

MACHINE LEARNING FINAL PROJECT PART A

NAME: Gazi Aman Khan BATCH: 15th May(DataScience)

EXPLANATION VIDEO LINK:https://www.youtube.com/watch?v=70330f13e3e43981a5c9d1d2e0192

PROBLEM STATEMENT The objective of this analysis is to predict the listing prices of Airbnb properties based on various factors such as location, property type, number of bedrooms, amenities, and other listing features. By building and evaluating regression models, the goal is to identify key factors influencing price variations and develop a model capable of accurately estimating the price of new Airbnb listings.

This will help hosts understand pricing patterns and enable data-driven pricing decisions, improving competitiveness and revenue potential.

LIBRARIES USED: pandas, numpy, scikit-learn, matplotlib, seaborn, XGBoost.

```
In [7]: !pip install xgboost

Requirement already satisfied: xgboost in c:\users\hpb\anaconda3\lib\site-packages (2.1.1)
Requirement already satisfied: numpy in c:\users\hpb\anaconda3\lib\site-packages (from xgboost) (2.1.3)
Requirement already satisfied: scipy in c:\users\hpb\anaconda3\lib\site-packages (from xgboost) (1.15.3)

AIRBNB PRICE PREDICTION AND INSIGHTS

In [11]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
from xgboost import XGBRegressor
from xgboost import plot_importance

In [12]: #MODEL 1- LINEAR REGRESSION

In [13]: #TASK-1 :Data Exploration and Preprocessing

In [11]: df=pd.read_csv("E:\Course datasets\Airbnb\data.csv")
df.head(5)

Out[11]:
```

	id	log_price	property_type	room_type	amenities	accommodates	bathrooms	bed_type	cancellation_policy	cleaning_fee	...	latitude	longitude	name	neighbourhood	number_of_reviews	review_scores_rating	thumbnail_url
0	6901257	5.010635	Apartment	Entire home/apt	["Wireless Internet","Air conditioning","Kitch...	3	1.0	Real Bed	strict	True	...	40.686524	-73.981617	Beautiful brooklyn 1-bedroom	Brooklyn Heights	2	100.0	https://a0.muscache.com/im/pictures/6d7dbd67-c...
1	6304928	5.129899	Apartment	Entire home/apt	["Wireless Internet","Air conditioning","Kitch...	7	1.0	Real Bed	strict	True	...	40.766115	-73.989040	Superb 3BR Apt Located Near Times Square	Half's Kitchen	6	93.0	https://a0.muscache.com/im/pictures/948a559e-4...
2	7919400	4.976734	Apartment	Entire home/apt	["TV","Cable TV","Wireless Internet","Air condit...	5	1.0	Real Bed	moderate	True	...	40.808110	-73.943756	The Garden Oasis	Harlem	10	92.0	https://a0.muscache.com/im/pictures/6d7dbd67-c...
3	13418779	6.620073	House	Entire home/apt	["TV","Cable TV","Internet","Wireless Internet","Ki...	4	1.0	Real Bed	flexible	True	...	37.772004	-122.431619	Beautiful Flat in the Heart of SFI	Lower Haight	0	NaN	https://a0.muscache.com/im/pictures/72208bde-9...
4	3808709	4.744932	Apartment	Entire home/apt	["TV","Internet","Wireless Internet","Air conditio...	2	1.0	Real Bed	moderate	True	...	38.925627	-77.034596	Great studio in mdtown DC	Columbia Heights	4	40.0	NaN

5 rows × 19 columns

```
In [12]: df.info()

Out[12]:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 74111 entries, 0 to 74110
Data columns (total 19 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   id           74111 non-null   int64
1   log_price    74111 non-null   float64
2   property_type 74111 non-null   object
3   room_type    74111 non-null   object
4   amenities     74111 non-null   list64
5   accommodates 73911 non-null   int64
6   bathrooms    73911 non-null   float64
7   bed_type     74111 non-null   object
8   cancellation_policy 74111 non-null object
9   cleaning_fee 74111 non-null   bool
10  city         74111 non-null   object
11  description  74111 non-null   object
12  first_review 58247 non-null   object
13  host_has_profile_pic 73923 non-null object
14  host_identity_verified 73923 non-null object
15  host_response_rate 58810 non-null object
16  host_since   73911 non-null   object
17  instant_bookable 74111 non-null object
18  last_review 58247 non-null   object
19  latitude     74111 non-null   float64
20  longitude    74111 non-null   float64
21  name         74111 non-null   object
22  neighbourhood 67239 non-null object
23  number_of_reviews 74111 non-null int64
24  review_scores_rating 57389 non-null float64
25  thumbnail_url 65810 non-null object
26  zipcode      73143 non-null   object
27  bedrooms     74020 non-null   float64
28  beds         73980 non-null   float64
dtypes: bool(1), float64(17), int64(3), object(18)
memory usage: 15.9+ MB

In [13]: df.isnull().sum()

Out[13]:
id              0
log_price        0
property_type    0
room_type        0
amenities        0
accommodates    200
bathrooms        0
bed_type         0
cancellation_policy 0
cleaning_fee     0
city             0
description      0
first_review    15864
host_has_profile_pic 188
host_identity_verified 188
host_response_rate 16299
host_since      188
instant_bookable 0
last_review     15817
latitude        0
longitude       0
name            0
neighbourhood   6872
number_of_reviews 0
review_scores_rating 16722
thumbnail_url   6216
zipcode         968
bedrooms        91
beds           131
dtype: int64

In [14]: df.describe()

Out[14]:
```

	id	log_price	accommodates	bathrooms	latitude	longitude	number_of_reviews	review_scores_rating	bedrooms	beds
count	7.411100e+04	74111.000000	74111.000000	73911.000000	74111.000000	74111.000000	74111.000000	57389.000000	74020.000000	73980.000000
mean	1.126662e+07	4.762069	3.155146	1.235263	38.465958	-92.397525	20.900568	94.067365	1.265793	1.710868
std	6.081735e+06	0.717394	2.153589	0.562044	3.080167	21.705322	37.828641	7.836556	0.852143	1.254142
min	1.440000e+02	0.000000	1.000000	0.000000	33.338905	-122.511500	0.000000	20.000000	0.000000	0.000000
25%	6.261964e+06	4.317488	2.000000	1.000000	34.127908	-118.342374	1.000000	92.000000	1.000000	1.000000
50%	1.225415e+07	4.709630	2.000000	1.000000	40.862138	-76.996965	6.000000	96.000000	1.000000	1.000000
75%	1.640225e+07	5.220356	4.000000	1.000000	40.746096	-73.954660	23.000000	100.000000	1.000000	2.000000
max	2.122090e+07	7.600402	16.000000	8.000000	42.390437	-70.985047	605.000000	100.000000	10.000000	16.000000

```
In [15]: #Outliers
df['price'] = np.exp(df['log_price'])

sns.boxplot(x=df['price'])
plt.title('Outliers in price')
plt.show()

Outliers in price

In [16]: #HANDLING MISSING VALUES
num_cols = df.select_dtypes(include=['int64','float64']).columns #numerical missing values to median
df[num_cols] = df[num_cols].fillna(df[num_cols].median())

cat_cols = df.select_dtypes(include=['object']).columns #categorical missing values to "unknown"
df[cat_cols] = df[cat_cols].fillna('Unknown')

df.isnull().sum()

Out[16]:
id              0
log_price        0
property_type    0
room_type        0
amenities        0
accommodates    200
bathrooms        0
bed_type         0
cancellation_policy 0
cleaning_fee     0
city             0
description      0
first_review    15864
host_has_profile_pic 188
host_identity_verified 188
host_response_rate 16299
host_since      188
instant_bookable 0
last_review     15817
latitude        0
longitude       0
name            0
neighbourhood   6872
number_of_reviews 0
review_scores_rating 16722
thumbnail_url   6216
zipcode         968
bedrooms        91
beds           131
price
dtype: int64

In [17]: df['host_since'] = pd.to_datetime(df['host_since'], errors='coerce')
df['last_review'] = pd.to_datetime(df['last_review'], errors='coerce')

C:\Users\hpb\AppData\Local\Temp\ipykernel_3912\2637232071.py:1: UserWarning: Parsing dates in %d-%m-%Y format when dayfirst=False (the default) was specified. Pass 'dayfirst=True' or specify a format to silence this warning.
df['host_since'] = pd.to_datetime(df['host_since'], errors='coerce')
C:\Users\hpb\AppData\Local\Temp\ipykernel_3912\2637232071.py:2: UserWarning: Parsing dates in %d-%m-%Y format when dayfirst=False (the default) was specified. Pass 'dayfirst=True' or specify a format to silence this warning.
df['last_review'] = pd.to_datetime(df['last_review'], errors='coerce')

In [18]: df['host_experience_days'] = (pd.Timestamp.today() - df['host_since']).dt.days #host experience indays
df['n_amenities'] = df['amenities'].apply(lambda x: len(str(x).split(','))) #count of amenities in each row

In [19]: #STANDARDIZING THE DATA
scaler=StandardScaler()
df[['accommodates','bedrooms','beds','n_amenities']] = scaler.fit_transform(
df[['accommodates','bedrooms','beds','n_amenities']])

In [20]: #TASK-2:Model Development.

In [21]: cat_cols = ['property_type', 'room_type', 'neighbourhood'] #SELECTING ONLY MEANINGFUL COLUMNS BECAUSE OF MEMORY ISSUES
num_cols = ['accommodates','bathrooms','bedrooms','beds','number_of_reviews',
'review_scores_rating','n_amenities','host_experience_days']

In [22]: x=df[cat_cols+num_cols]
y=df['log_price']

In [23]: #SPLIT THE DATA
X_train,X_test,y_train,y_test= train_test_split(x, y, test_size=0.2, random_state=42)

In [24]: #HANDLING CATEGORICAL COLUMNS
X_train = pd.get_dummies(X_train, columns=cat_cols, drop_first=True)
X_test = pd.get_dummies(X_test, columns=cat_cols, drop_first=True)
X_train=X_test.reindex(columns=X_train.columns, fill_value=0) #ensuring both datasets have same columns in order

In [25]: X_train = X_train.fillna(X_train.median()) #handling remaining NaN values in numeric columns
X_test = X_test.fillna(X_train.median())

LINEAR REGRESSION

In [26]: #TRAIN MODEL
model=LinearRegression()
model.fit(X_train,y_train)

Out[26]:
LinearRegression()

In [27]: #PREDICTIONS
y_pred=model.predict(X_test)

In [28]: #EVALUATING THE MODEL'S PERFORMANCE
mse=mean_squared_error(y_test, y_pred)
rmse=np.sqrt(mse)
mae=mean_absolute_error(y_test, y_pred)
r2=r2_score(y_test, y_pred)

print("RMSE:",rmse)
print("MAE:",mae)
print("R2 SCORE:",r2)

RMSE: 0.423663210941468
MAE: 0.307548217263002
R2 SCORE: 0.645402237110424

In [29]: y_pred_price = np.exp(y_pred)
X_test_price = np.exp(X_test)

pred_df = pd.DataFrame({
    'Actual Price': y_test_price,
    'Predicted Price': y_pred_price})

pred_df.head(10)

Out[29]:
```

	Actual Price	Predicted Price
4079	45.0	52.367957
33735	140.0	121.137535
69475	140.0	102.768853
454	130.0	268.934771
25153	39.0	72.356799
42304	126.0	144.191619
22686	185.0	199.441368
20776	296.0	170.844461
14934	35.0	47.790085
39634	69.0	199.058672

```
In [30]: #SCATTER PLOT OF ACTUAL VS PREDICTED PRICES
plt.figure(figsize=(8,6))
sns.scatterplot(x=y_test_price, y=y_pred_price)
plt.title('Actual VS Predicted Price')
plt.xlabel('Actual Price')
plt.ylabel('Predicted Price')
plt.plot((y_test_price.min(), y_test_price.max()),
        (y_test_price.min(), y_test_price.max()), 'r--')
plt.show()

Actual VS Predicted Price

MODEL-2: XGBOOST MODEL

In [32]: xgb_model = XGBRegressor(n_estimators=300, learning_rate=0.05, max_depth=6, random_state=42, n_jobs=-1)

In [33]: #train the data
xgb_model.fit(X_train, y_train)

Out[33]:
XGBRegressor(base_score=None, booster=None, callbacks=None,
             colsample_bylevel=None, colsample_bynode=None,
             colsample_hytree=None, device=None, early_stopping_rounds=None,
             enable_categorical=False, eval_metric=None, feature_types=None,
             feature_weights=None, gamma=None, grow_policy=None,
             importance_type=None, interaction_constraints=None,
             learning_rate=0.05, max_bin=None, max_cat_threshold=None,
             max_cat_to_onehot=None, max_delta_step=None, max_depth=6,
             max_leaves=None, min_child_weight=None, missing=nan,

In [34]: #PREDICTIONS
yqg_pred= xgb_model.predict(X_test)

In [35]: #EVALUATE THE DATA
RMSE=mean_squared_error(y_test, yqg_pred)
RMSE=np.sqrt(RMSE)
MAE=mean_absolute_error(y_test, yqg_pred)
R2=r2_score(y_test, yqg_pred)

print("RMSE:", RMSE)
print("MAE:", MAE)
print("R2 score:", R2)

RMSE: 0.428640944848881
MAE: 0.316514580956707
R2 score: 0.6423463699193158

In [36]: yqg_pred_price = np.exp(yqg_pred)
X_test_price = np.exp(X_test)

preds_df = pd.DataFrame({
    'Actual Price': y_test_price,
    'Predicted price': yqg_pred_price})

preds_df.head(10)

Out[36]:
```

	Actual price	Predicted price
4079	45.0	61.249636
33735	140.0	134.211304
69475	140.0	140.262390
454	130.0	213.955887
25153	39.0	66.718208
42304	126.0	154.585922
22686	185.0	205.670850
20776	296.0	130.928757
14934	35.0	62.748131
39634	69.0	94.732292

```
In [37]: #SCATTER PLOT OF ACTUAL VS PREDICTED PRICES
plt.figure(figsize=(8,6))
sns.scatterplot(x=y_test_price, y=yqg_pred_price)
plt.title('Actual VS Predicted Price XGBOOST Model')
plt.xlabel('Actual Price')
plt.ylabel('Predicted Price')
plt.plot((y_test_price.min(), y_test_price.max()),
        (y_test_price.min(), y_test_price.max()), 'r--')
plt.show()

Actual VS Predicted Price XGBOOST Model

In [38]: #METRIC COMPARISON
result=
{
    'Model':('Linear Regression', 'XGBoost'),
    'RMSE':(rmse,RMSE),
    'MAE':(mae,MAE),
    'R2 Score':(r2,R2)}

result_df=pd.DataFrame(result)
print(result_df)

Model      RMSE      MAE      R2 Score
0  Linear Regression  0.423663  0.307548  0.650610
1      XGBoost       0.428645  0.316515  0.642344

In [42]: #IMPORTANCE
plot_importance(xgb_model, max_num_features=10)

Out[42]:
<Axes: title='Feature importance', xlabel='Importance score', ylabel='Features'>

Feature importance

host_experience_days 1894.0
n_amenities          1799.0
number_of_reviews    1593.0
accommodates         1532.0
bathrooms            916.0
review_scores_rating 865.0
bedrooms             696.0
beds                 344.0
beds                 344.0
room_type_Private room 264.0
property_type_House   235.0
Importance score
```