

# REGRESSION MODEL FOR HOUSE PRICE PREDICTION

Dataset was obtain from Kaggle platform where the price of a house was predicted based on the following variables:  
number of bedrooms,number of bathrooms sqft\_living size, sqft\_lot size, number of floors, number of waterfront,view,house condition and grade amongst others.  
R\_squared score was used to measure model accuracy.

Linear Regression model R2\_squared score was (77%),Extra Tree Regressor model (89%),CatBoost Regressor(91%) and Gradient Boosting Regressor(88%).

NeuN CatBoost Regressor model was the best in predicting house prices for this dataset.

## STEPS

- 1>Loading the dataset
- 2.Exploratory data analysis
- 3.Data pre\_processing
- 4.Building and Training the models
- 5.Fitting the models
- 7.Model Evaluation

```
In [1]: #Importing relevant packages
import pandas as pd
import seaborn as sns
import numpy as np
import matplotlib.pyplot as plt
sns.set()
import matplotlib
import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: #Loading dataset
data=pd.read_csv('C:/Users/User/Desktop/DATA/house_data.csv')
```

```
In [3]: #read of data
data.head()
```

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	__	grade	sqft_above	sqft_basement	yr_built	yr_renovated	zipcode
0	7129300520	20141013T000000	219000.0	3	1.00	1180	5650	1.0	0	0	—	7	1180	0	1955	0	98178 47.
1	6414100102	20141209T000000	138000.0	3	2.25	2570	7242	2.0	0	0	—	7	2170	400	1951	1991	98125 47.
2	5603150400	20141020T000000	180000.0	2	1.00	770	10000	1.0	0	0	—	6	770	0	1963	0	98028 47.
3	2487200875	20141207T000000	640000.0	4	3.00	1960	5000	1.0	0	0	—	7	1050	910	1955	0	98136 47.
4	1954400151	20150218T000000	510000.0	3	2.00	1680	8080	1.0	0	0	—	8	1680	0	1987	0	98074 47.

5 rows x 17 columns

```
In [4]: #shape of data
data.shape
```

```
Out[4]: (21613, 17)
```

```
In [5]: #columns in data
data.columns
```

```
Out[5]: Index(['id', 'date', 'price', 'bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floors', 'waterfront', 'view', 'condition', 'grade', 'sqft_above', 'sqft_basement', 'yr_built', 'yr_renovated', 'zipcode', 'lat', 'long', 'sqft_living15', 'sqft_lot15'],
      dtype='object')
```

```
In [6]: #descriptive statistics
data.describe()
```

	count	mean	std	min	25%	50%	75%	max
id	21613.0	4.580302e+08	2.876566e+09	1.000102e+06	2.123049e+09	3.904800e+09	7.309800e+09	9.900000e+09
price	21613.0	5.401822e+05	3.673622e+05	7.500000e+04	3.218500e+05	4.500000e+05	6.450000e+05	7.700000e+06
bedrooms	21613.0	3.370842e+00	0.930613e+01	0.000000e+00	3.000000e+00	3.000000e+00	4.000000e+00	3.300000e+01
bathrooms	21613.0	2.114757e+00	0.701832e+01	0.000000e+00	1.750000e+00	2.250000e+00	2.500000e+00	8.000000e+00
sqft_living	21613.0	2.079900e+03	9.184409e+02	2.900000e+02	1.427000e+03	1.910000e+03	2.500000e+03	1.354000e+04
sqft_lot	21613.0	1.510067e+04	4.142051e+04	5.200000e+02	5.040000e+03	7.618000e+03	1.068800e+04	1.631359e+06
floors	21613.0	1.494309e+00	0.599889e+01	1.000000e+00	1.000000e+00	1.500000e+00	2.000000e+00	3.500000e+00
waterfront	21613.0	7.541757e-03	8.651729e-02	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	1.000000e+00
view	21613.0	2.343034e-01	7.663176e-01	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	1.000000e+00
condition	21613.0	3.409430e+00	6.507439e-01	1.000000e+00	3.000000e+00	4.000000e+00	4.000000e+00	5.000000e+00
grade	21613.0	7.656873e-00	1.175459e+00	1.000000e+00	7.000000e+00	7.000000e+00	8.000000e+00	1.300000e+01
sqft_above	21613.0	1.788391e+03	8.280910e+02	2.900000e+02	1.190000e+03	1.560000e+03	2.210000e+03	9.410000e+03
sqft_basement	21613.0	2.915050e+02	4.425750e+02	0.000000e+00	0.000000e+00	0.000000e+02	4.820000e+02	4.820000e+03
yr_built	21613.0	1.971005e+01	2.937341e+01	1.900000e+03	1.951000e+03	1.975000e+03	1.997000e+03	2.015000e+03
yr_renovated	21613.0	8.440226e-01	4.016792e+02	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	1.631359e+04
zipcode	21613.0	9.807794e+04	5.350503e+01	9.800100e+04	9.803300e+04	9.806500e+04	9.818000e+04	9.819900e+04
lat	21613.0	4.750005e+01	1.385631e-01	4.715590e+01	4.747100e+01	4.757100e+01	4.767800e+01	4.777600e+01
long	21613.0	-1.222139e+02	1.408283e-01	-1.225190e+02	-1.223280e+02	-1.223250e+02	-1.221250e+02	-1.213150e+02
sqft_living15	21613.0	1.986652e+03	6.653813e+02	3.990000e+02	1.490000e+03	1.840000e+03	2.360000e+03	6.210000e+03
sqft_lot15	21613.0	1.276846e+04	2.730418e+04	6.510000e+02	5.100000e+03	7.620000e+03	1.008300e+04	8.712000e+05

```
In [7]: #Information on data
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21613 entries, 0 to 21612
Data columns (total 17 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   id          21613 non-null   int64
 1   date        21613 non-null   object
 2   price       21613 non-null   float64
 3   bedrooms    21613 non-null   int64
 4   bathrooms   21613 non-null   float64
 5   sqft_living 21613 non-null   float64
 6   sqft_lot    21613 non-null   int64
 7   floors      21613 non-null   float64
 8   waterfront  21613 non-null   int64
 9   view        21613 non-null   int64
10   condition   21613 non-null   int64
11   grade       21613 non-null   int64
12   sqft_above  21613 non-null   int64
13   sqft_basement 21613 non-null   int64
14   yr_built    21613 non-null   int64
15   yr_renovated 21613 non-null   int64
16   zipcode     21613 non-null   int64
17   lat         21613 non-null   float64
18   long        21613 non-null   float64
19   sqft_living15 21613 non-null   int64
20   sqft_lot15  21613 non-null   int64
dtypes: float64(5), int64(15), object(1)
memory usage: 3.5+ MB
```

```
In [8]: #Convert 'zipcode' to string type
data['zipcode'] = data['zipcode'].astype(str)
```

```
In [9]: #types on data
data.dtypes
```

```
Out[9]: id          int64
       date        object
       price       float64
       bedrooms    int64
       bathrooms   float64
       sqft_living  float64
       sqft_lot     int64
       floors      float64
       waterfront  int64
       view        int64
       condition   int64
       grade       int64
       sqft_above  int64
       sqft_basement int64
       yr_built    int64
       yr_renovated int64
       zipcode     object
       lat         float64
       long        float64
       sqft_living15 int64
       sqft_lot15  int64
       dtype: object
```

```
In [10]: #checking null values
data.isnull().sum().sum()
```

```
Out[10]: 0
```

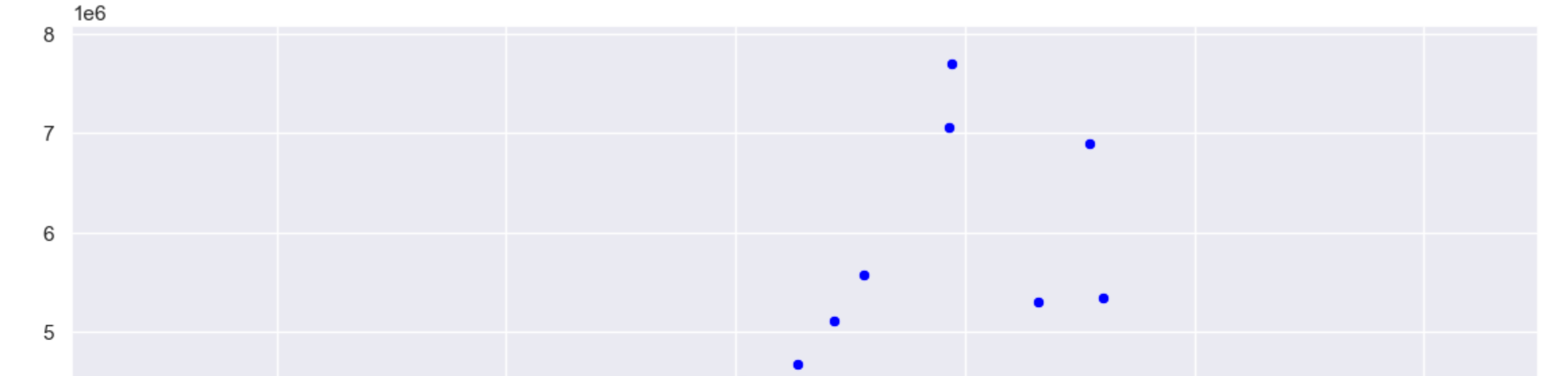
```
In [11]: #data duplicated().sum()
data.duplicated().sum()
```

```
Out[11]: 0
```

## EXPLORATORY DATA ANALYSIS

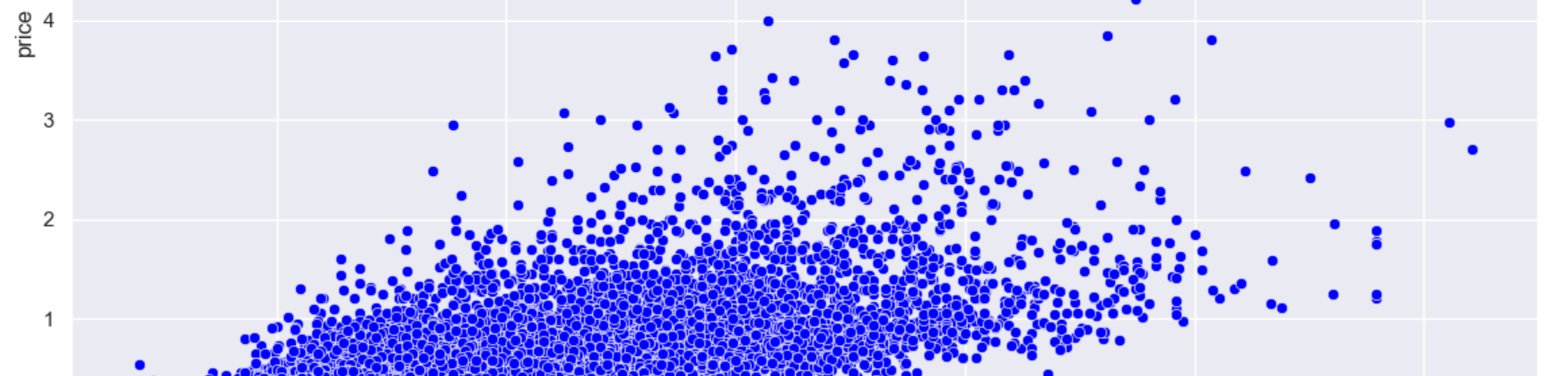
```
In [12]: plt.figure(figsize=(14,8))
sns.scatterplot(x='sqft_living15',y='price',data=data,color='blue')
```

```
Out[12]: <Axes: xlabel='sqft_living15', ylabel='price'>
```



```
In [13]: sns.histplot(x='price',data=data,color='blue')
```

```
Out[13]: <Axes: xlabel='price', ylabel='Count'>
```

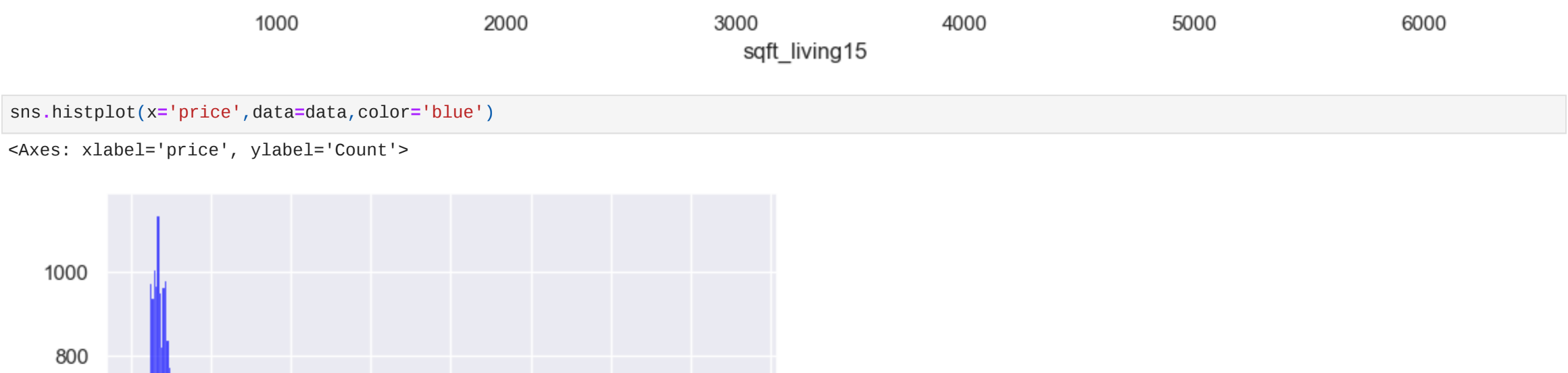


```
In [14]: #transforming price using np.log
log_price=np.log(data['price'])
```

```
In [15]: data['log_price']=log_price
```

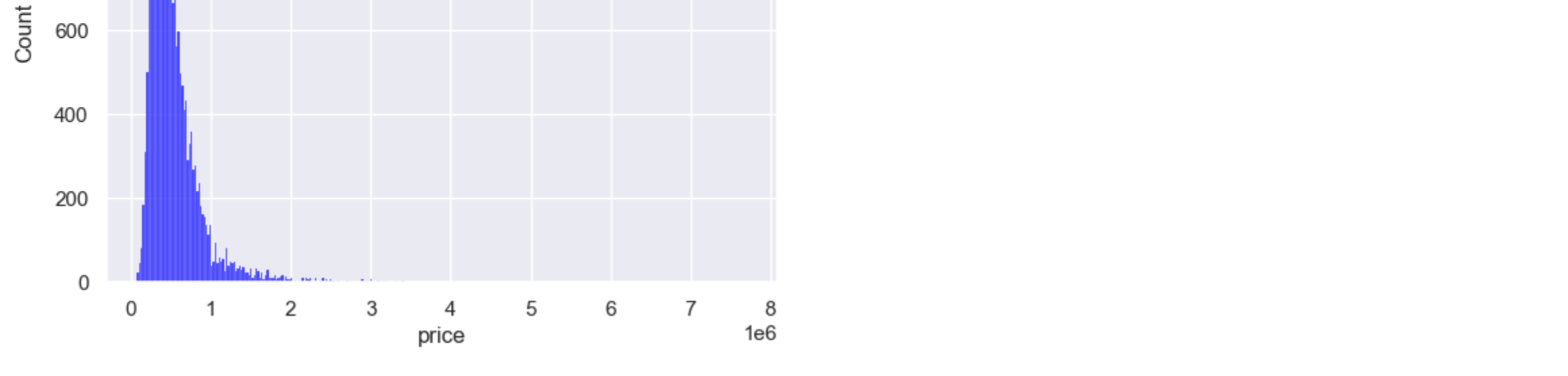
```
In [16]: sns.histplot(x='log_price',data=data,color='blue')
```

```
Out[17]: <Axes: xlabel='log_price', ylabel='Count'>
```



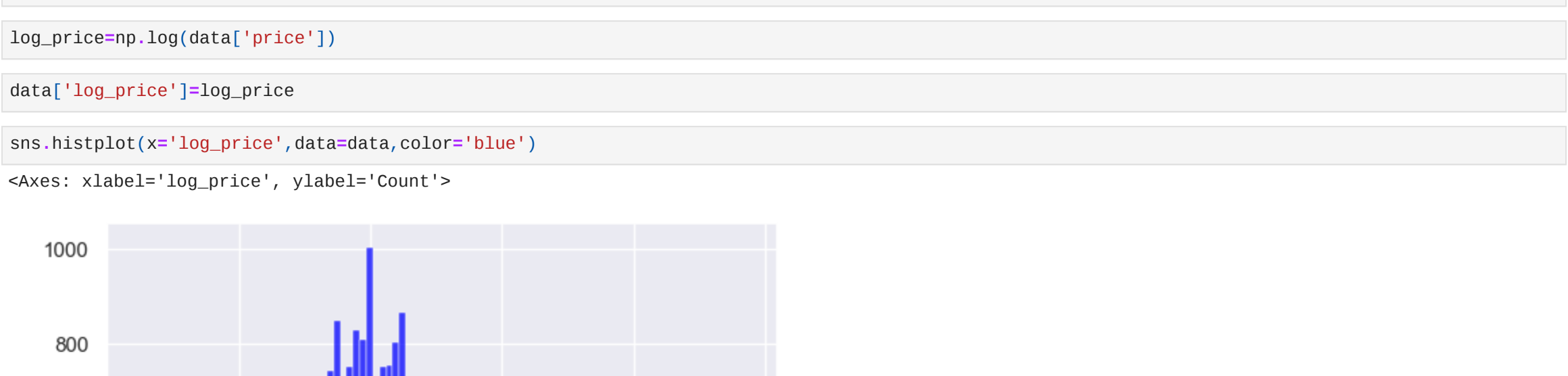
```
In [18]: #barplot on bedrooms against price
plt.figure(figsize=(15,6))
ax=sns.barplot(x=data['bedrooms'], y=data['price'],palette='bright')
for i in ax.containers:
    ax.bar_label(i,
plt.title('Bedrooms VS Price', fontsize=14)
Text(0.5, 1.0, 'Bedrooms VS Price')
```

```
Out[18]:
```



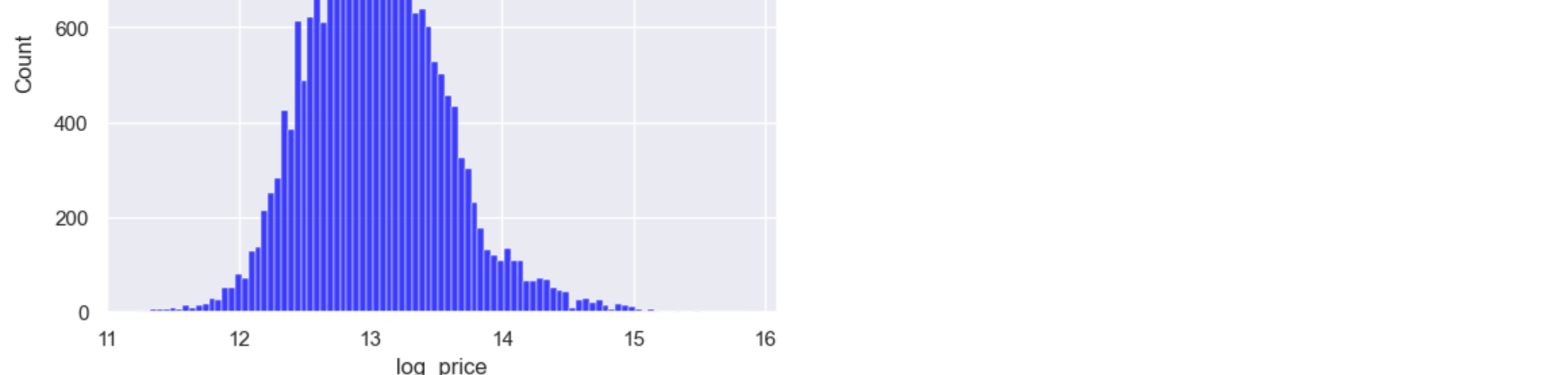
```
In [19]: #barplot on bedrooms against sqft_living
plt.figure(figsize=(15,6))
ax=sns.barplot(x=data['bedrooms'], y=data['sqft_living'],palette='Set1')
for i in ax.containers:
    ax.bar_label(i,
plt.title('Sqft_living VS Price', fontsize=14)
Text(0.5, 1.0, 'Sqft_living VS Price')
```

```
Out[19]:
```



```
In [20]: #barplot on floors against price
plt.figure(figsize=(15,6))
ax=sns.barplot(x=data['floors'], y=data['price'],palette='bright')
for i in ax.containers:
    ax.bar_label(i,
plt.title('Floors VS Price', fontsize=14)
Text(0.5, 1.0, 'Floors VS Price')
```

```
Out[20]:
```



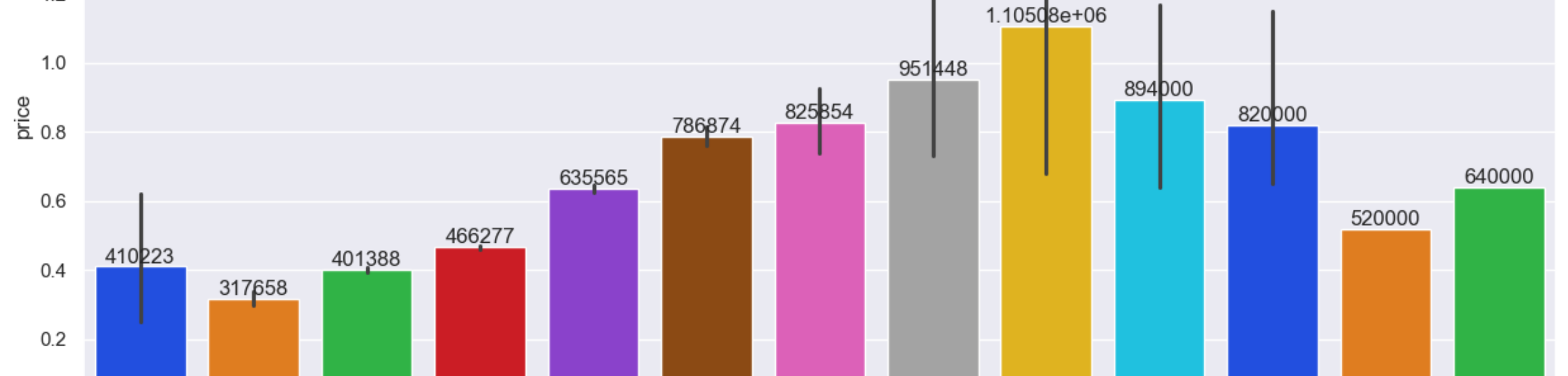
```
In [21]: #boxplot on waterfront against price
sns.boxplot(x='waterfront',y='price',data=data,palette='plasma')
```

```
Out[21]: <Axes: xlabel='waterfront', ylabel='price'>
```



```
In [22]: #violin plot on view against price
sns.violinplot(x='view',y='price',data=data,palette='plasma')
```

```
Out[22]: <Axes: xlabel='view', ylabel='price'>
```



```
In [23]: #boxplot on condition against price
sns.boxplot(x='condition',y='price',data=data,palette='plasma')
```

```
Out[23]: <Axes: xlabel='condition', ylabel='price'>
```



## DATA CLEANING

```
In [24]: #dropping irrelevant columns
data.drop(columns='date',axis=1,inplace=True)
```

```
Out[24]:
```

```
In [25]: data.shape
```

```
Out[25]: (21613, 15)
```

```
In [26]: #correlation
data.corr()
```

	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	condition	grade	sqft_above	sqft_basement	yr_built	yr_renovated	zipcode
price	1.000000	0.390338	0.523134	0.702044	0.089956	0.256766	0.266321	0.397348	0.036392	0.687463	0.323937	0.039362	0.125442	-0.053136	
bedrooms	0.390338	1.000000	0.515894	0.576671	0.031703	0.175429	0.056582	0.078523	0.028472	0.355967	0.477600	0.302093	0.154178	0.018841	-0.152686
bathrooms	0.523134	0.515894	1.000000	0.754665	0.087740	0.500053	0.063744	0.187737	-0.124982	0.654963	0.696342	0.329770	0.509019	0.059739	-0.20386
sqft_living	0.702044	0.576671	0.754665	1.000000	0.172826	0.353041	0.193818	0.284611	-0.056753	0.762704	0.876597	0.493043	0.319408	0.055363	-0.19943
sqft_lot	0.089956	0.031703	0.087740	0.172826	1.000000	-0.005201	0.021604	0.074170	-0.008958	0.113821	0.139312	0.015296	0.053090	0.007644	-0.12967
floors	0.256766	0.175429	0.500053	0.353041	0.005201	1.000000	0.026986	0.029444	-0.008958	0.458183	0.523885	-0.245765	0.489319	0.006338	-0.05912
waterfront	0.266321	0.063744	0.193818	0.172826	0.005201	0.026986	1.000000	0.401857	0.016653	0.082775	0.072075	0.005958	-0.026161	0.092885	0.02028
view	0.397348	0.079532	0.187737	0.284611	0.074170	0.029444	0.401857	1.000000	0.009990	0.251321	0.167649	0.279947	-0.354140	0.103917	0.08482
condition	0.036392	0.028472	-0.124982	-0.056753	-0.005201	-0.026986	0.016653	0.009990	1.000000	-0.144674	-0.158214	0.174105	-0.361417	-0.060618	0.00302
grade	0.687463	0.355967	0.654963	0.762704	0.113821	0.458183	0.082775	0.251321	-0.144674	1.000000	0.759223	0.500000	-0.051943	0.144414	-0.18486
sqft_above	0.696565	0.477600	0.685342	0.876597	0.139312	0.523885	0.072075	0.167649	-0.158214	0.759223	1.000000	0.000000	-0.051943	0.423998	-0.26119
sqft_basement	0.323937	0.302093	0.329770	0.509019	0.015296	0.006338	0.005958	0.276947	0.174105	0.168392	-0.051943	1.000000	-0.133124	-0.071323	0.07484
yr_built	0.053082	0.154178	0.509019	0.435043	0.030600	0.489319	-0.026161	-0.053440	-0.361417	0.446963	0.423998	-0.133124	1.000000	-0.244874	-0.04646
yr_renovated	0.125442	0.018841	0.059739	0.055063	0.007644	0.006338	0.002885	0.103917	-0.060618	0.144414	0.232385	0.071323	-0.244874	1.000000	0.06435
zipcode	-0.053136	-0.152686	-0.20386	-0.199430	-0.129574	-0.059121	0.030295	0.048627	0.003026	-0.184862	-0.261190	0.071323	-0.244874	0.064357	1.00000
lat	0.020519	0.128473	0.223042	0.240223	0.229521	0.125419	-0.041910	0.074000	-0.106500	0.198372	0.343803	-0.144765	0.409356	-0.068372	-0.56407
sqft_living15	0.585374	0.391638	0.568634	0.756420	0.144608	0.279885	0.066403	0.280439	-0.092824	0.713202	0.731870	0.020756	0.326229	-0.002073	-0.27903
sqft_lot15	0.082456	0.029244	0.087175	0.183260	0.718557	-0.011269	0.030703	0.072575	-0.003406	0.119248	0.194050	0.013795	0.070958	0.007854	-0.14722
log_price	0.891680	0.343563	0.550819	0.695380	0.099625	0.310566	0.174597	0.340576	0.039574	0.703675	0.601380	0.311698	0.080338	0.114512	-0.03829

```
In [27]: #heatmap on correlation
plt.figure(figsize=(14,10))
sns.heatmap(data.corr(),annot=True)
```

```
Out[27]: &lt
```