

In [1]: *#importation of necessary libraries*

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

import warnings
warnings.filterwarnings('ignore')
```

In [2]: *#Load Housing Data*

```
raw_data = pd.read_csv('dataset/Housing_dataset_train.csv')
test_data = pd.read_csv('dataset/Housing_dataset_test.csv')
submission = pd.read_csv("dataset/Sample_submission.csv")
```

Basic Data Exploration

In [3]: raw_data.head()

Out[3]:

	ID	loc	title	bedroom	bathroom	parking_space	price
0	3583	Katsina	Semi-detached duplex	2.0	2.0	1.0	1149999.565
1	2748	Ondo	Apartment	NaN	2.0	4.0	1672416.689
2	9261	Ekiti	NaN	7.0	5.0	NaN	3364799.814
3	2224	Anambra	Detached duplex	5.0	2.0	4.0	2410306.756
4	10300	Kogi	Terrace duplex	NaN	5.0	6.0	2600700.898

In [4]: raw_data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14000 entries, 0 to 13999
Data columns (total 7 columns):
#   Column          Non-Null Count  Dtype
---  -
0   ID              14000 non-null  int64
1   loc             12187 non-null  object
2   title           12278 non-null  object
3   bedroom         12201 non-null  float64
4   bathroom        12195 non-null  float64
5   parking_space   12189 non-null  float64
6   price           14000 non-null  float64
dtypes: float64(4), int64(1), object(2)
memory usage: 765.8+ KB
```

In [5]: test_data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6000 entries, 0 to 5999
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  -
0   ID              6000 non-null   int64
1   loc             6000 non-null   object
2   title           6000 non-null   object
3   bedroom         6000 non-null   int64
4   bathroom        6000 non-null   int64
5   parking_space   6000 non-null   int64
dtypes: int64(4), object(2)
memory usage: 281.4+ KB
```

In [6]: raw_data.describe()

Out[6]:

	ID	bedroom	bathroom	parking_space	price
count	14000.000000	12201.000000	12195.000000	12189.000000	1.400000e+04
mean	4862.700357	4.308171	3.134235	3.169825	2.138082e+06
std	3818.348214	2.441165	2.035950	1.599415	1.083057e+06
min	0.000000	1.000000	1.000000	1.000000	4.319673e+05
25%	1672.750000	2.000000	1.000000	2.000000	1.393990e+06
50%	3527.000000	4.000000	2.000000	3.000000	1.895223e+06
75%	8011.250000	6.000000	5.000000	4.000000	2.586699e+06
max	12999.000000	9.000000	7.000000	6.000000	1.656849e+07

Data Cleaning

```
In [7]: #drop columns not needed
raw_data.dropna(subset=['loc', 'title'], inplace=True)
raw_data.reset_index()
```

Out[7]:

	index	ID	loc	title	bedroom	bathroom	parking_space	price
0	0	3583	Katsina	Semi-detached duplex	2.0	2.0	1.0	1149999.565
1	1	2748	Ondo	Apartment	NaN	2.0	4.0	1672416.689
2	3	2224	Anambra	Detached duplex	5.0	2.0	4.0	2410306.756
3	4	10300	Kogi	Terrace duplex	NaN	5.0	6.0	2600700.898
4	5	1733	Borno	Mansion	NaN	1.0	3.0	1341750.867
...
10526	13994	10477	Taraba	Detached duplex	8.0	1.0	6.0	2837199.086
10527	13995	6175	Edo	Bungalow	NaN	7.0	NaN	2367927.861
10528	13996	9704	Kaduna	Apartment	NaN	7.0	5.0	2228516.471
10529	13997	11190	Plateau	Bungalow	8.0	6.0	5.0	2406812.693
10530	13998	9256	Delta	Flat	NaN	6.0	1.0	3348918.718

10531 rows × 8 columns

```
In [8]: #check for duplicated rows
raw_data.duplicated().sum()
```

Out[8]: 0

```
In [9]: #check for null values
raw_data.isnull().sum()
```

```
Out[9]: ID                0
loc                  0
title                0
bedroom            1675
bathroom           1672
parking_space      1671
price                0
dtype: int64
```

```
In [10]: test_data.isnull().sum()
```

```
Out[10]: ID                0
loc                  0
title                0
bedroom            0
bathroom           0
parking_space      0
dtype: int64
```

```
In [11]: raw_data.groupby(['title', 'loc'])[['bedroom', 'bathroom', 'parking_space']].mo
```

Out[11]:

		bedroom	bathroom	parking_space
title	loc			
Apartment	Abia	5.0	3.0	4.0
	Adamawa	4.0	3.0	3.0
	Akwa Ibom	4.0	3.0	4.0
	Anambra	4.0	2.0	3.0
	Bauchi	3.0	3.0	3.0
...
Townhouse	Rivers	4.0	3.0	3.0
	Sokoto	4.0	3.0	3.0
	Taraba	6.0	3.0	3.0
	Yobe	3.0	3.0	3.0

```
In [12]: raw_data.groupby(['title', 'loc'])[['bedroom', 'bathroom', 'parking_space']].mo
```

Out[12]:

		bedroom	bathroom	parking_space
title	loc			
Apartment	Abia	4.0	2.0	4.0
	Adamawa	4.0	2.0	4.0
	Akwa Ibom	4.0	2.0	4.0
	Anambra	4.5	2.0	3.0
	Bauchi	3.0	2.0	3.0
...
Townhouse	Rivers	5.0	2.0	2.5
	Sokoto	4.0	3.0	3.0
	Taraba	6.0	3.0	4.0
	Yobe	3.0	2.0	2.0
	Zamfara	2.5	2.0	3.0

360 rows × 3 columns

```
In [13]: #Proceeded to fill na values with the median subset by location and type of house
filled_median = raw_data
na_cols = ['bedroom', 'bathroom', 'parking_space']

for col in na_cols:
    median_data = filled_median.groupby(['title', 'loc'])[col].transform('median')
    filled_median[col] = filled_median[col].fillna(median_data)
```

```
In [14]: filled_median
```

Out[14]:

	ID	loc	title	bedroom	bathroom	parking_space	price
0	3583	Katsina	Semi-detached duplex	2.0	2.0	1.0	1149999.565
1	2748	Ondo	Apartment	3.5	2.0	4.0	1672416.689
3	2224	Anambra	Detached duplex	5.0	2.0	4.0	2410306.756
4	10300	Kogi	Terrace duplex	5.0	5.0	6.0	2600700.898
5	1733	Borno	Mansion	4.0	1.0	3.0	1341750.867
...
13994	10477	Taraba	Detached duplex	8.0	1.0	6.0	2837199.086
13995	6175	Edo	Bungalow	4.0	7.0	4.0	2367927.861
13996	9704	Kaduna	Apartment	4.0	7.0	5.0	2228516.471
13997	11190	Plateau	Bungalow	8.0	6.0	5.0	2406812.693
13998	9256	Delta	Flat	3.0	6.0	1.0	3348918.718

10531 rows × 7 columns

```
In [15]: filled_mean = raw_data
na_cols = ['bedroom', 'bathroom', 'parking_space']

for col in na_cols:
    mean_data = filled_mean.groupby(['title', 'loc'])[col].transform('mean')
    filled_mean[col] = filled_mean[col].fillna(mean_data)
```

```
In [16]: filled_mean
```

```
Out[16]:
```

	ID	loc	title	bedroom	bathroom	parking_space	price
0	3583	Katsina	Semi-detached duplex	2.0	2.0	1.0	1149999.565
1	2748	Ondo	Apartment	3.5	2.0	4.0	1672416.689
3	2224	Anambra	Detached duplex	5.0	2.0	4.0	2410306.756
4	10300	Kogi	Terrace duplex	5.0	5.0	6.0	2600700.898
5	1733	Borno	Mansion	4.0	1.0	3.0	1341750.867
...
13994	10477	Taraba	Detached duplex	8.0	1.0	6.0	2837199.086
13995	6175	Edo	Bungalow	4.0	7.0	4.0	2367927.861
13996	9704	Kaduna	Apartment	4.0	7.0	5.0	2228516.471
13997	11190	Plateau	Bungalow	8.0	6.0	5.0	2406812.693
13998	9256	Delta	Flat	3.0	6.0	1.0	3348918.718

10531 rows × 7 columns

```
In [17]: final_data = filled_median
```

Data Processing

```
In [18]: final_data['bedbath_ratio'] = final_data['bathroom'] / final_data['bedroom']
test_data['bedbath_ratio'] = test_data['bathroom'] / test_data['bedroom']
final_data['total_rooms'] = final_data['bathroom'] + final_data['bedroom']
test_data['total_rooms'] = test_data['bathroom'] + test_data['bedroom']
```

```
In [19]: # creating a log transformation
final_data['log_price'] = np.log(final_data['price'])
```

```
In [20]: zones = {
    'Benue': 'North Central',
    'Kogi': 'North Central',
    'Kwara': 'North Central',
    'Nasarawa': 'North Central',
    'Niger': 'North Central',
    'Plateau': 'North Central',
    'Adamawa': 'North East',
    'Bauchi': 'North East',
    'Borno': 'North East',
    'Gombe': 'North East',
    'Taraba': 'North East',
    'Yobe': 'North East',
    'Jigawa': 'North West',
    'Kaduna': 'North West',
    'Kano': 'North West',
    'Katsina': 'North West',
    'Kebbi': 'North West',
    'Sokoto': 'North West',
    'Zamfara': 'North West',
    'Abia': 'South East',
    'Anambra': 'South East',
    'Ebonyi': 'South East',
    'Enugu': 'South East',
    'Imo': 'South East',
    'Akwa Ibom': 'South South',
    'Bayelsa': 'South South',
    'Cross River': 'South South',
    'Delta': 'South South',
    'Edo': 'South South',
    'Rivers': 'South South',
    'Ekiti': 'South West',
    'Lagos': 'South West',
    'Ogun': 'South West',
    'Ondo': 'South West',
    'Osun': 'South West',
    'Oyo': 'South West',
}


# Map the location to values based on zone
final_data['zone'] = final_data['loc'].map(zones)
test_data['zone'] = test_data['loc'].map(zones)

# One-hot encode the 'Region' column
final_data = pd.get_dummies(final_data, columns=['zone'], prefix='', prefix_sep='')
test_data = pd.get_dummies(test_data, columns=['zone'], prefix='', prefix_sep='')

# Print the updated dataframe
final_data.head()
```

Out[20]:

	ID	loc	title	bedroom	bathroom	parking_space	price	bedbath_ratio	t
0	3583	Katsina	Semi-detached duplex	2.0	2.0	1.0	1149999.565	1.000000	
1	2748	Ondo	Apartment	3.5	2.0	4.0	1672416.689	0.571429	
3	2224	Anambra	Detached duplex	5.0	2.0	4.0	2410306.756	0.400000	
4	10300	Kogi	Terrace duplex	5.0	5.0	6.0	2600700.898	1.000000	
5	1733	Borno	Mansion	4.0	1.0	3.0	1341750.867	0.250000	

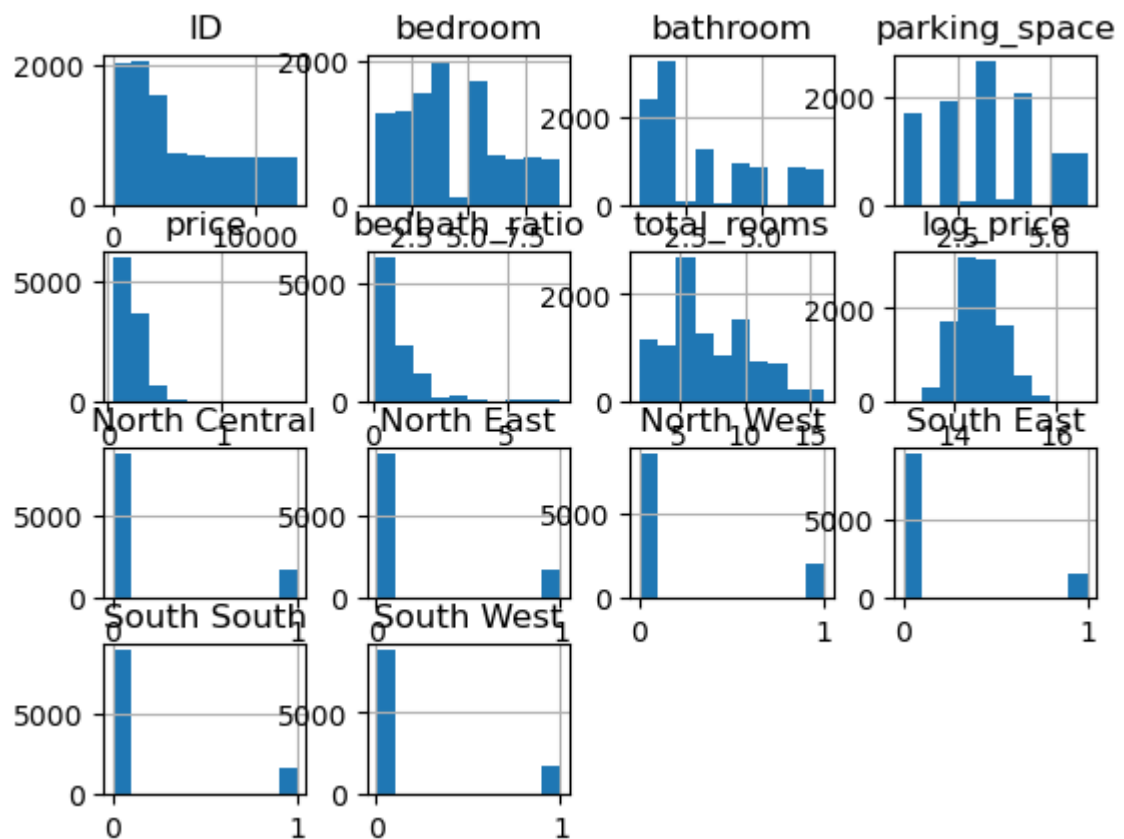


In [21]: `columns = final_data.columns`

Explanatory Data Analysis


```
In [22]: final_data.hist()
```

```
Out[22]: array([[<Axes: title={'center': 'ID'}>,
  <Axes: title={'center': 'bedroom'}>,
  <Axes: title={'center': 'bathroom'}>,
  <Axes: title={'center': 'parking_space'}>],
 [ <Axes: title={'center': 'price'}>,
  <Axes: title={'center': 'bedbath_ratio'}>,
  <Axes: title={'center': 'total_rooms'}>,
  <Axes: title={'center': 'log_price'}>],
 [ <Axes: title={'center': 'North Central'}>,
  <Axes: title={'center': 'North East'}>,
  <Axes: title={'center': 'North West'}>,
  <Axes: title={'center': 'South East'}>],
 [ <Axes: title={'center': 'South South'}>,
  <Axes: title={'center': 'South West'}>, <Axes: >, <Axes: >]],
 dtype=object)
```



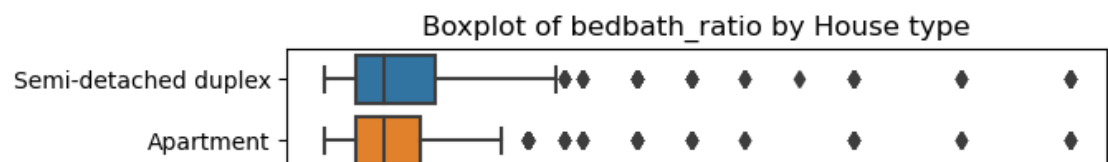
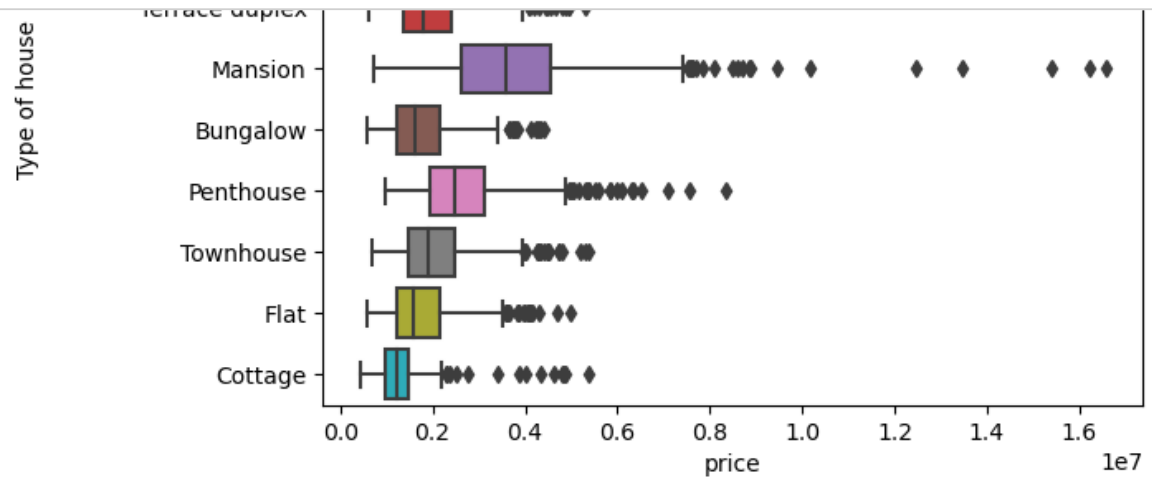
```
In [23]: # Get all possible categories for the categorical columns
cat_col = ['loc', 'title']
for name in cat_col:
    print(name, ':')
    print(raw_data[name].value_counts(), '\n')
```

```
loc :
Cross River      317
Imo              311
Anambra          310
Benue            309
Kaduna           309
Yobe             307
Zamfara          307
Borno            306
Plateau          301
Kano             301
Oyo              301
Ondo             300
Ogun             298
Ebonyi           298
Niger            297
Gombe           296
Kebbi            295
Nasarawa         295
Katsina          292
Jigawa           289
Enugu            288
Bauchi           288
Sokoto           287
Ekiti            286
Osun             286
Adamawa          285
Kwara            285
Bayelsa          284
Taraba           281
Kogi             279
Rivers           278
Delta            277
Abia             276
Lagos            274
Akwa Ibom        273
Edo              265
Name: loc, dtype: int64
```

```
title :
Flat              1182
Apartment         1147
Townhouse         1139
Semi-detached duplex 1133
Mansion           1125
Detached duplex   1115
Penthouse         1103
Bungalow          1102
Terrace duplex    1095
Cottage           390
Name: title, dtype: int64
```

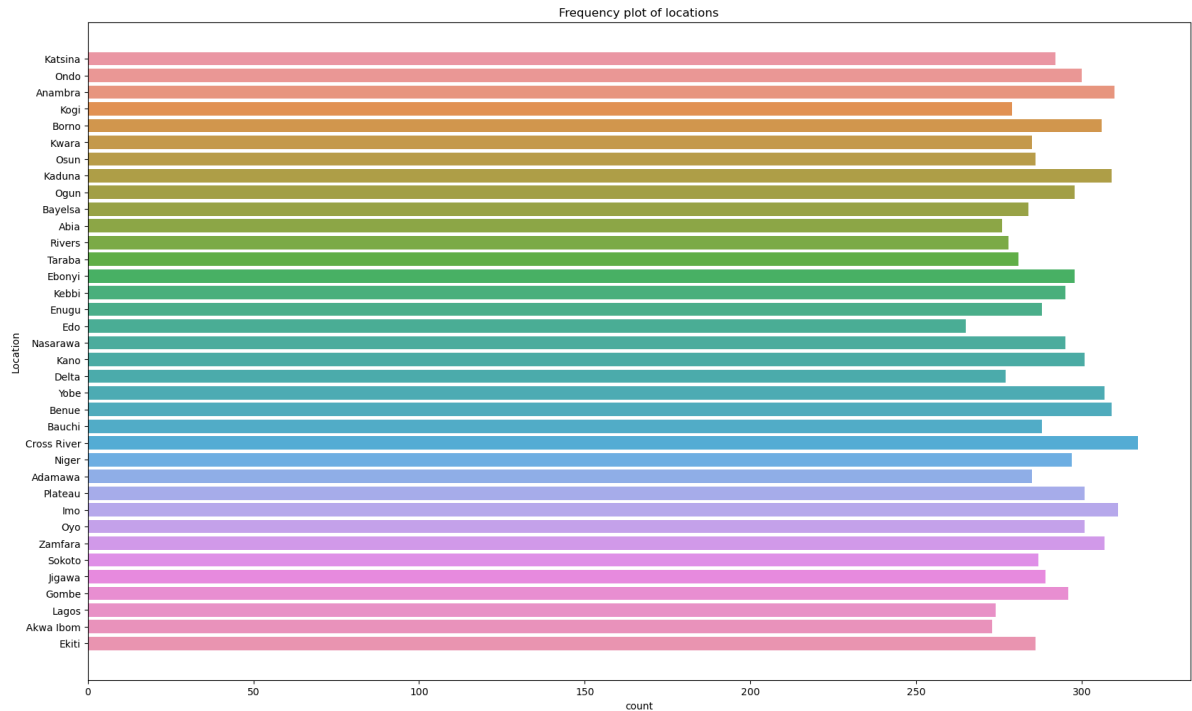
In [24]: *#shows relationship between other columns and type of house*

```
for col in columns[3:9]:  
    sns.boxplot(data=final_data, x=col, y='title')  
    plt.title('Boxplot of {} by House type'.format(col))  
    plt.xlabel('{}'.format(col))  
    plt.ylabel('Type of house')  
    plt.show()
```



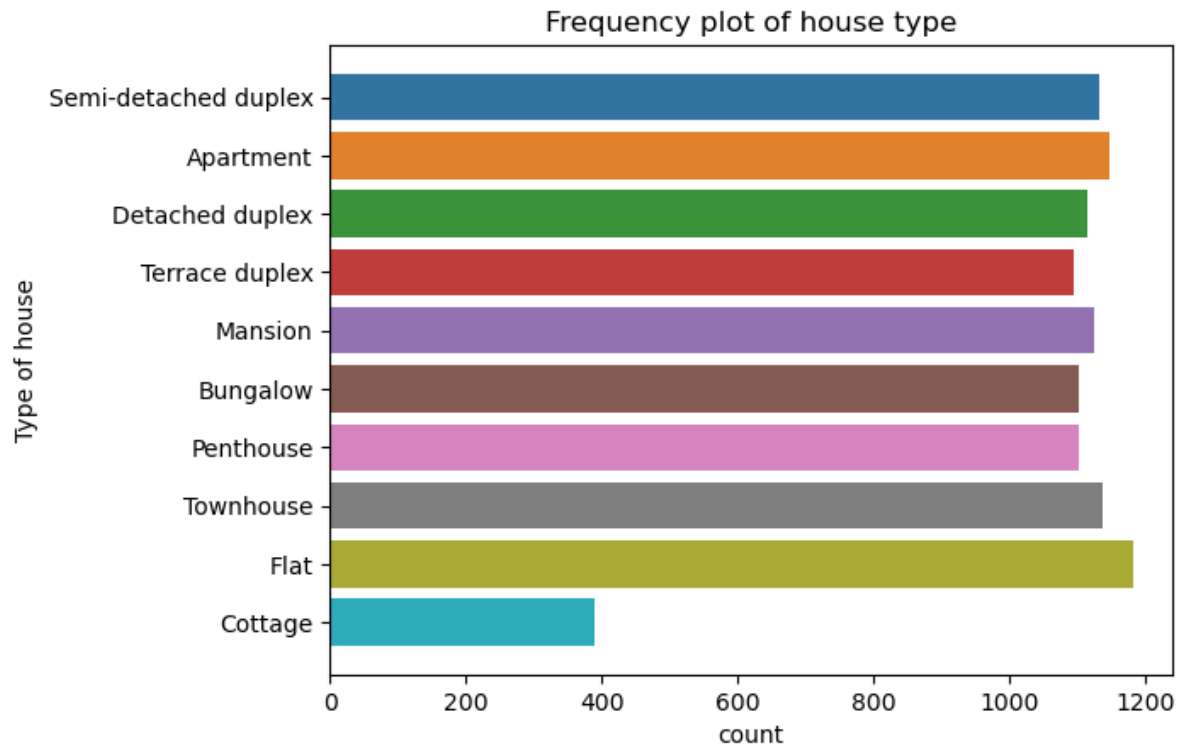
```
In [25]: #count of rows by location
plt.figure(figsize=(20, 12));
sns.countplot(data=final_data, y='loc' )
plt.title("Frequency plot of locations")
plt.ylabel("Location")
plt.plot()
```

Out[25]: []



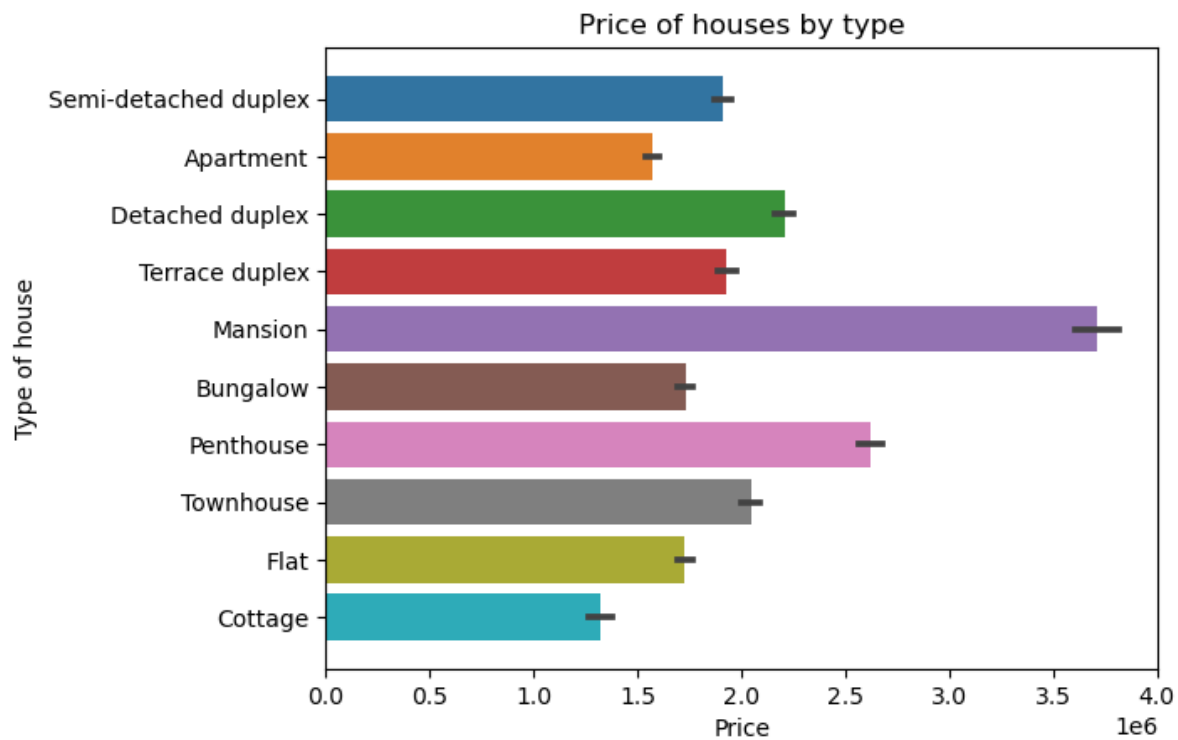
```
In [26]: #count of rows by type of house
sns.countplot(data=final_data, y='title')
plt.title("Frequency plot of house type")
plt.ylabel("Type of house")
plt.plot()
```

Out[26]: []



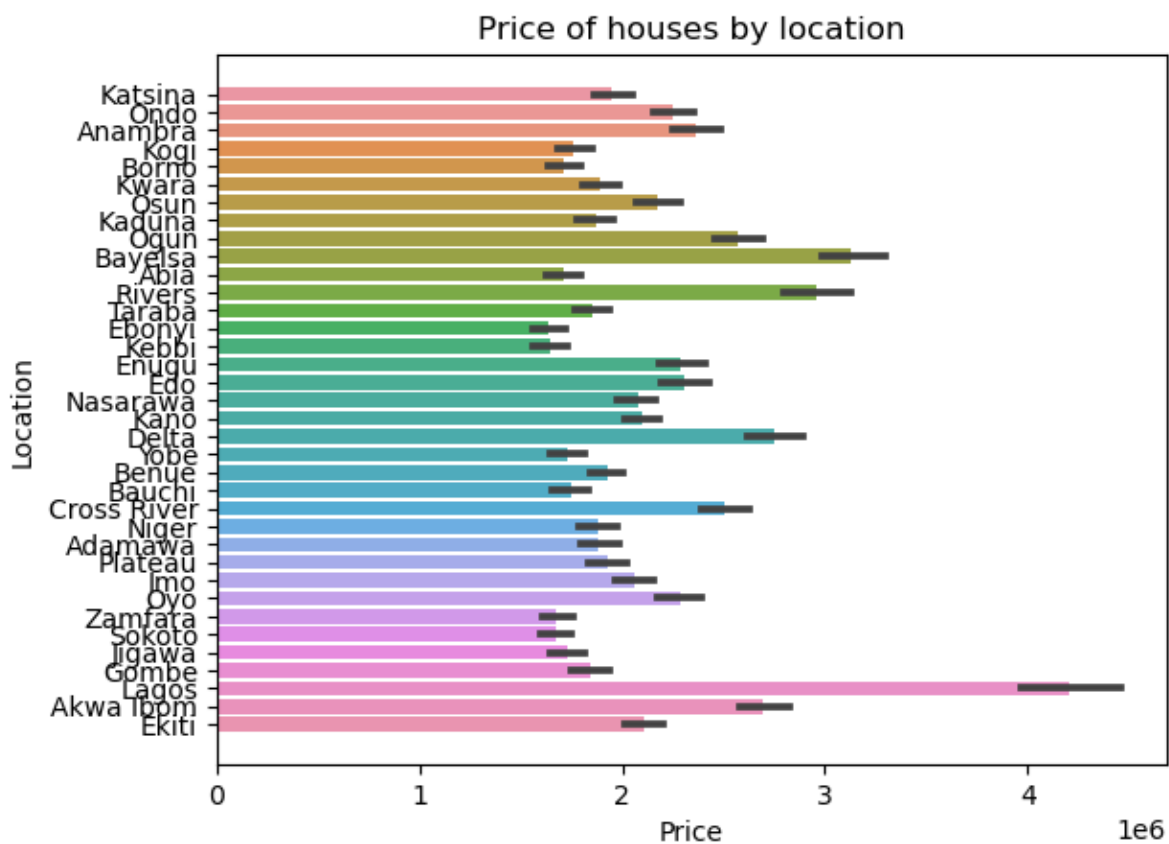
```
In [27]: #relationship between house type and price
sns.barplot(data=final_data, y='title', x='price')
plt.title("Price of houses by type")
plt.ylabel("Type of house")
plt.xlabel("Price")
plt.plot()
```

Out[27]: []



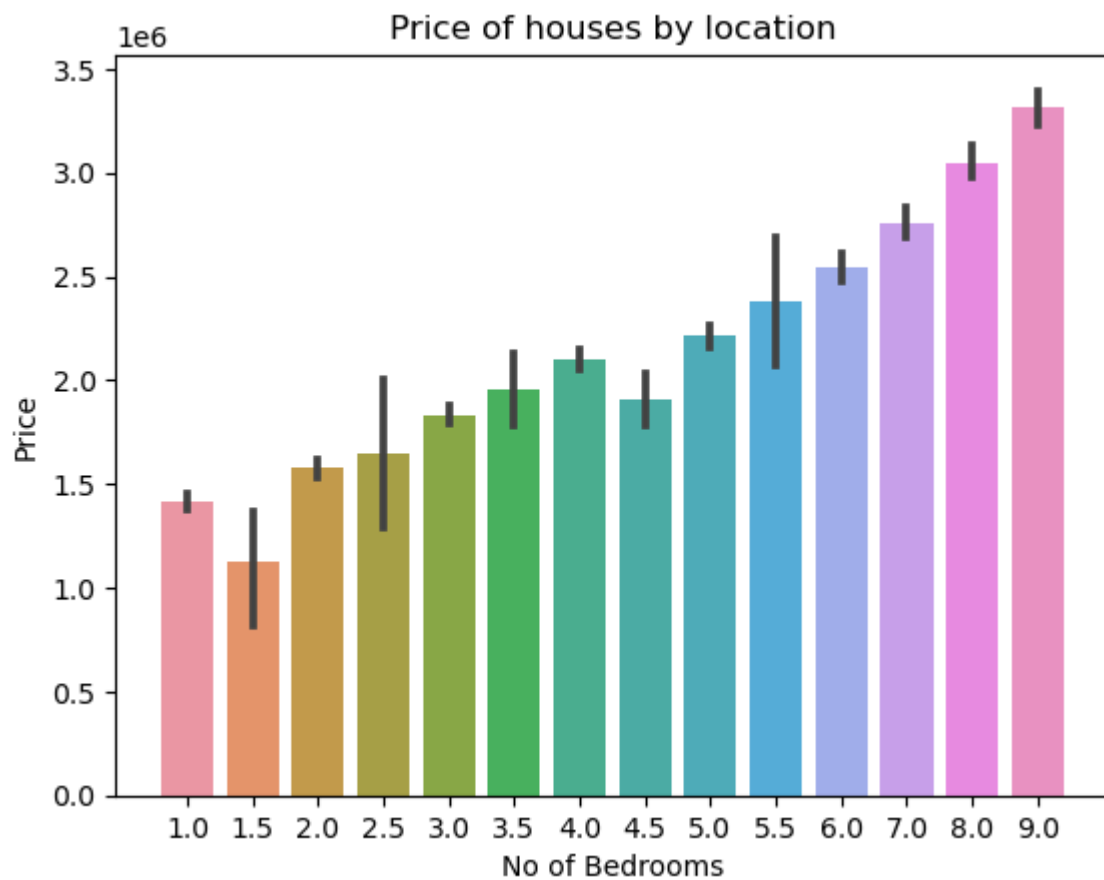
```
In [28]: #Relationship between location and price
sns.barplot(data=final_data, y='loc', x='price')
plt.title("Price of houses by location")
plt.ylabel("Location")
plt.xlabel("Price")
plt.plot()
```

Out[28]: []



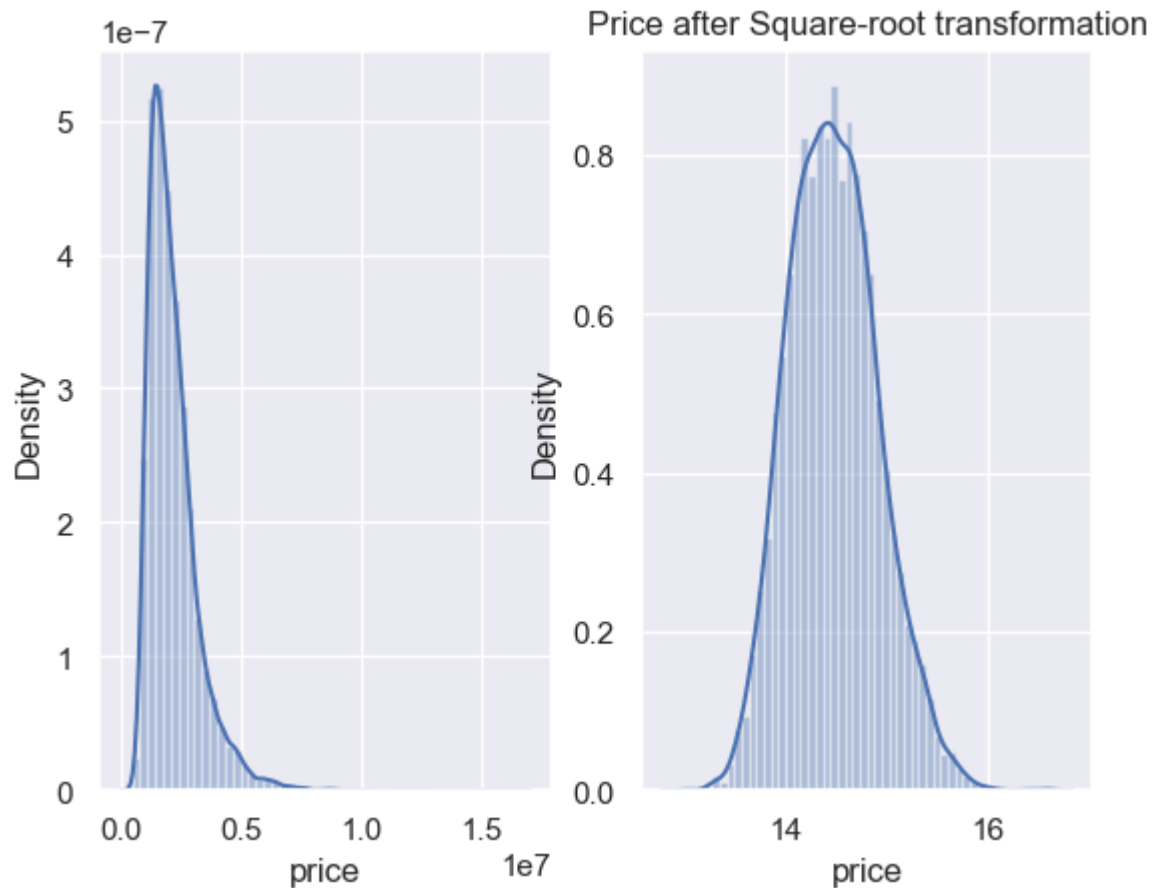

```
In [29]: #relationship between no of bedroom and price
sns.barplot(data=final_data, x='bedroom', y='price')
plt.title("Price of houses by location")
plt.xlabel("No of Bedrooms")
plt.ylabel("Price")
plt.plot()
```

Out[29]: []



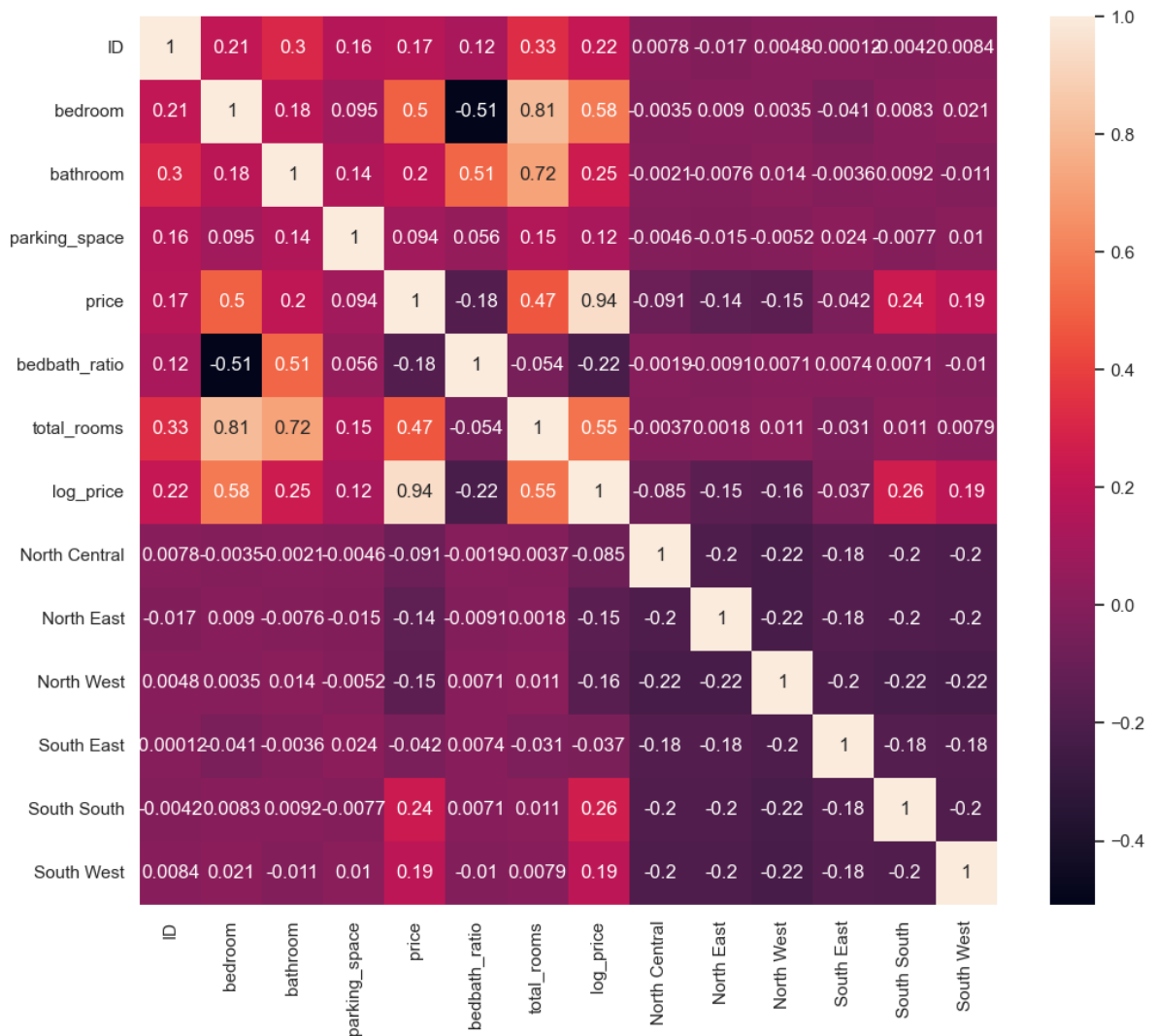
```
In [30]: sns.set()
y = final_data.price
y_transformed = pd.Series(np.log(y))

fig, ax = plt.subplots(1, 2)
sns.distplot(y, ax=ax[0])
plt.title("Price after Square-root transformation")
# ax[0].axvline(y_transformed)
sns.distplot(y_transformed, ax=ax[1])
plt.show()
```



```
In [31]: #plotting corr map
corr = final_data.corr()
plt.figure(figsize = (12,10))
sns.heatmap(corr, annot=True)
```

Out[31]: <Axes: >



Data Processing

```
In [32]: all_data= final_data.drop(columns=['price'], axis=1).append(test_data)
all_data.shape
```

Out[32]: (16531, 15)

```

In [33]: # Define the ranking based on size (arranged from smallest to biggest)
house_type_ranks = {
    "Apartment":1,
    "Flat":2,
    "Cottage":3,
    "Terrace duplex":4,
    "Bungalow":5,
    "Semi-detached duplex":6,
    "Townhouse":7,
    "Detached duplex":8,
    "Penthouse":9,
    "Mansion":10,
}

# Map the house types to numerical values based on size ranking
all_data['title'] = all_data['title'].map(house_type_ranks)

# Print the updated dataframe
final_data.head()

```

Out[33]:

	ID	loc	title	bedroom	bathroom	parking_space	price	bedbath_ratio	t
0	3583	Katsina	Semi-detached duplex	2.0	2.0	1.0	1149999.565	1.000000	
1	2748	Ondo	Apartment	3.5	2.0	4.0	1672416.689	0.571429	
3	2224	Anambra	Detached duplex	5.0	2.0	4.0	2410306.756	0.400000	
4	10300	Kogi	Terrace duplex	5.0	5.0	6.0	2600700.898	1.000000	
5	1733	Borno	Mansion	4.0	1.0	3.0	1341750.867	0.250000	

```
In [34]: # Calculate the frequency of each category in the 'loc' column
loc_frequencies = all_data['loc'].value_counts(normalize=True)

# Create a dictionary to map each category to its corresponding frequency
loc_frequency_mapping = loc_frequencies.to_dict()

# Map the 'loc' and column to its corresponding frequency values
all_data['loc'] = all_data['loc'].map(loc_frequency_mapping)

# Print the updated dataframe
all_data.head()
```

Out[34]:

	ID	loc	title	bedroom	bathroom	parking_space	bedbath_ratio	total_rooms	log_pr
0	3583	0.028250	6	2.0	2.0	1.0	1.000000	4.0	13.955
1	2748	0.028250	1	3.5	2.0	4.0	0.571429	5.5	14.329
3	2224	0.029641	8	5.0	2.0	4.0	0.400000	7.0	14.695
4	10300	0.027585	4	5.0	5.0	6.0	1.000000	10.0	14.771
5	1733	0.029883	10	4.0	1.0	3.0	0.250000	5.0	14.109

Modeling

```
In [35]: #importing necessary libraries
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from catboost import CatBoostRegressor
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.model_selection import KFold
from sklearn.metrics import mean_squared_error
from lightgbm import LGBMRegressor

import warnings
warnings.filterwarnings('ignore')
```

```
In [36]: # dropping columns not needed and setting the feature and label
not_needed = ['log_price']
# splitting all data into x, y and test_df
X= all_data[:final_data.shape[0]].drop(columns = not_needed, axis = 1)
y= final_data['price']
test_data= all_data[final_data.shape[0]:]

#checking the outcome
X.shape, y.shape, test_data.shape
```

```
Out[36]: ((10531, 14), (10531,), (6000, 15))
```

```
In [37]: # split data for training and testing with ratio 7:3
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, ra
```

CatBoost Regressor

```
In [38]: #creating the model function
catb_model= CatBoostRegressor(verbose=0)

#model fitting
catb_model.fit(X_train, y_train)

#prediction
catb_pred= catb_model.predict(X_test)

#checking the mean squared error
print(f'mse = {mean_squared_error(y_test, catb_pred, squared=False)}')
```

```
catb_pred
```

```
mse = 445773.6317913167
```

```
Out[38]: array([2077811.04427041, 3566937.78881224, 755214.83587719, ...,
1632577.26935423, 1474284.42791369, 3266851.06509875])
```

LGBMRegressor

```

In [39]: #creating the model function
lgb_model = LGBMRegressor()

#model fitting
lgb_model.fit(X_train, y_train)

#prediction
lgb_pred= lgb_model.predict(X_test)

#checking the mean squared error
print(f'mse = {mean_squared_error(y_test, lgb_pred, squared=False)}')

#printing the prediction
lgb_pred

[LightGBM] [Warning] Found whitespace in feature_names, replace with underlines
[LightGBM] [Warning] Auto-choosing row-wise multi-threading, the overhead of testing was 0.000506 seconds.
You can set `force_row_wise=true` to remove the overhead.
And if memory is not enough, you can set `force_col_wise=true`.
[LightGBM] [Info] Total Bins 430
[LightGBM] [Info] Number of data points in the train set: 7898, number of used features: 14
[LightGBM] [Info] Start training from score 2133402.757573
mse = 454624.8693543962

Out[39]: array([2135734.79100528, 3407041.45513572, 822074.4006056 , ...,
                1652112.74659283, 1549887.53296109, 3330870.81250884])

```

```

In [40]: params = {
    'n_estimators': 500,
    'colsample_bytree': 0.86,
    'learning_rate': 0.032,
    'max_depth': 7,
    'subsample': 0.85}

test_pred=[]
y_pred = []

fold = KFold(n_splits=8, shuffle=True)#15#5#10
i=1
for train_index, test_index in fold.split(X,y):

    X_train, X_test = X.iloc[train_index], X.iloc[test_index]
    y_train, y_test = np.log1p(y.iloc[train_index]), y.iloc[test_index]

    model = LGBMRegressor(**params, objective = "rmse")
    model.fit(X_train,y_train,eval_set=[(X_train,y_train),(X_test, y_test)])#e

    preds= model.predict(X_test)
    print("err: ",(mean_squared_error(y_test,np.expm1(preds), squared=False)))
    y_pred.append(mean_squared_error(y_test,np.expm1(preds),squared=False))
    t_pred = model.predict(test_data[X.columns])
    test_pred.append(np.expm1(t_pred))

print(np.mean(y_pred))

```

[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
 [LightGBM] [Warning] No further splits with positive gain, best gain: -inf
 [LightGBM] [Warning] No further splits with positive gain, best gain: -inf
 [LightGBM] [Warning] No further splits with positive gain, best gain: -inf
 [LightGBM] [Warning] No further splits with positive gain, best gain: -inf
 [LightGBM] [Warning] No further splits with positive gain, best gain: -inf
 [LightGBM] [Warning] No further splits with positive gain, best gain: -inf
 [LightGBM] [Warning] No further splits with positive gain, best gain: -inf
 [LightGBM] [Warning] No further splits with positive gain, best gain: -inf
 [LightGBM] [Warning] No further splits with positive gain, best gain: -inf
 [LightGBM] [Warning] No further splits with positive gain, best gain: -inf
 [LightGBM] [Warning] Accuracy may be bad since you didn't explicitly set num_leaves OR 2^max_depth > num_leaves. (num_leaves=31).
 err: 447723.8307158369
 [LightGBM] [Warning] Accuracy may be bad since you didn't explicitly set num_leaves OR 2^max_depth > num_leaves. (num_leaves=31).
 427809.17916420946


```
In [41]: submission.head()
```

Out[41]:

	ID
0	845
1	1924
2	10718
3	12076
4	12254

```
In [42]: submission['price'] = np.mean(test_pred, axis = 0)
```

```
In [43]: submission.head()
```

Out[43]:

	ID	price
0	845	2.319931e+06
1	1924	1.013887e+06
2	10718	1.231032e+06
3	12076	8.366577e+06
4	12254	1.915268e+06

```
In [44]: submission.to_csv('Submission.csv', index=False)
```

```
In [ ]:
```