null int64
calendar_updated      3818 non-null object
has_availability      3818 non-null object
availability_30      3818 non-null int64
availability_60      3818 non-null int64
availability_90      3818 non-null int64
availability_365      3818 non-null int64
calendar_last_scraped      3818 non-null object
number_of_reviews      3818 non-null int64
first_review      3191 non-null object
last_review      3191 non-null object
review_scores_rating      3171 non-null float64
review_scores_accuracy      3169 non-null float64
review_scores_cleanliness      3165 non-null float64
review_scores_communication      3167 non-null float64
review_scores_location      3163 non-null float64
review_scores_value      3162 non-null float64
requires_license      3818 non-null object
license      0 non-null float64
jurisdiction_names      3818 non-null object
instant_bookable      3818 non-null object
cancellation_policy      3818 non-null object
require_guest_profile_picture      3818 non-null object
picture_url      3818 non-null object
calculated_host_listings_count      3818 non-null int64
reviews_per_month      3191 non-null float64
dtypes: float64(17), int64(13), object(62)
memory usage: 2.7+ MB
None
```

```
In [4]: print(listings.describe())

count      3818.000000      id      3818.000000      host_id      host_listings_count \
mean      5.58011e+06      3.13800e+02      3.13800e+02      3816.000000
std      2.96208e+06      0.00000e+00      1.45038e+07      28.620149
min      3.35000e+06      2.01601e+13      1.57850e+07      7.107650
25%      3.26825e+06      2.01601e+13      3.27294e+06      1.000000
50%      6.13244e+06      2.01601e+13      1.05534e+07      1.000000
75%      8.09377e+06      2.01601e+13      2.59039e+07      3.800000
max      1.03401e+07      2.01601e+13      5.228861e+07      16.000000

count      host_total_listings_count      latitude      longitude      accommodates \
mean      3816.000000      3616.000000      3816.000000      3816.000000
std      28.620149      0.843852      0.031745      1.977599
min      1.000000      47.609808      122.331803      3.349389
25%      1.000000      47.699418      122.354320      2.000000
50%      1.000000      47.623891      122.326874      3.000000
75%      3.000000      47.602204      122.318889      4.000000
max      502.000000      47.733358      122.240687      16.000000

count      bathrooms      bedrooms      review_scores_rating \
mean      3802.000000      3812.000000      3171.000000
std      2.96208e+06      0.00000e+00      1.45038e+07
min      3.35000e+06      2.01601e+13      1.57850e+07
25%      3.26825e+06      2.01601e+13      3.27294e+06
50%      6.13244e+06      2.01601e+13      1.05534e+07
75%      8.09377e+06      2.01601e+13      2.59039e+07
max      1.03401e+07      2.01601e+13      5.228861e+07

count      review_scores_accuracy      review_scores_cleanliness \
mean      3169.000000      3165.000000
std      0.688031      0.586398
min      1.000000      0.000000
25%      1.000000      0.000000
50%      1.000000      1.000000
75%      1.000000      1.000000
max      1.000000      1.000000

count      review_scores_location      review_scores_value \
mean      3169.000000      3162.000000
std      0.688031      0.586398
min      1.000000      0.000000
25%      1.000000      0.000000
50%      1.000000      1.000000
75%      1.000000      1.000000
max      1.000000      1.000000

count      calculated_host_listings_count      reviews_per_month \
mean      3816.000000      3191.000000
std      2.96208e+06      0.00000e+00
min      1.000000      0.000000
25%      1.000000      0.000000
50%      1.000000      1.000000
75%      1.000000      1.000000
max      502.000000      47.733358

[8 rows x 30 columns]
```

```
In [5]: # To keep the graph of everything uniform
plt.rcParams['figure.figsize'] = (20,7)
```

```
In [ ]:
```

Performing few data manipulation to ensure data is fit for modeling

```
In [6]: listings.price = listings.price.str.strip(" ").str.replace(" ", "")
listings.monthly_price = listings.monthly_price.str.strip(" ").str.replace(" ", "")
listings.weekly_price = listings.weekly_price.str.strip(" ").str.replace(" ", "")
print(listings.price.head(3))
print(listings.monthly_price.head(3))
print(listings.weekly_price.head(3))
```

```
0      85.00
1     150.00
2      NaN
Name: price, dtype: object
0      NaN
1     3000.00
2      NaN
Name: monthly_price, dtype: object
0      NaN
1     1000.00
2      NaN
3     650.00
4      NaN
Name: weekly_price, dtype: object
```

```
In [7]: columns = ['price','bathrooms','bedrooms','beds','last_review']
null_columns =listings[columns].isnull().sum()
```

```
In [8]: # mean is calculated for price and then replaced the mean in place of nulls
listings.price = listings.price.astype(np.float64)
price.mean = listings.price.mean()
listings["price"] =listings["price"].fillna(listings["price"].mean())

#mode is calculated for last review and then replaced with the mode in place of nulls
last_review_mode = listings.last_review.mode()[0]
print(list.last_review.mode())
listings["last_review"] =listings["last_review"].fillna(listings["last_review"].mode()[0])

#mode is calculated for bathrooms and all the nulls are replaced by nulls
bathrooms_mode = listings.bathrooms.mode()
print(bathrooms.mode())
listings["bathrooms"] =listings["bathrooms"].fillna(listings["bathrooms"].mode()[0])

#mode is calculated for bedrooms and all the nulls are replaced by nulls
bedrooms_mode = listings.bedrooms.mode()
print(bedrooms.mode())
listings["bedrooms"] =listings["bedrooms"].fillna(listings["bedrooms"].mode()[0])

#mode is calculated for bedrooms and all the nulls are replaced by nulls
beds_mode = listings.beds.mode()
print(beds.mode())
listings["beds"] =listings["beds"].fillna(listings["beds"].mode()[0])

127.97616553169199
2016-01-02
0      1.0
dtype: float64
0      1.0
dtype: float64
0      1.0
dtype: float64
```

Converting all the necessary columns to appropriate data types

```
In [10]: listings.monthly_price = listings.monthly_price.astype(np.float64)
listings.weekly_price = listings.weekly_price.astype(np.float64)
listings.bathrooms = listings.bathrooms.astype(int)
listings.bedrooms = listings.bedrooms.astype(int)
listings.beds = listings.beds.astype(int)
listings.accommodates = listings.accommodates.astype(int)
```

```
In [11]: Room =listings[['room_type','review_scores_cleanliness','review_scores_communication']]
print(Room)

0      room_type      review_scores_cleanliness      review_scores_communication
1  Entire home/apt      10.0      10.0
2  Entire home/apt      10.0      10.0
3  Entire home/apt      NaN      NaN
4  Entire home/apt      9.0      10.0

3813 Entire home/apt      10.0      10.0
3814 Entire home/apt      10.0      10.0
3815 Entire home/apt      NaN      NaN
3816 Entire home/apt      NaN      NaN
3817 Entire home/apt      NaN      NaN

[3818 rows x 3 columns]
```

```
In [12]: #grouping the dataset by room to perform analysis
room =Room.groupby("room_type").mean()

#removing the index of the column
room = room.reset_index()
```

```
In [13]: room.review_scores_communication

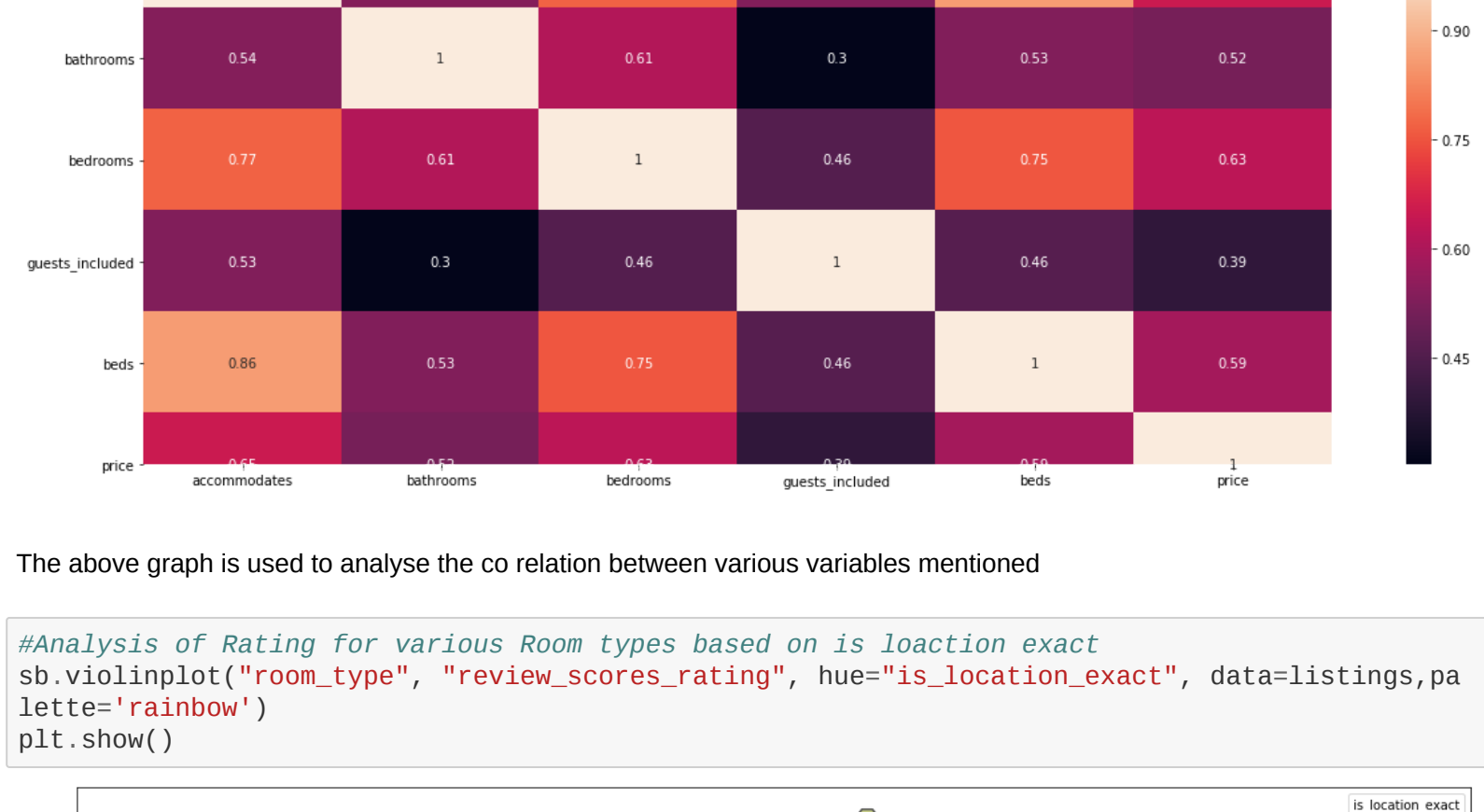
Out[13]:
0      9.808580
1      9.822292
2      9.706522
Name: review_scores_communication, dtype: float64
```

```
In [14]: plt.plot('room_type','review_scores_cleanliness', data=room, marker='o', markerfacecolor='orange', color='darkblue', linestyle='dashed', linewidth=4)
plt.plot('room_type','review_scores_communication', data=room, marker='o', markerfacecolor='orange', color='olive', linestyle='dashed', linewidth=4)
plt.legend()
plt.xlabel("Room Type")
plt.ylabel("Average Rating out of 5")
plt.title("Average Rating of Cleanliness & Communication By Room Type")
plt.show()
```



We can see from the graph that the shared room has overall Cleanliness & Communication rating less

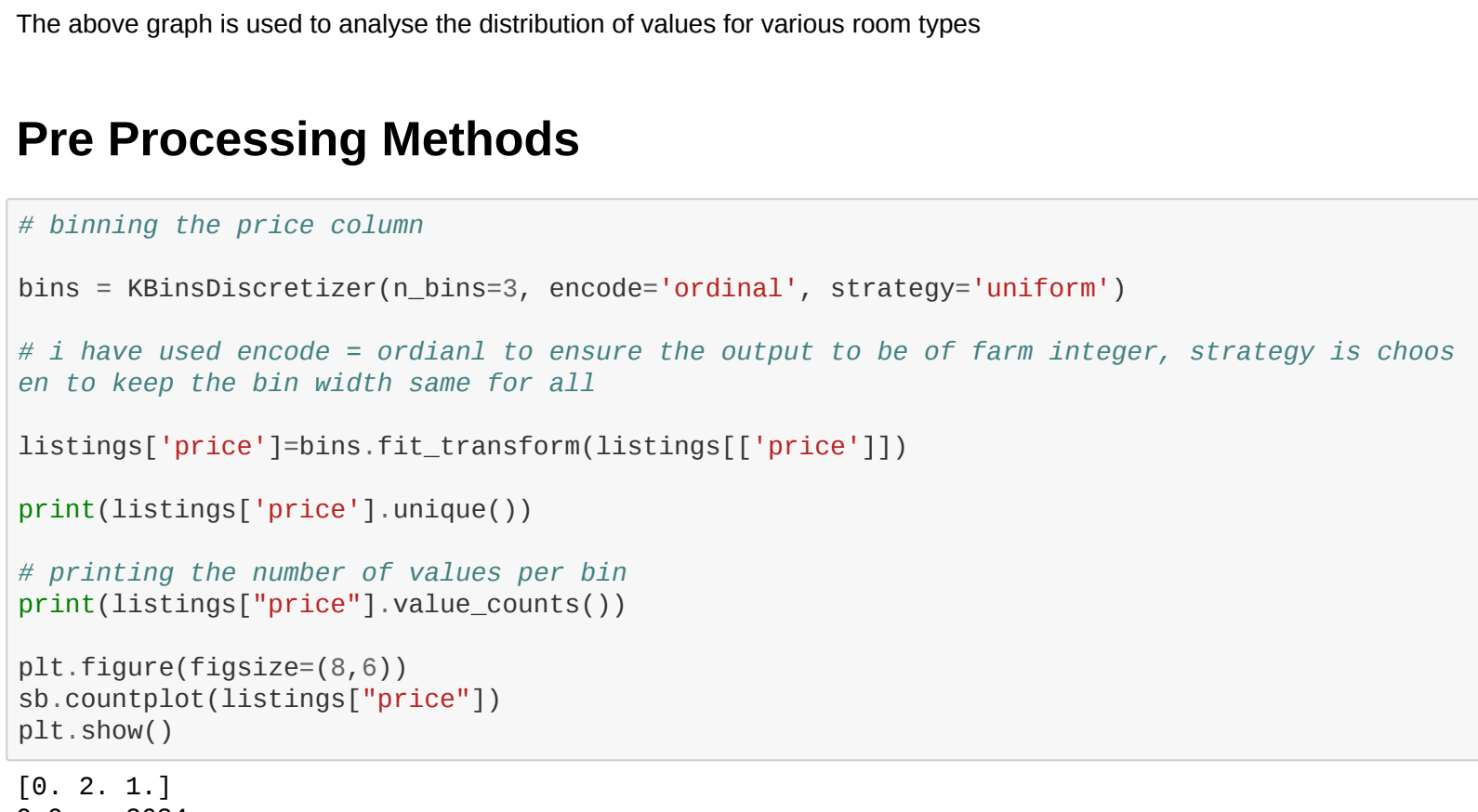
```
In [15]: fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(10,10))
# graph to analyse the count of reviews received for below mentioned categories rated
# on the scale of 1 to 10
sb.countplot(listings['review_scores_cleanliness'], ax=axes[0][0])
sb.countplot(listings['review_scores_communication'], ax=axes[0][1])
sb.countplot(listings['review_scores_location'], ax=axes[1][0])
sb.countplot(listings['review_scores_value'], ax=axes[1][1])
plt.title("Analysis of Rating provided for various factors")
plt.show()
```



This above graphs is used to analyse the total no of records that are having very less ratings

```
In [16]: property =listings[['beds','price','property_type']]
price= property.groupby('property_type', 'property_type').agg({'beds': 'mean', 'price': 'mean'})

plt.figure(figsize=(13,8))
sb.scatterplot(x='beds', y='price', hue='property_type', data=price, s=150, palette='rainbow')
plt.title("Analysis of no of beds vs Avg Price by property type")
plt.xlabel("No of beds")
plt.ylabel("Average price")
plt.show()
```



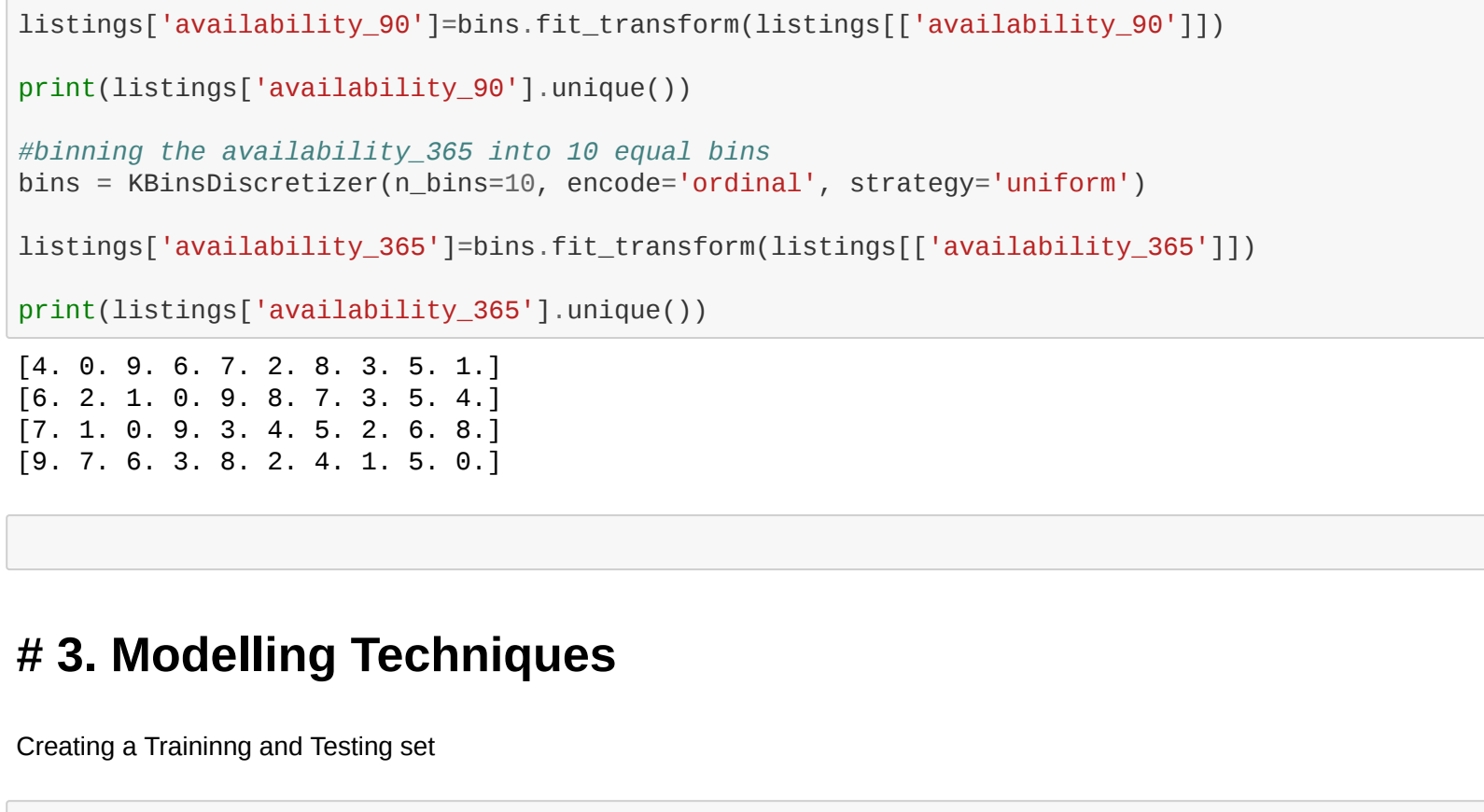
The above graph is used to Analysis of no of beds vs Avg Price by property type

```
In [17]: columns = ['accommodates','bathrooms','bedrooms','guests_included','beds','price']
sb.heatmap(listings[columns].corr(), annot=True)
plt.show()
```



The above graph is used to analyse the co relation between various variables mentioned

```
In [18]: #Analysis of Rating for various Room types based on is location exact
sb.violinplot('room_type', review_scores_rating, hue='is_location_exact', data=listings, palette='rainbow')
plt.show()
```



The above graph is used to analyse the distribution of values for various room types

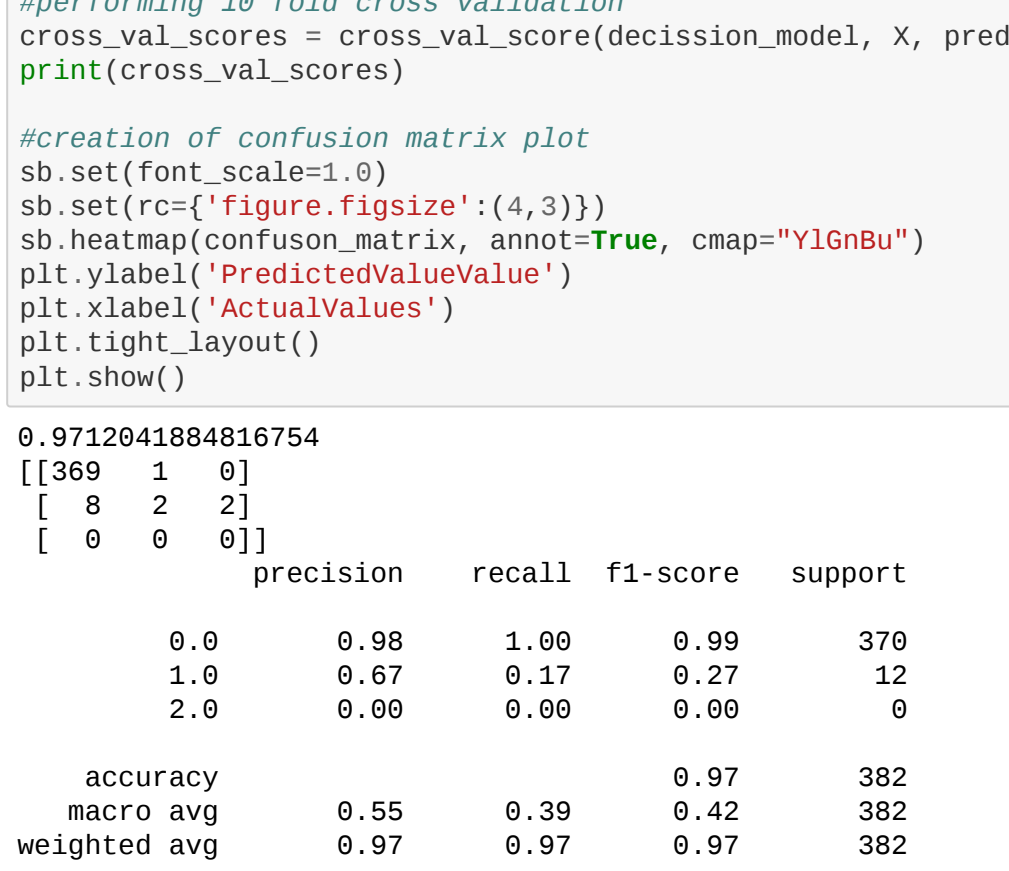
Pre Processing Methods

```
In [19]: # binning the price column
bins = KBinsDiscretizer(n_bins=3, encode='ordinal', strategy='uniform')
# i have used encode = ordinal to ensure the output to be of farm integer, strategy is choos
# en to keep the bin width same for all
listings['price'] =bins.fit_transform(listings[['price']])

print(listings['price'].unique())

# printing the number of values per bin
print(listings['price'].value_counts())

plt.figure(figsize=(8,6))
sb.countplot(listings['price'])
plt.show()
```



```
In [20]: #first converting the last review column to date farm and then subtracting it from todays da
# te to calculate no of days
listings['last_review'] = pd.to_datetime(listings['last_review']).dt.date
listings['no_days'] = (pd.datetime.now().date()-listings['last_review']).dt.days

print(listings[['last_review','no_days']].head())

last_review      no_days
0  2016-01-02      1591
1  2015-12-29      1555
2  2015-09-03      1672
3  2016-01-02      1551
4  2015-10-24      1621
```

```
In [21]: #binning the availability_30 into 10 equal bins
bins = KBinsDiscretizer(n_bins=10, encode='ordinal', strategy='uniform')
listings['availability_30'] =bins.fit_transform(listings[['availability_30']])

#binning the availability_60 into 10 equal bins
bins = KBinsDiscretizer(n_bins=10, encode='ordinal', strategy='uniform')
listings['availability_60'] =bins.fit_transform(listings[['availability_60']])

print(listings['availability_60'].unique())

#binning the availability_90 into 10 equal bins
bins = KBinsDiscretizer(n_bins=10, encode='ordinal', strategy='uniform')
listings['availability_90'] =bins.fit_transform(listings[['availability_90']])

print(listings['availability_90'].unique())

#binning the availability_365 into 10 equal bins
bins = KBinsDiscretizer(n_bins=10, encode='ordinal', strategy='uniform')
listings['availability_365'] =bins.fit_transform(listings[['availability_365']])

print(listings['availability_365'].unique())

[0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
[0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
```

```
In [ ]:
```

3. Modelling Techniques

Creating a Training and Testing set

```
In [22]: X=listings[['accommodates','bathrooms','bedrooms','beds']]
predictor=KBinsDiscretizer(n_bins=3, encode='ordinal', strategy='uniform')

X_train,X_test,y_train,y_test = train_test_split(X,predictor, test_size=0.1, random_stat
e=10)

print("number of training samples (%)" %format(X_train.shape))
print("number of testing (%)" %format(y_train.shape))

number of training samples (3436, 4)
number of testing (3436, 4)
```

```
In [23]: # 1. Applying Decision Tree algorithm to predict the price based on the given training varia
# bles
decision=DecisionTreeClassifier()
sb.fit(X_train,y_train,y_train,X_test,y_test)
predict=decision.predict(X_test)
print(accuracy_score(y_test, predict))

# creating confusion matrix
confusion_matrix=confusion_matrix(y_test, predict)
print(confusion_matrix)
print(classification_report(y_test, predict))

#performing 10 fold cross validation
cross_val_scores = cross_val_score(decision_model, X, predictor, cv=10)
print(cross_val_scores)

#creation of confusion matrix plot
sb.set(font_scale=1.0)
sb.set(rc={'figure.figsize':(4,3)})
sb.heatmap(confusion_matrix, annot=True, cmap='YlGnBu')
plt.xlabel('ActualValues')
plt.ylabel('PredictedValue')
plt.tight_layout()
plt.show()

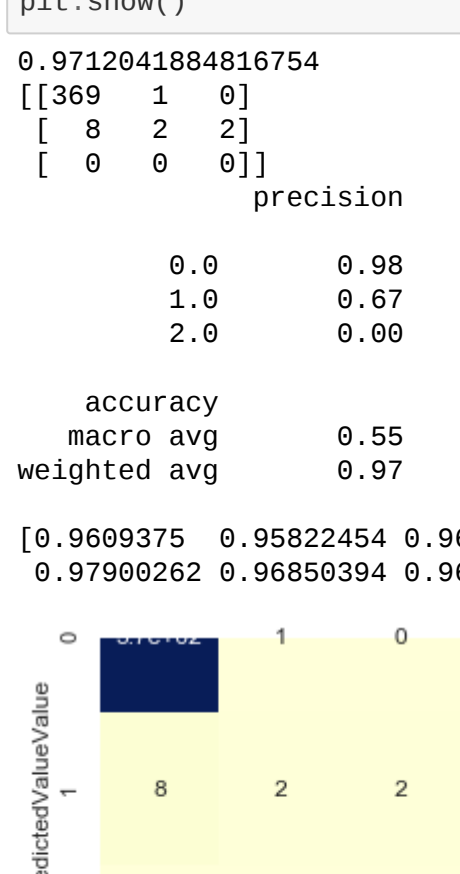
0.9712041884816754
[[369 1 0]
 [ 0 2 1]
 [ 0 0 0]]

precision      recall      f1-score      support

0.0      0.00      0.00      1.00      0.09      370
1.0      0.07      0.17      0.27      12      12
2.0      0.00      0.00      0.00      0.00      0

accuracy      0.97      0.97      0.97      382
macro avg      0.55      0.39      0.42      382
weighted avg      0.97      0.97      0.97      382

[0.96934167 0.95822454 0.96605744 0.97087759 0.96689394 0.96587927 0.96325459
0.97908262 0.96589394 0.96587927 0.97112861]
```



```
In [24]: #Random Forest Classifier
random = RandomForestClassifier(random.state = 0)
# training the model with training data
rand_forest = random.fit(X_train,y_train)
# prediction for test data
rand_forest_pred = rand_forest.predict(X_test)
print(accuracy_score(y_test, rand_forest_pred))

# 10 fold validation
cv_randforest = cross_val_validate(rand_forest, X, predictor, cv=10)
print(cv_randforest)

confusion_matrix_randforest = confusion_matrix(y_test, rand_forest_pred)
print(confusion_matrix_randforest)
# plotting of the confusion matrix
sb.set(font_scale=1.0)
sb.set(rc={'figure.figsize':(4,3)})
sb.heatmap(confusion_matrix_randforest, annot=True, cmap='YlGnBu')
plt.xlabel('PredictedValue')
plt.ylabel('ActualValues')
plt.tight_layout()
plt.show()

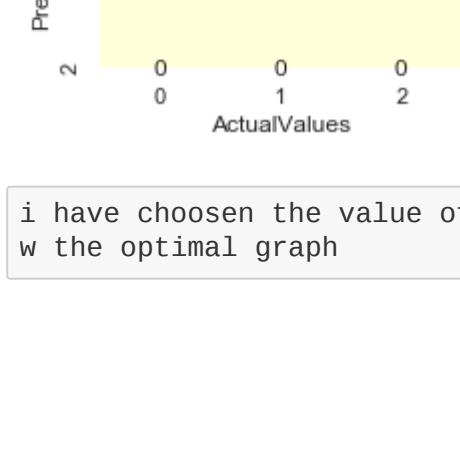
0.9712041884816754
[[369 1 0]
 [ 0 2 1]
 [ 0 0 0]]

precision      recall      f1-score      support

0.0      0.00      0.00      1.00      0.09      370
1.0      0.07      0.17      0.27      12      12
2.0      0.00      0.00      0.00      0.00      0

accuracy      0.97      0.97      0.97      382
macro avg      0.55      0.39      0.42      382
weighted avg      0.97      0.97      0.97      382

[0.96934167 0.95822454 0.96605744 0.97087759 0.96689394 0.96587927 0.96325459
0.97908262 0.96589394 0.96587927 0.97112861]
```



```
In [ ]:
```

```
# 1. Applying Knn to predict the price based on the given training variables
k=KNeighborsClassifier(n_neighbors = 11)
knn_model= decision.fit(X_train,y_train)
predict= decision.predict(X_test)
print(accuracy_score(y_test, predict))

# creating confusion matrix
confusion_matrix=confusion_matrix(y_test, predict)
print(confusion_matrix)
print(classification_report(y_test, predict))

#performing 10 fold cross validation
cross_val_scores = cross_val_score(knn_model, X, predictor, cv=10)
print(cross_val_scores)

#creation of confusion matrix plot
sb.set(font_scale=1.0)
sb.set(rc={'figure.figsize':(4,3)})
sb.heatmap(confusion_matrix, annot=True, cmap='YlGnBu')
plt.xlabel('PredictedValue')
plt.ylabel('ActualValues')
plt.tight_layout()
plt.show()

0.9712041884816754
[[369 1 0]
 [ 0 2 1]
 [ 0 0 0]]

precision      recall      f1-score      support

0.0      0.00      0.00      1.00      0.09      370
1.0      0.07      0.17      0.27      12      12
2.0      0.00      0.00      0.00      0.00      0

accuracy      0.97      0.97      0.97      382
macro avg      0.55      0.39      0.42      382
weighted avg      0.97      0.97      0.97      382

[0.96934167 0.95822454 0.96605744 0.97087759 0.96689394 0.96587927 0.96325459
0.97908262 0.96589394 0.96587927 0.97112861]
```



```
In [ ]:
```

```
# 1. Applying Knn to predict the price based on the given training variables
k=KNeighborsClassifier(n_neighbors = 11)
knn_model= decision.fit(X_train,y_train)
predict= decision.predict(X_test)
print(accuracy_score(y_test, predict))

# creating confusion matrix
confusion_matrix=confusion_matrix(y_test, predict)
print(confusion_matrix)
print(classification_report(y_test, predict))

#performing 10 fold cross validation
cross_val_scores = cross_val_score(knn_model, X, predictor, cv=10)
print(cross_val_scores)

#creation of confusion matrix plot
sb.set(font_scale=1.0)
sb.set(rc={'figure.figsize':(4,3)})
sb.heatmap(confusion_matrix, annot=True, cmap='YlGnBu')
plt.xlabel('PredictedValue')
plt.ylabel('ActualValues')
plt.tight_layout()
plt.show()

0.9712041884816754
[[369 1 0]
 [ 0 2 1]
 [ 0 0 0]]

precision      recall      f1-score      support

0.0      0.00      0.00      1.00      0.09      370
1.0      0.07      0.17      0.27      12      12
2.0      0.00      0.00      0.00      0.00      0

accuracy      0.97      0.97      0.97      382
macro avg      0.55      0.39      0.42      382
weighted avg      0.97      0.97      0.97      382

[0.96934167 0.95822454 0.96605744 0.97087759 0.96689394 0.96587927 0.96325459
0.97908262 0.96589394 0.96587927 0.97112861]
```