

```
import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
%matplotlib inline
import matplotlib
matplotlib.rcParams["figure.figsize"]=(20,10)
```

```
df=pd.read_csv("bengaluru_house_prices.csv")
df.head()
```

	area_type	availability	location	size	society	total_sqft	bath	balcony	price
0	Super built-up Area	19-Dec	Electronic City Phase II	2 BHK	Coomee	1056	2.0	1.0	39.07
1	Plot Area	Ready To Move	Chikka Tirupathi	4 Bedroom	Theanmp	2600	5.0	3.0	120.00
2	Built-up Area	Ready To Move	Uttarahalli	3 BHK	NaN	1440	2.0	3.0	62.00
3	Super built-up Area	Ready To Move	Lingadheeranahalli	3 BHK	Soiewre	1521	3.0	1.0	95.00
4	Super built-up Area	Ready To Move	Kothanur	2 BHK	NaN	1200	2.0	1.0	51.00

```
df.shape
```

(13320, 9)

```
df.groupby('area_type')['area_type'].agg('count')
```

```
area_type
Built-up Area      2418
Carpet Area         87
Plot Area          2025
Super built-up Area 8790
Name: area_type, dtype: int64
```

```
df1=df.drop(['area_type','society','balcony','availability'],axis='columns')
df1.head()
```

	location	size	total_sqft	bath	price
0	Electronic City Phase II	2 BHK	1056	2.0	39.07
1	Chikka Tirupathi	4 Bedroom	2600	5.0	120.00
2	Uttarahalli	3 BHK	1440	2.0	62.00
3	Lingadheeranahalli	3 BHK	1521	3.0	95.00
4	Kothanur	2 BHK	1200	2.0	51.00

```
df1.isnull().sum()
```

```
location      1
size          16
total_sqft     0
bath          73
price         0
dtype: int64
```

```
df2=df1.dropna()
df2.isnull().sum()
```

```
location      0
size          0
total_sqft     0
bath          0
price         0
dtype: int64
```

```
df2['size'].unique()
```

```
array(['2 BHK', '4 Bedroom', '3 BHK', '4 BHK', '6 Bedroom', '3 Bedroom',
      '1 BHK', '1 RK', '1 Bedroom', '8 Bedroom', '2 Bedroom',
      '7 Bedroom', '5 BHK', '7 BHK', '6 BHK', '5 Bedroom', '11 BHK',
      '9 BHK', '9 Bedroom', '27 BHK', '10 Bedroom', '11 Bedroom',
      '10 BHK', '19 BHK', '16 BHK', '43 Bedroom', '14 BHK', '8 BHK',
      '12 Bedroom', '13 BHK', '18 Bedroom'], dtype=object)
```

```
df2['bhk']=df2['size'].apply(lambda x: int(x.split(' ')[0]))
df2.head()
```

```
<ipython-input-16-4b3156990b73>:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-df2\['bhk'\]=df2\['size'\].apply\(lambda x: int\(x.split\(' '\)\[0\]\)\)](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-df2['bhk']=df2['size'].apply(lambda x: int(x.split(' ')[0])))

	location	size	total_sqft	bath	price	bhk
0	Electronic City Phase II	2 BHK	1056	2.0	39.07	2
1	Chikka Tirupathi	4 Bedroom	2600	5.0	120.00	4

```
df2['bhk'].unique()
```

```
array([ 2,  4,  3,  6,  1,  8,  7,  5, 11,  9, 27, 10, 19, 16, 43, 14, 12,
        13, 18])
```

```
df2[df2.bhk>20]
```

	location	size	total_sqft	bath	price	bhk
1718	2Electronic City Phase II	27 BHK	8000	27.0	230.0	27
4684	Munnekollal	43 Bedroom	2400	40.0	660.0	43

```
def is_float(x):
```

```
    try:
        float(x)
    except:
        return False
    return True
```

```
df2[~df2['total_sqft'].apply(is_float)].head()
```

	location	size	total_sqft	bath	price	bhk
30	Yelahanka	4 BHK	2100 - 2850	4.0	186.000	4
122	Hebbal	4 BHK	3067 - 8156	4.0	477.000	4
137	8th Phase JP Nagar	2 BHK	1042 - 1105	2.0	54.005	2
165	Sarjapur	2 BHK	1145 - 1340	2.0	43.490	2
188	KR Puram	2 BHK	1015 - 1540	2.0	56.800	2

```
def convert_sqft_to_num(x):
```

```
    tokens=x.split('-')
    if len(tokens)==2:
        return (float(tokens[0])+float(tokens[1]))/2
    try :
        return float(x)
    except:
        return None
```

```
df3=df2.copy()
```

```
df3['total_sqft']=df3['total_sqft'].apply(convert_sqft_to_num)
df3.head()
```

	location	size	total_sqft	bath	price	bhk
0	Electronic City Phase II	2 BHK	1056.0	2.0	39.07	2
1	Chikka Tirupathi	4 Bedroom	2600.0	5.0	120.00	4
2	Uttarahalli	3 BHK	1440.0	2.0	62.00	3
3	Lingadheeranahalli	3 BHK	1521.0	3.0	95.00	3
4	Kothanur	2 BHK	1200.0	2.0	51.00	2

```
df3.loc[30]
```

```
location    Yelahanka
size        4 BHK
total_sqft  2475.0
bath        4.0
price       186.0
bhk         4
Name: 30, dtype: object
```

```
df3.head()
```

	location	size	total_sqft	bath	price	bhk
0	Electronic City Phase II	2 BHK	1056.0	2.0	39.07	2
1	Chikka Tirupathi	4 Bedroom	2600.0	5.0	120.00	4
2	Uttarahalli	3 BHK	1440.0	2.0	62.00	3

```
df4=df3.copy()
df4['price_per_sqft']=df4['price']*100000/df4['total_sqft']
df4.head()
```

	location	size	total_sqft	bath	price	bhk	price_per_sqft
0	Electronic City Phase II	2 BHK	1056.0	2.0	39.07	2	3699.810606
1	Chikka Tirupathi	4 Bedroom	2600.0	5.0	120.00	4	4615.384615
2	Uttarahalli	3 BHK	1440.0	2.0	62.00	3	4305.555556
3	Lingadheeranahalli	3 BHK	1521.0	3.0	95.00	3	6245.890861
4	Kothanur	2 BHK	1200.0	2.0	51.00	2	4250.000000

```
len(df4['location'].unique())
```

```
1304
```

```
df4.location=df4.location.apply(lambda x: x.strip())
location_stats=df4.groupby('location')['location'].agg('count').sort_values(ascending=False)
location_stats
```

```
location
Whitefield      535
Sarjapur Road   392
Electronic City 304
Kanakpura Road  266
Thanisandra     236
...
1 Giri Nagar    1
Kanakapura Road, 1
Kanakapura main Road 1
Karnataka Shabarimala 1
whitefiled      1
Name: location, Length: 1293, dtype: int64
```

```
len(location_stats[location_stats<=10])
```

```
1052
```

```
location_stats_less_than_10=location_stats[location_stats<=10]
location_stats_less_than_10
```

```
location
Basapura      10
1st Block Koramangala 10
Gunjur Palya  10
Kalkere       10
Sector 1 HSR Layout 10
..
1 Giri Nagar    1
Kanakapura Road, 1
Kanakapura main Road 1
Karnataka Shabarimala 1
whitefiled      1
Name: location, Length: 1052, dtype: int64
```

```
df4.location=df4.location.apply(lambda x: 'other' if x in location_stats_less_than_10 else x)
len(df4.location.unique())
df4.head()
```

	location	size	total_sqft	bath	price	bhk	price_per_sqft
0	Electronic City Phase II	2 BHK	1056.0	2.0	39.07	2	3699.810606
1	Chikka Tirupathi	4 Bedroom	2600.0	5.0	120.00	4	4615.384615
2	Uttarahalli	3 BHK	1440.0	2.0	62.00	3	4305.555556
3	Lingadheeranahalli	3 BHK	1521.0	3.0	95.00	3	6245.890861
4	Kothanur	2 BHK	1200.0	2.0	51.00	2	4250.000000

```
df4[df4.total_sqft/df4.bhk<300].head()
```

	location	size	total_sqft	bath	price	bhk	price_per_sqft
9	other	6 Bedroom	1020.0	6.0	370.0	6	36274.509804
45	HSR Layout	8 Bedroom	600.0	9.0	200.0	8	33333.333333

```
df5=df4[~(df4.total_sqft/df4.bhk<300)]
df5.shape
```

```
(12502, 7)
```

```
df5.price_per_sqft.describe()
```

```
count    12456.000000
mean      6308.502826
std       4168.127339
min       267.829813
25%      4210.526316
50%      5294.117647
75%      6916.666667
max      176470.588235
Name: price_per_sqft, dtype: float64
```

```
def remove_pps_outliers(df):
    df_out=pd.DataFrame()
    for key ,subdf in df.groupby('location'):
        m=np.mean(subdf.price_per_sqft)
        st=np.std(subdf.price_per_sqft)
        reduced_df=subdf[(subdf.price_per_sqft>(m-st)) & (subdf.price_per_sqft<=(m+st))]
        df_out=pd.concat([df_out,reduced_df],ignore_index=True)
    return df_out
```

```
df6=remove_pps_outliers(df5)
df6.shape
```

```
(10241, 7)
```

```
def plot_scatter_chart(df,location):
    bhk2=df[(df.location==location) & (df.bhk==2)]
    bhk3=df[(df.location==location) & (df.bhk==3)]
    matplotlib.rcParams['figure.figsize']=(15,10)
    plt.scatter(bhk2.total_sqft,bhk2.price_per_sqft,color='blue',label='2 bhk ',s=50)
    plt.scatter(bhk3.total_sqft,bhk3.price_per_sqft,color='green',marker='+',label='3 bhk ',s=50)
    plt.xlabel("total square feet area")
    plt.ylabel("price per square feet area")
    plt.title(location)
    plt.legend()
plot_scatter_chart(df6,"Hebbal")
```

```
def remove_bhk_outliers(df):
    exclude_indices=np.array([])
    for location,location_df in df.groupby('location'):
```

```

bhk_stats={}
for bhk,bhk_df in location_df.groupby('bhk'):
    bhk_stats[bhk]={
        'mean':np.mean(bhk_df.price_per_sqft),
        'std':np.std(bhk_df.price_per_sqft),
        'count':bhk_df.shape[0]
    }
for bhk,bhk_df in location_df.groupby('bhk'):
    stats=bhk_stats.get(bhk-1)
    if stats and stats['count']>5:
        exclude_indices=np.append(exclude_indices,bhk_df[bhk_df.price_per_sqft<(stats['mean'])].index.values)
return df.drop(exclude_indices,axis='index')

```

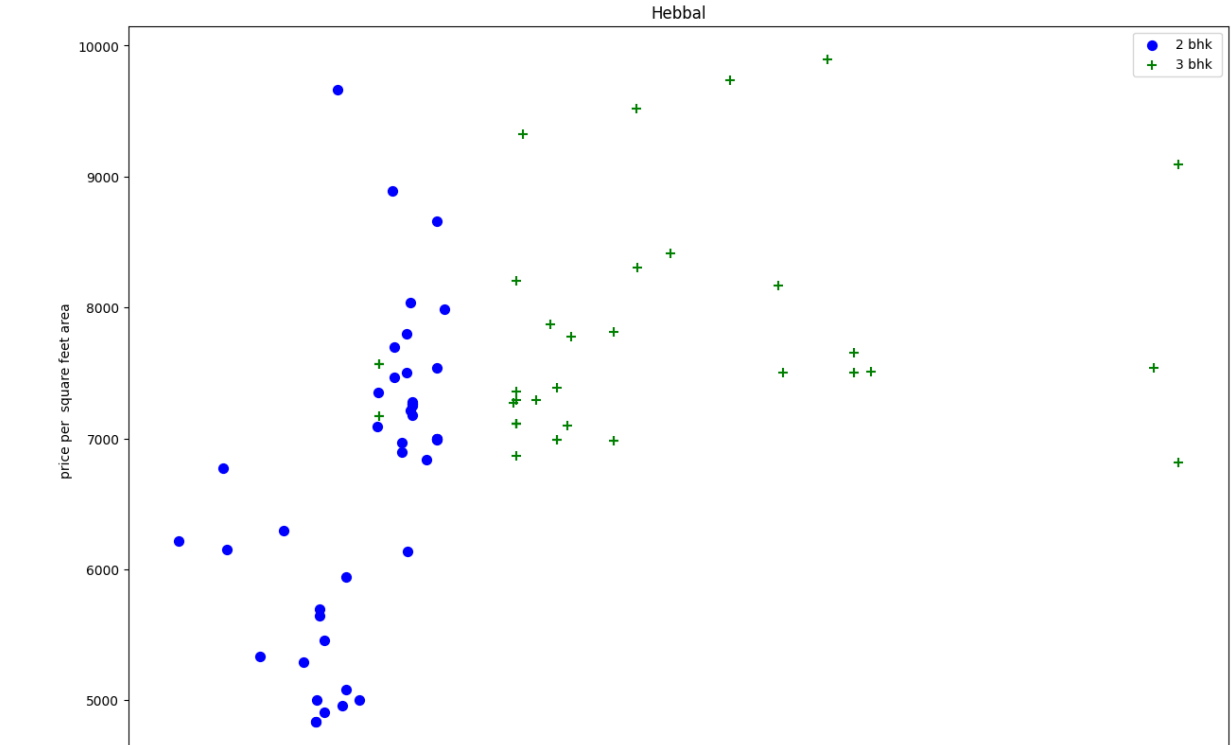
```

df7=remove_bhk_outliers(df6)
df7.shape

(7329, 7)

plot_scatter_chart(df7,"Hebbal")

```



```

import matplotlib
matplotlib.rcParams["figure.figsize"]=(20,10)
plt.hist(df7.price_per_sqft,rwidth=0.8)
plt.xlabel("price per square feet")
plt.ylabel("count")

```

```

df7.bath.unique()

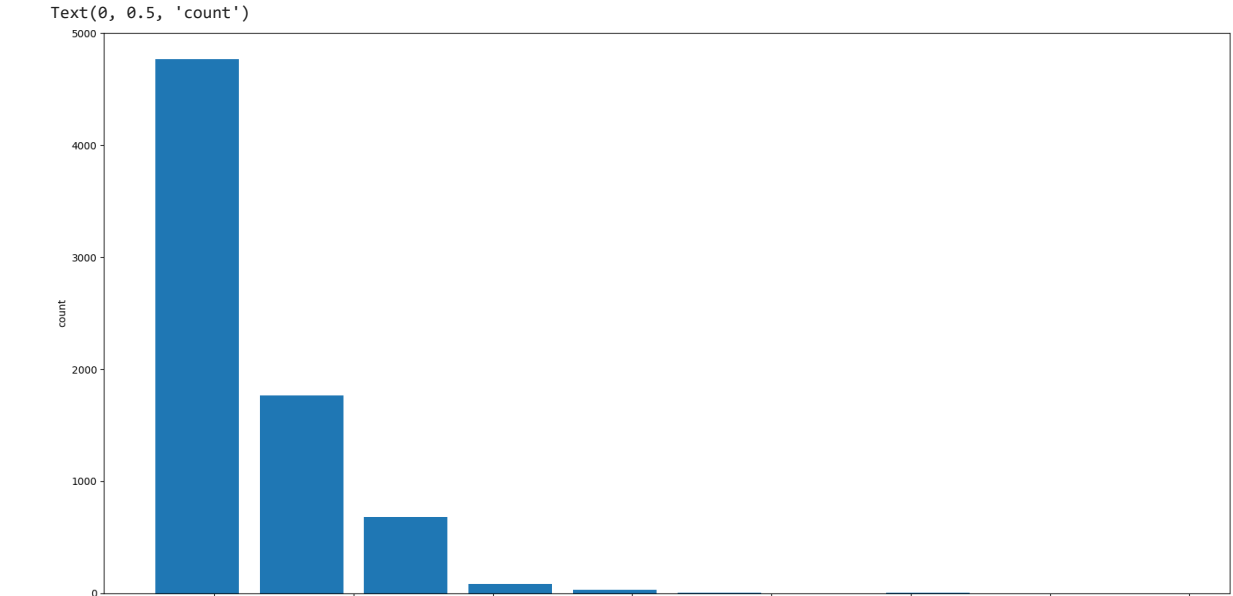
```

```
array([ 4.,  3.,  2.,  5.,  8.,  1.,  6.,  7.,  9., 12., 16., 13.])
|
df7[df.bath>10]
```

```
<ipython-input-42-d2192480d2f0>:1: UserWarning: Boolean Series key will be reindexed to match DataFrame index.
df7[df.bath>10]
```

	location	size	total_sqft	bath	price	bhk	price_per_sqft
938	Bannerghatta Road	2 BHK	1200.0	2.0	78.00	2	6500.000000
1078	Begur Road	2 BHK	1200.0	2.0	44.00	2	3666.666667
1718	Dasanapura	3 BHK	1286.0	2.0	68.00	3	5287.713841
1768	Devarachikkanahalli	2 BHK	947.0	2.0	43.00	2	4540.654699
1953	Electronic City	3 BHK	1563.0	3.0	91.84	3	5875.879718
1979	Electronic City	2 BHK	1128.0	2.0	65.50	2	5806.737589
3096	Hennur Road	3 BHK	2264.0	3.0	168.00	3	7420.494700
4684	Kudlu Gate	3 BHK	1535.0	3.0	85.00	3	5537.459283
6937	Uttarahalli	2 BHK	1025.0	2.0	35.88	2	3500.487805
8106	other	3 BHK	1976.0	3.0	184.00	3	9311.740891
8636	other	2 BHK	900.0	2.0	70.00	2	7777.777778

```
plt.hist(df7.bath,rwidth=0.8)
plt.xlabel("Num of bathrooms")
plt.ylabel("count")
```



```
df7[df7.bath>df7.bhk+2]
```

	location	size	total_sqft	bath	price	bhk	price_per_sqft
1626	Chikkabanavar	4 Bedroom	2460.0	7.0	80.0	4	3252.032520
5238	Nagasandra	4 Bedroom	7000.0	8.0	450.0	4	6428.571429
6711	Thanisandra	3 BHK	1806.0	6.0	116.0	3	6423.034330
8411	other	6 BHK	11338.0	9.0	1000.0	6	8819.897689

```
df8=df7[df7.bath<df7.bhk+2]
df8.shape
```

(7251, 7)

```
df9=df8.drop(['size','price_per_sqft'],axis='columns')
df9.head()
```

```
location total sqft bath price bhk
dummies=pd.get_dummies(df9.location)
dummies.head()
```

	1st Block Jayanagar	1st Phase JP Nagar	2nd Phase Judicial Layout	2nd Stage Nagarbhavi	5th Block Hbr Layout	5th Phase JP Nagar	6th Phase JP Nagar	7th Phase JP Nagar	8th Phase JP Nagar	9th Phase JP Nagar	...	Vishveshwarya Layout	Vishwapriya Layout	Vittasa
0	1	0	0	0	0	0	0	0	0	0	...	0	0	
1	1	0	0	0	0	0	0	0	0	0	...	0	0	
2	1	0	0	0	0	0	0	0	0	0	...	0	0	
3	1	0	0	0	0	0	0	0	0	0	...	0	0	
4	1	0	0	0	0	0	0	0	0	0	...	0	0	

5 rows × 242 columns

```
df10=pd.concat([df9,dummies.drop('other',axis='columns')],axis='columns')
df10.head()
```

	location	total_sqft	bath	price	bhk	1st Block Jayanagar	1st Phase JP Nagar	2nd Phase Judicial Layout	2nd Stage Nagarbhavi	5th Block Hbr Layout	...	Vijayanagar	Vishveshwarya Layout	Vi
0	1st Block Jayanagar	2850.0	4.0	428.0	4	1	0	0	0	0	...	0	0	
1	1st Block Jayanagar	1630.0	3.0	194.0	3	1	0	0	0	0	...	0	0	
2	1st Block Jayanagar	1875.0	2.0	235.0	3	1	0	0	0	0	...	0	0	
3	1st Block Jayanagar	1200.0	2.0	130.0	3	1	0	0	0	0	...	0	0	
4	1st Block Jayanagar	1235.0	2.0	148.0	2	1	0	0	0	0	...	0	0	

5 rows × 246 columns

```
df11=df10.drop('location',axis='columns')
df11.head()
```

	total_sqft	bath	price	bhk	1st Block Jayanagar	1st Phase JP Nagar	2nd Phase Judicial Layout	2nd Stage Nagarbhavi	5th Block Hbr Layout	5th Phase JP Nagar	...	Vijayanagar	Vishveshwarya Layout	Vishw
0	2850.0	4.0	428.0	4	1	0	0	0	0	0	...	0	0	
1	1630.0	3.0	194.0	3	1	0	0	0	0	0	...	0	0	
2	1875.0	2.0	235.0	3	1	0	0	0	0	0	...	0	0	
3	1200.0	2.0	130.0	3	1	0	0	0	0	0	...	0	0	
4	1235.0	2.0	148.0	2	1	0	0	0	0	0	...	0	0	

5 rows × 245 columns

```
df11.shape

(7251, 245)

X=df11.drop('price',axis='columns')
X.head()
```

	1st Block Phase	1st Phase	2nd Phase	2nd Stage	5th Block	5th Phase	6th Phase	Highwayway	Highway
y=df11.price									
y.head()									
0	428.0								
1	194.0								
2	235.0								
3	130.0								
4	148.0								
Name: price, dtype: float64									

```

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.2,random_state=10)

from sklearn.linear_model import LinearRegression
lr_clf = LinearRegression()
lr_clf.fit(X_train,y_train)
lr_clf.score(X_test,y_test)

0.8452277697874376

from sklearn.model_selection import ShuffleSplit
from sklearn.model_selection import cross_val_score

cv = ShuffleSplit(n_splits=5, test_size=0.2, random_state=0)

cross_val_score(LinearRegression(), X, y, cv=cv)

array([0.82430186, 0.77166234, 0.85089567, 0.80837764, 0.83653286])

from sklearn.model_selection import GridSearchCV

from sklearn.linear_model import Lasso
from sklearn.tree import DecisionTreeRegressor

def find_best_model_using_gridsearchcv(X,y):
    algos = {
        'linear_regression': {
            'model': LinearRegression(),
            'params': {

            }
        },
        'lasso': {
            'model': Lasso(),
            'params': {
                'alpha': [1,2],
                'selection': ['random', 'cyclic']
            }
        },
        'decision_tree': {
            'model': DecisionTreeRegressor(),
            'params': {
                'criterion' : ['mse','friedman_mse'],
                'splitter': ['best','random']
            }
        }
    }
    scores = []
    cv = ShuffleSplit(n_splits=5, test_size=0.2, random_state=0)
    for algo_name, config in algos.items():
        gs = GridSearchCV(config['model'], config['params'], cv=cv, return_train_score=False)
        gs.fit(X,y)
        scores.append({
            'model': algo_name,
            'best_score': gs.best_score_,
            'best_params': gs.best_params_
        })

    return pd.DataFrame(scores,columns=['model', 'best_score', 'best_params'])
find_best_model_using_gridsearchcv(X,y)

```



```
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_validation.py:378: FitFailedWarning:
10 fits failed out of a total of 20.
The score on these train-test partitions for these parameters will be set to nan.
If these failures are not expected, you can try to debug them by setting error_score='raise'.
```

Below are more details about the failures:

-----  
10 fits failed with the following error:

Traceback (most recent call last):

```
File "/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_validation.py", line 686, in _fit_and_score
    estimator.fit(X_train, y_train, **fit_params)
File "/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py", line 1247, in fit
    super().fit(
File "/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py", line 177, in fit
    self._validate_params()
File "/usr/local/lib/python3.10/dist-packages/sklearn/base.py", line 600, in _validate_params
    validate_parameter_constraints(
File "/usr/local/lib/python3.10/dist-packages/sklearn/utils/_param_validation.py", line 97, in validate_parameter_constraints
    raise InvalidParameterError(
sklearn.utils._param_validation.InvalidParameterError: The 'criterion' parameter of DecisionTreeRegressor must be a str among
```

X.columns

```
Index(['total_sqft', 'bath', 'bhk', '1st Block Jayanagar',
      '1st Phase JP Nagar', '2nd Phase Judicial Layout',
      '2nd Stage Nagarbhavi', '5th Block Hbr Layout', '5th Phase JP Nagar',
      '6th Phase JP Nagar',
      ...,
      'Vijayanagar', 'Vishveshwarya Layout', 'Vishwapriya Layout',
      'Vittasandra', 'Whitefield', 'Yelachenahalli', 'Yelahanka',
      'Yelahanka New Town', 'Yelenahalli', 'Yeshwanthpur'],
      dtype='object', length=244)
```

```
def predict_price(location,sqft,bath,bhk):
    loc_index = np.where(X.columns==location)[0][0]
```

```
    x = np.zeros(len(X.columns))
    x[0] = sqft
    x[1] = bath
    x[2] = bhk
    if loc_index >= 0:
        x[loc_index] = 1

    return lr_clf.predict([x])[0]
```

```
predict_price('1st Phase JP Nagar',1000,3,4)
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegression was fitted
warnings.warn(
85.03231618809112
```



```
predict_price('Indira Nagar',1000, 2, 2)
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegression was fitted
warnings.warn(
181.2781548400639
```



```
predict_price('Indira Nagar',1000, 3, 3)
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LinearRegression was fitted
warnings.warn(
184.5843020203317
```



```
import pickle
with open ('bangalore_home_prices_model.pickle','wb') as f:
    pickle.dump(lr_clf,f)
```

```
import json
columns={
    'data_columns':[col.lower() for col in X.columns]
}
with open ("columns.json","w") as f:
    f.write(json.dumps(columns))
```

