

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5 from scipy import stats
        6 %matplotlib inline
```

```
In [2]: 1 df = pd.read_csv('//Users//omniaelmenhawy//Desktop//machinfy//
        2 df.head(5)
```

Out[2]:

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	house
0	-122.23	37.88	41.0	880	129.0	322.0	
1	-122.22	37.86	21.0	7099	1106.0	2401.0	
2	-122.24	37.85	52.0	1467	190.0	496.0	
3	-122.25	37.85	52.0	1274	235.0	558.0	
4	-122.25	37.85	NaN	1627	280.0	NaN	

Data Analysis:

```
In [3]: 1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20640 entries, 0 to 20639
Data columns (total 11 columns):
#   Column                Non-Null Count  Dtype
---  -
0   longitude              20640 non-null  float64
1   latitude               20640 non-null  float64
2   housing_median_age     20382 non-null  float64
3   total_rooms            20640 non-null  int64
4   total_bedrooms         15758 non-null  float64
5   population             20596 non-null  float64
6   households              19335 non-null  object
7   median_income          17873 non-null  float64
8   median_house_value     20640 non-null  int64
9   ocean_proximity        20640 non-null  object
10  gender                 16620 non-null  object
dtypes: float64(6), int64(2), object(3)
memory usage: 1.7+ MB
```

```
In [4]: 1 df = df.drop_duplicates()
```

In [5]: `1 print(df.shape)`

(20640, 11)

In [6]: `1 df.describe()`

Out[6]:

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	
count	20640.000000	20640.000000	20382.000000	20640.000000	15758.000000	20640
mean	-119.569704	35.631861	28.676283	2635.763081	539.920104	14
std	2.003532	2.135952	12.589284	2181.615252	419.834171	17
min	-124.350000	32.540000	1.000000	2.000000	1.000000	
25%	-121.800000	33.930000	18.000000	1447.750000	296.000000	7
50%	-118.490000	34.260000	29.000000	2127.000000	435.000000	17
75%	-118.010000	37.710000	37.000000	3148.000000	652.000000	17
max	-114.310000	41.950000	52.000000	39320.000000	6210.000000	356

In [7]: `1 missing_values = df.isnull().sum()`

`2 print(missing_values)`

```
longitude          0
latitude           0
housing_median_age 258
total_rooms        0
total_bedrooms    4882
population         44
households        1305
median_income      2767
median_house_value 0
ocean_proximity   0
gender            4020
dtype: int64
```

```
In [8]: 1 #sample of the missing values:
        2 df[df.isnull().any(axis=1)]
```

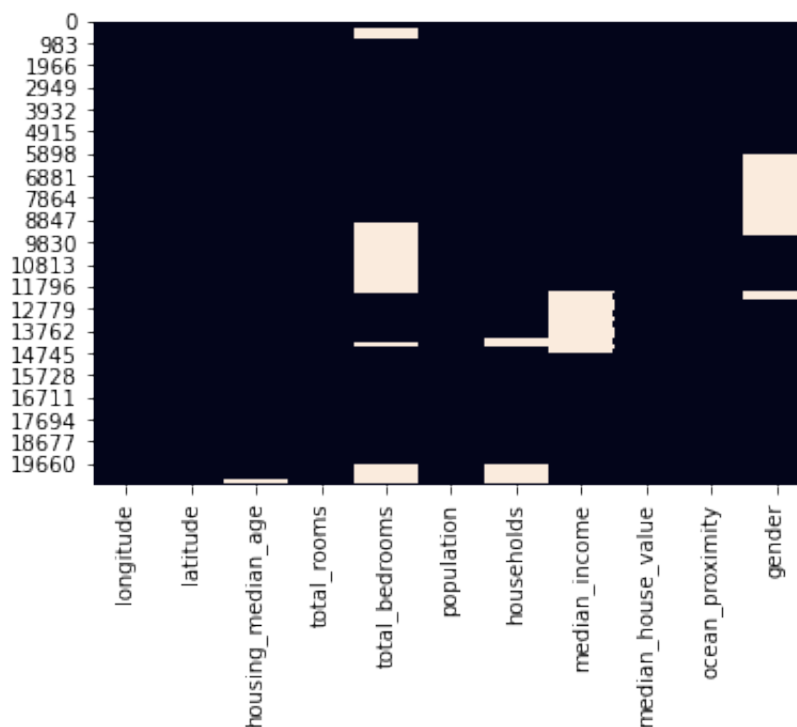
Out[8]:

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	h
4	-122.25	37.85	NaN	1627	280.0	NaN	
5	-122.25	37.85	NaN	919	213.0	NaN	
6	-122.25	37.84	NaN	2535	NaN	NaN	
7	-122.25	37.84	NaN	3104	NaN	NaN	
8	-122.26	37.84	42.0	2555	NaN	NaN	
...
20627	-121.32	39.13	NaN	358	NaN	169.0	
20628	-121.48	39.10	NaN	2043	NaN	1018.0	
20629	-121.39	39.12	NaN	10035	NaN	6912.0	
20630	-121.32	39.29	NaN	2640	NaN	1257.0	
20631	-121.40	39.33	15.0	2655	493.0	1200.0	

10463 rows × 11 columns

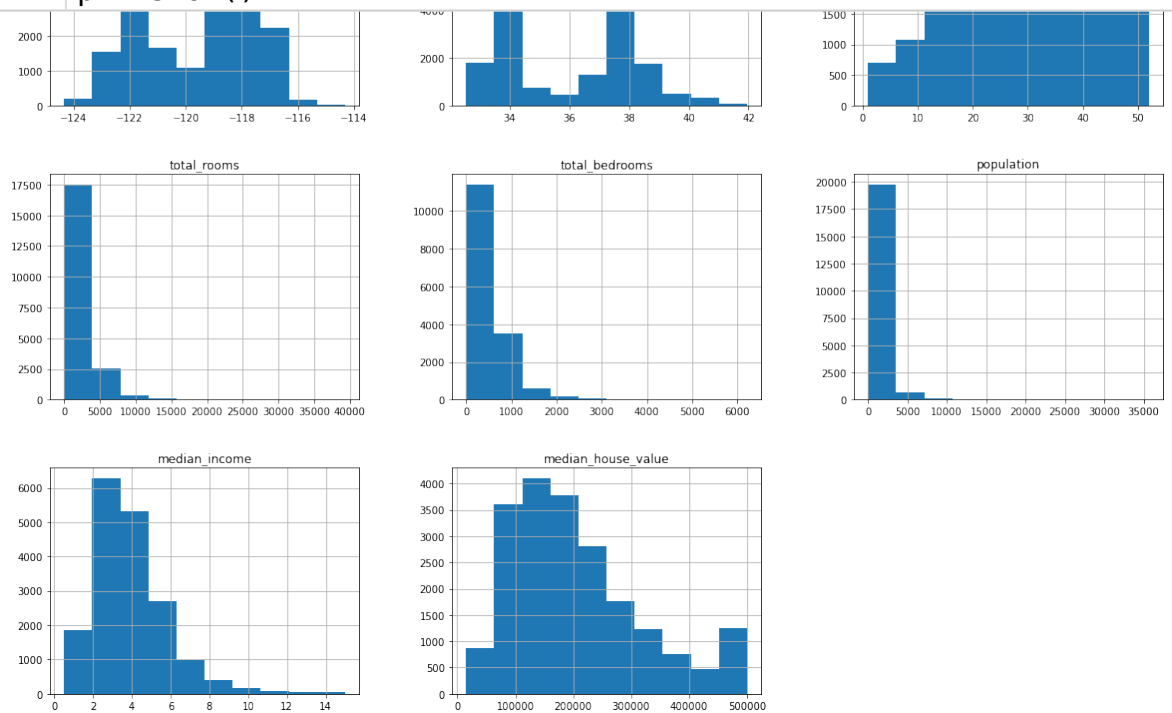
```
In [9]: 1 sns.heatmap(df.isnull(), cbar=False)
        2 print("Total Missing: ", df.isna().sum().values.sum())
```

Total Missing: 13276



Data Visualization Before:

```
In [10]: 1 df.hist(figsize=(20,15))  
        2 plt.show()
```

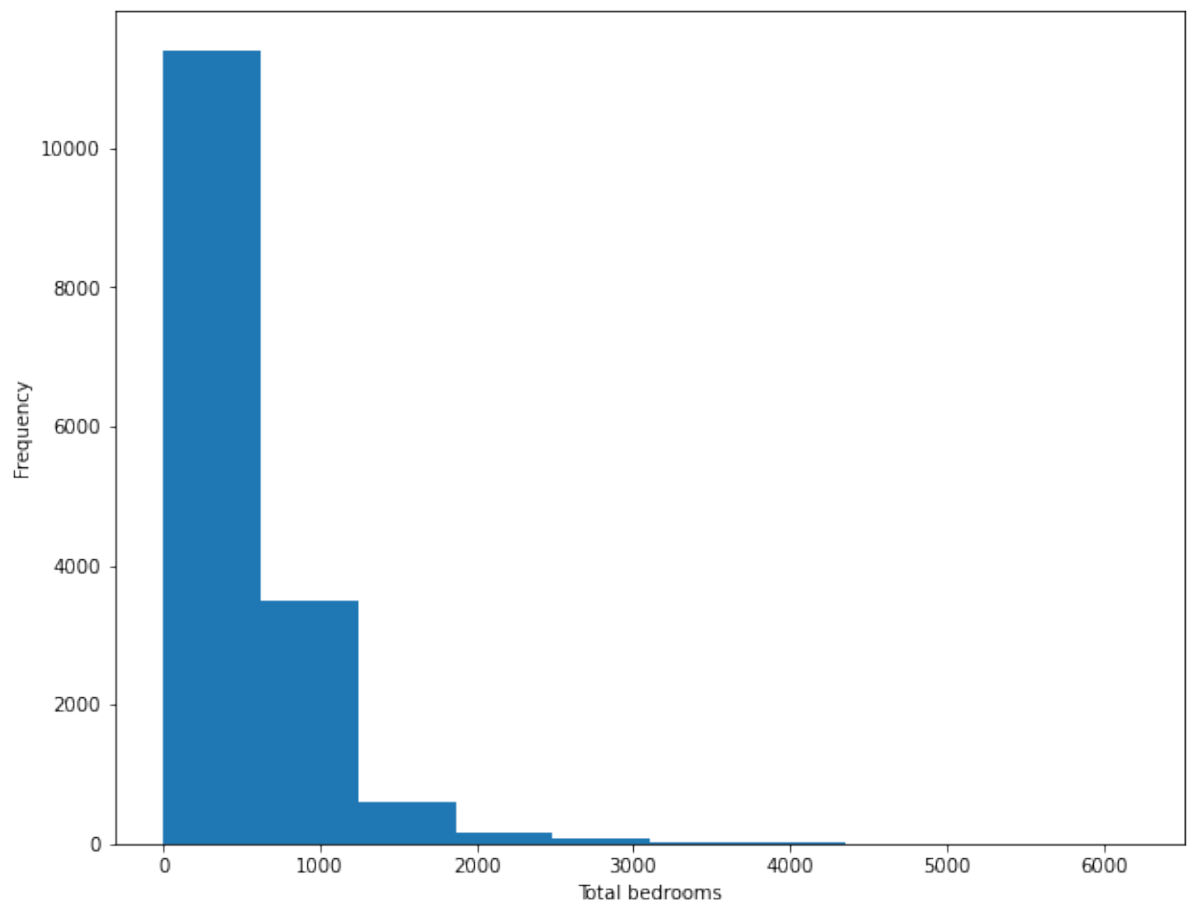


Most houses have bedrooms between 0 to 1000.

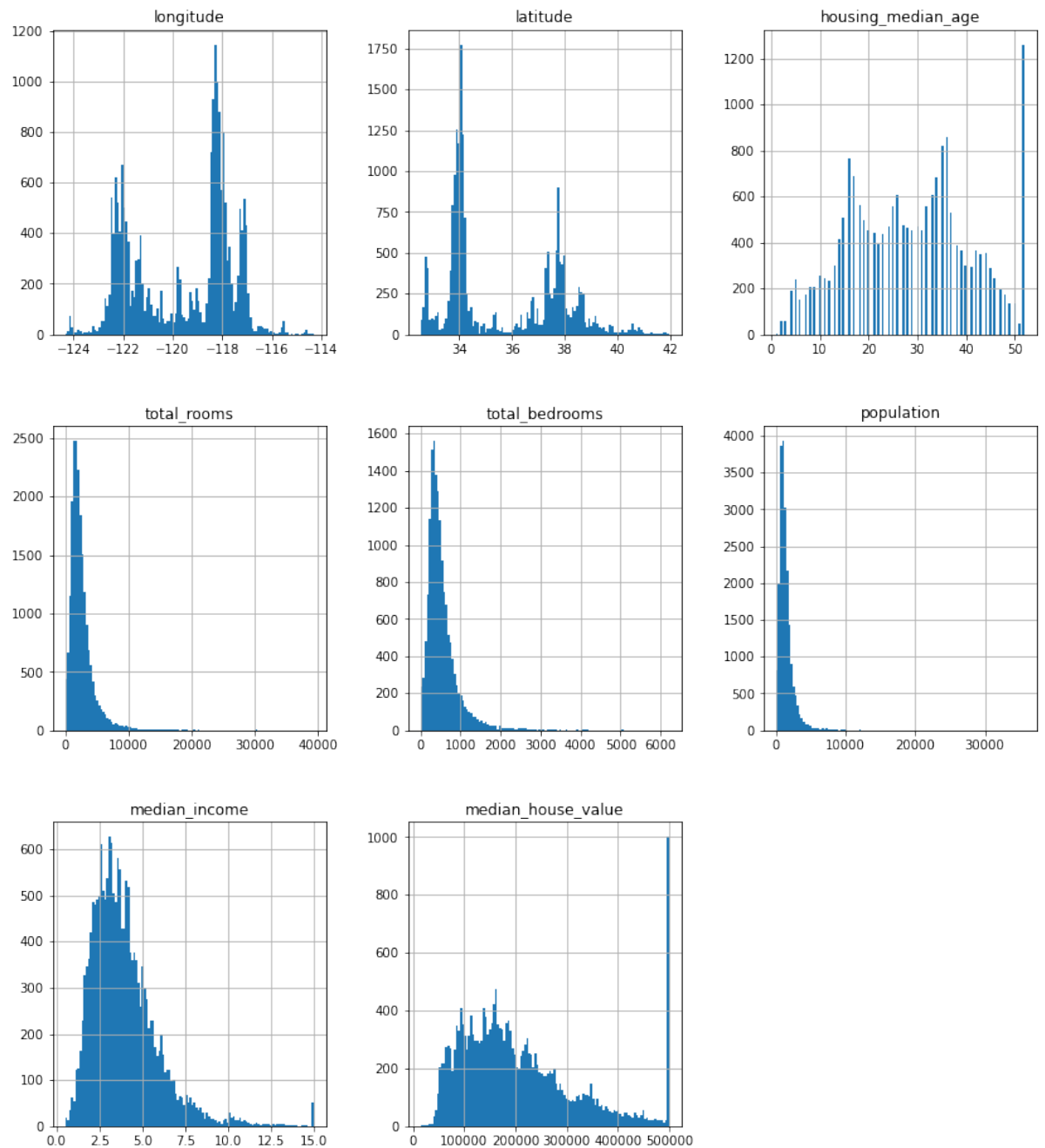
Houses with more than 1000 room are rare.

```
In [11]: 1 plt.figure(figsize=(10,8))
          2 df['total_bedrooms'].plot(kind='hist')
          3 plt.xlabel('Total bedrooms')
```

Out[11]: Text(0.5, 0, 'Total bedrooms')

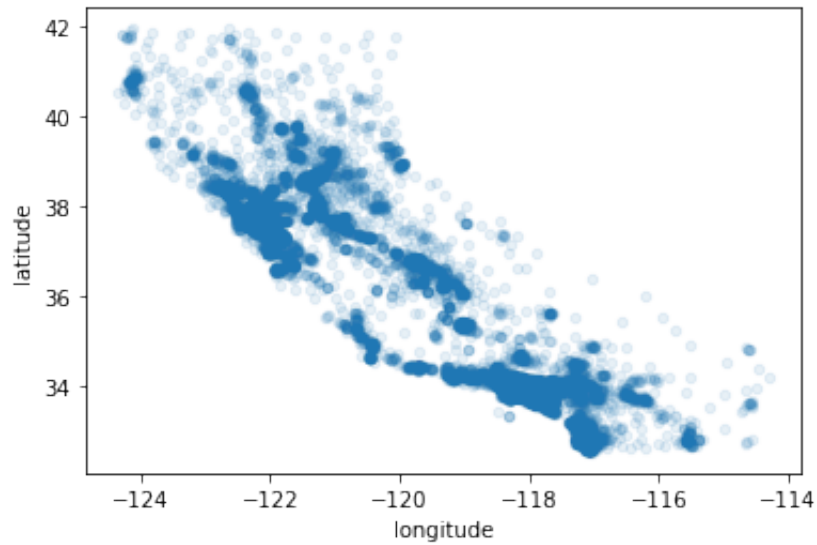


```
In [12]: 1 df.hist(bins=120, figsize = (14,16))  
        2 plt.show()
```



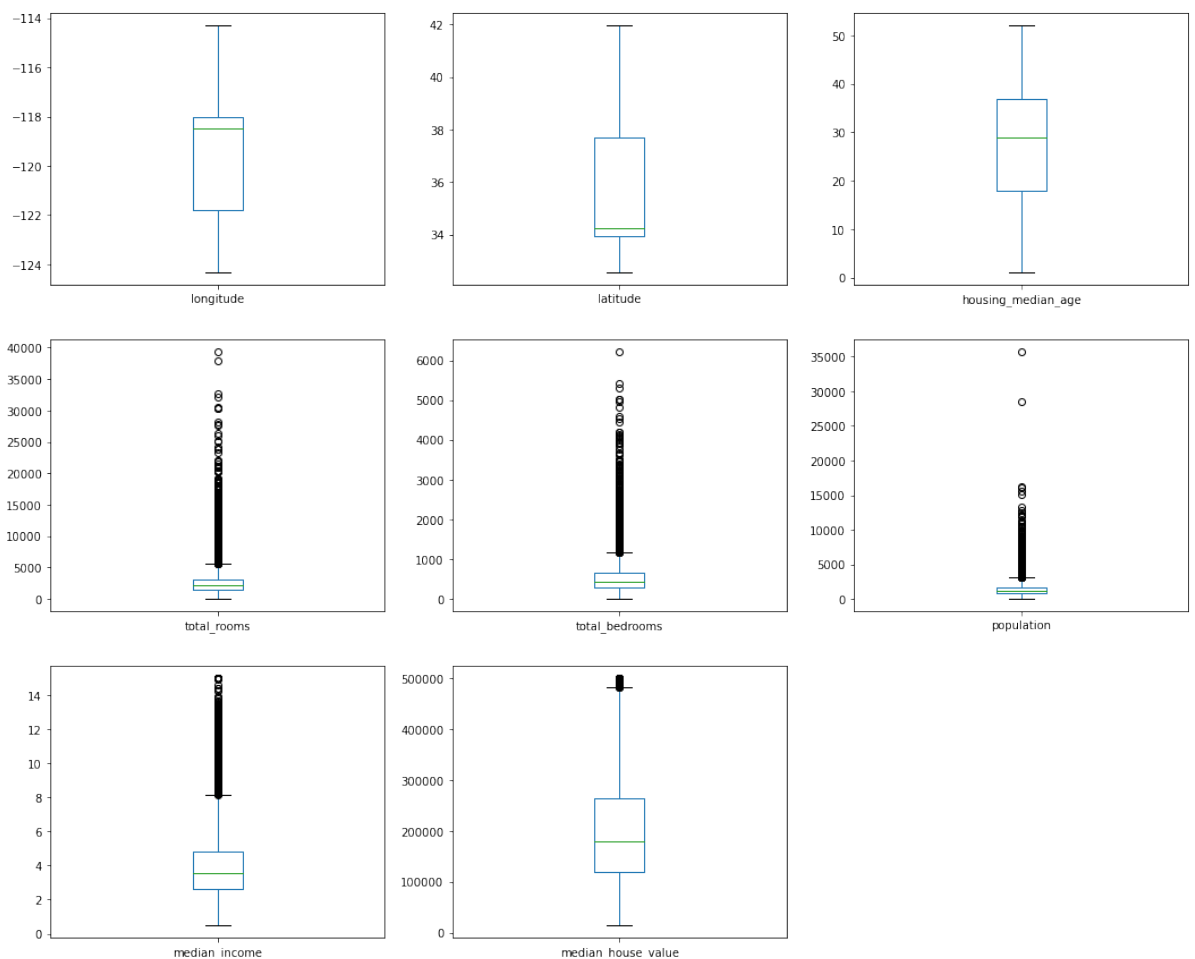
```
In [13]: 1 df.plot(kind='scatter', x='longitude', y='latitude', alpha = 0.
```

```
Out[13]: <AxesSubplot:xlabel='longitude', ylabel='latitude'>
```



Outliers:

```
In [14]: 1 df.plot(kind='box', subplots = True, figsize=(18,15),layout=(3,
2         plt.show())
```



Dropping unnecessary data:

In [15]: `1 df.columns`

Out[15]: Index(['longitude', 'latitude', 'housing_median_age', 'total_rooms',
'total_bedrooms', 'population', 'households', 'median_income',
'median_house_value', 'ocean_proximity', 'gender'],
dtype='object')

Dropping the gender column:

In [16]: `1 df = df[['longitude', 'latitude', 'housing_median_age', 'total_'
2 'total_bedrooms', 'population', 'households', 'median_in'
3 'median_house_value', 'ocean_proximity']]
4 df.head(1)`

Out[16]:

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	house
0	-122.23	37.88	41.0	880	129.0	322.0	

Handelling the missing data:

In [17]: `1 df['total_bedrooms'].mode()`

Out[17]: 0 280.0
dtype: float64

In [65]: `1 df.loc[df['total_bedrooms'].isna(), 'total_bedrooms'] = df['total_`

In [19]: `1 df.loc[df['median_income'].isna(), 'median_income'] = df['median_i`

In [20]: `1 df[(df['population']>4000)]`

Out[20]:

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	h
185	-122.23	37.79	43.0	5963	1344.000000	4367.0	
283	-122.16	37.79	22.0	12842	2048.000000	4985.0	
570	-122.24	37.72	5.0	18634	539.920104	7427.0	
576	-122.06	37.77	12.0	14316	539.920104	5781.0	
780	-122.10	37.63	18.0	9963	539.920104	5613.0	
...
20530	-121.76	38.57	NaN	15018	539.920104	7984.0	
20539	-121.71	38.56	NaN	8627	539.920104	4071.0	
20544	-121.76	38.55	NaN	8800	539.920104	6330.0	
20563	-121.75	38.67	NaN	12139	539.920104	6837.0	
20629	-121.39	39.12	NaN	10035	539.920104	6912.0	

582 rows × 10 columns

In [29]: `1 df.loc[df['population'].isna(), 'population']=df['population'].m`

In [30]: `1
2 df.loc[df['housing_median_age'].isna(), 'housing_median_age']=df`

In [31]: `1 df.head(3)`

Out[31]:

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	house
0	-122.23	37.88	41.0	880	129.0	322.0	
1	-122.22	37.86	21.0	7099	1106.0	2401.0	
2	-122.24	37.85	52.0	1467	190.0	496.0	

In [40]: `1 df = df[df.isnull().sum(axis=1) < 3]`

In [41]: 1 df

Out [41]:

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population
0	-122.23	37.88	41.000000	880	129.0	322.000000
1	-122.22	37.86	21.000000	7099	1106.0	2401.000000
2	-122.24	37.85	52.000000	1467	190.0	496.000000
3	-122.25	37.85	52.000000	1274	235.0	558.000000
4	-122.25	37.85	28.676283	1627	280.0	1424.928724
...
20635	-121.09	39.48	25.000000	1665	374.0	845.000000
20636	-121.21	39.49	18.000000	697	150.0	356.000000
20637	-121.22	39.43	17.000000	2254	485.0	1007.000000
20638	-121.32	39.43	18.000000	1860	409.0	741.000000
20639	-121.24	39.37	16.000000	2785	616.0	1387.000000

20640 rows × 10 columns

In [42]: 1 df["households"]=df['households'].fillna(df["households"].mode(

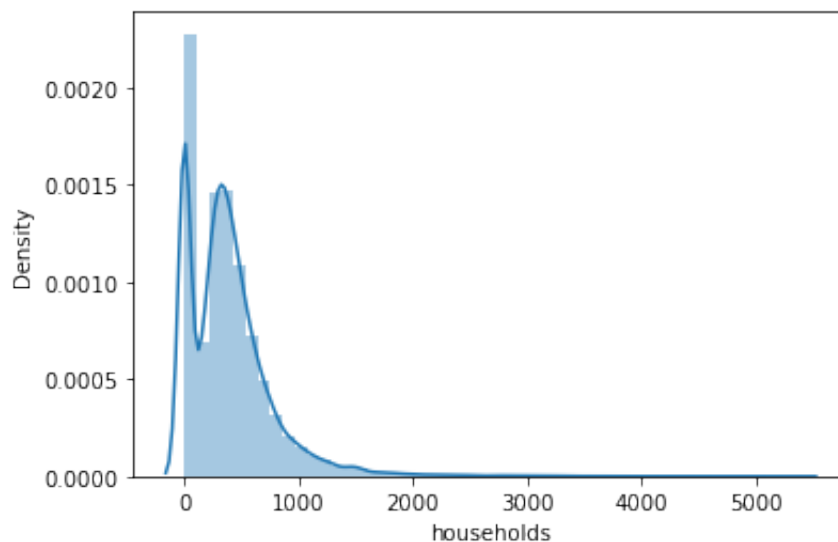
In [43]: 1 df.replace('no' ,0,inplace=True)

```
In [44]: 1 sns.distplot(df["households"])
```

/Users/omniaelmenshawy/opt/anaconda3/lib/python3.8/site-packages/seaborn/distributions.py:2551: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

```
warnings.warn(msg, FutureWarning)
```

```
Out [44]: <AxesSubplot:xlabel='households', ylabel='Density'>
```

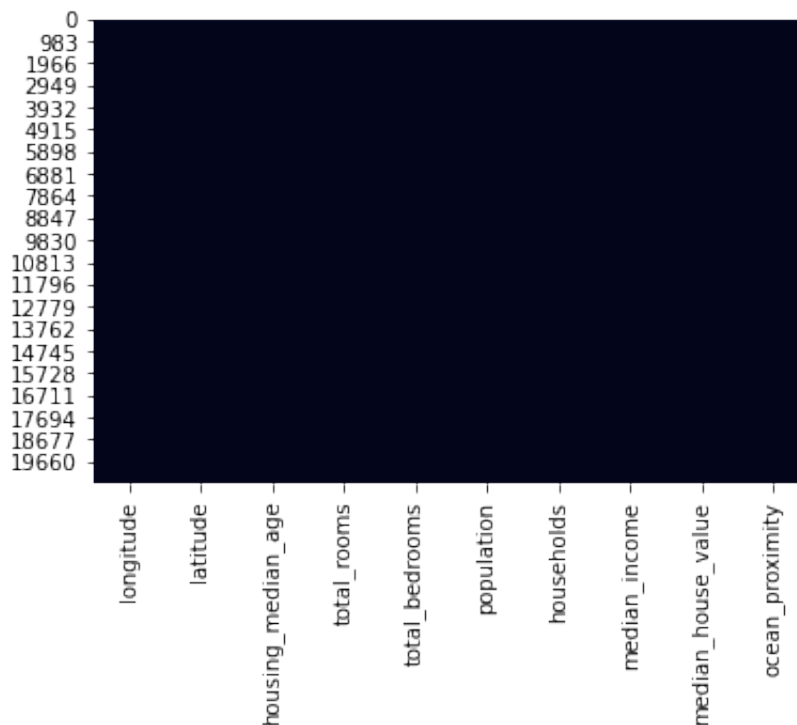


```
In [45]: 1 print(df.isnull().sum())
```

```
longitude      0
latitude       0
housing_median_age  0
total_rooms    0
total_bedrooms 0
population     0
households     0
median_income  0
median_house_value 0
ocean_proximity 0
dtype: int64
```

```
In [46]: 1 sns.heatmap(df.isnull(), cbar=False)
          2 print("Total Missing: ", df.isna().sum().values.sum())
```

Total Missing: 0



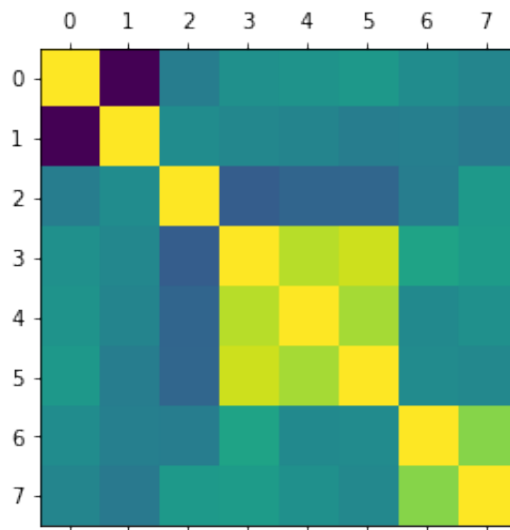
Correlation:

- the median house value is positively correlated with the median income and has a great impact
- while it is negatively correlated with the:
 - latitude: the value decreases when we go north
 - longitude: the value also decreases when we go west
 - population: the value decreases when the total number of people living in the block increases

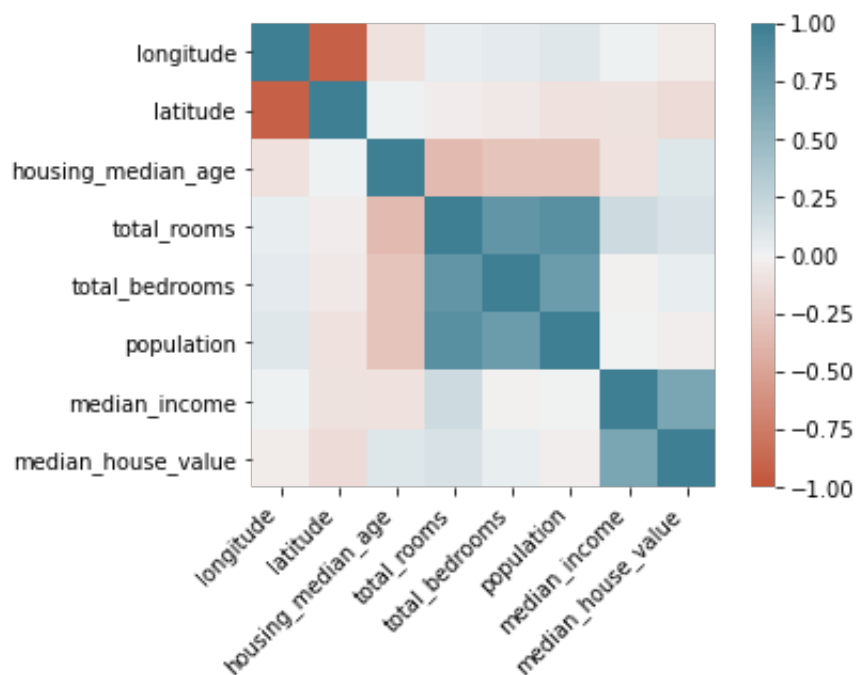
```
In [49]: 1 corr_matrix = df.corr()
          2 corr_matrix["median_house_value"].sort_values(ascending = False)
```

```
Out[49]: median_house_value    1.000000
          median_income      0.650304
          total_rooms        0.134153
          housing_median_age  0.106648
          total_bedrooms     0.044949
          population        -0.024351
          longitude         -0.045967
          latitude          -0.144160
          Name: median_house_value, dtype: float64
```

```
In [51]: 1 plt.matshow(df.corr())
          2 plt.show()
```

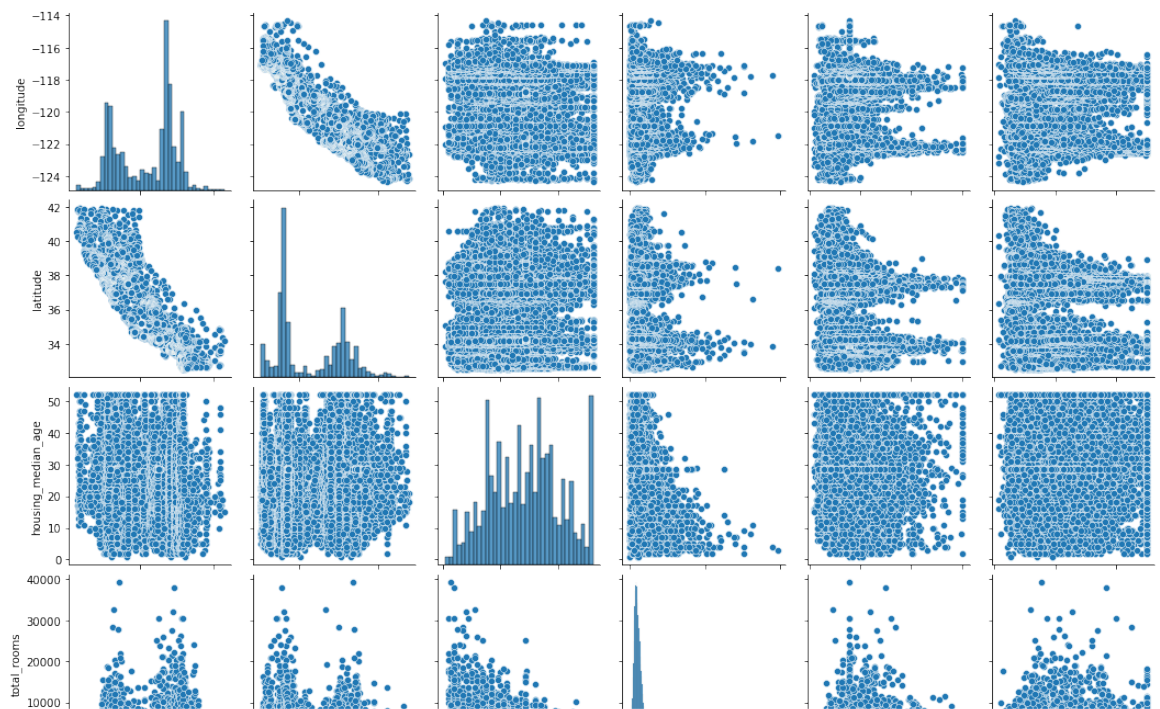


```
In [53]: 1 corr = df.corr()
          2 ax = sns.heatmap(
          3     corr,
          4     vmin=-1, vmax=1, center=0,
          5     cmap=sns.diverging_palette(20, 220, n=200),
          6     square=True
          7 )
          8 ax.set_xticklabels(
          9     ax.get_xticklabels(),
         10     rotation=45,
         11     horizontalalignment='right'
         12 );
```



Data Visualization after Handelling it:

```
In [58]: 1 column_study = ["longitude", "latitude", "housing_median_age",  
2          sns.pairplot(df[column_study], height = 2.5);  
3          plt.show()
```

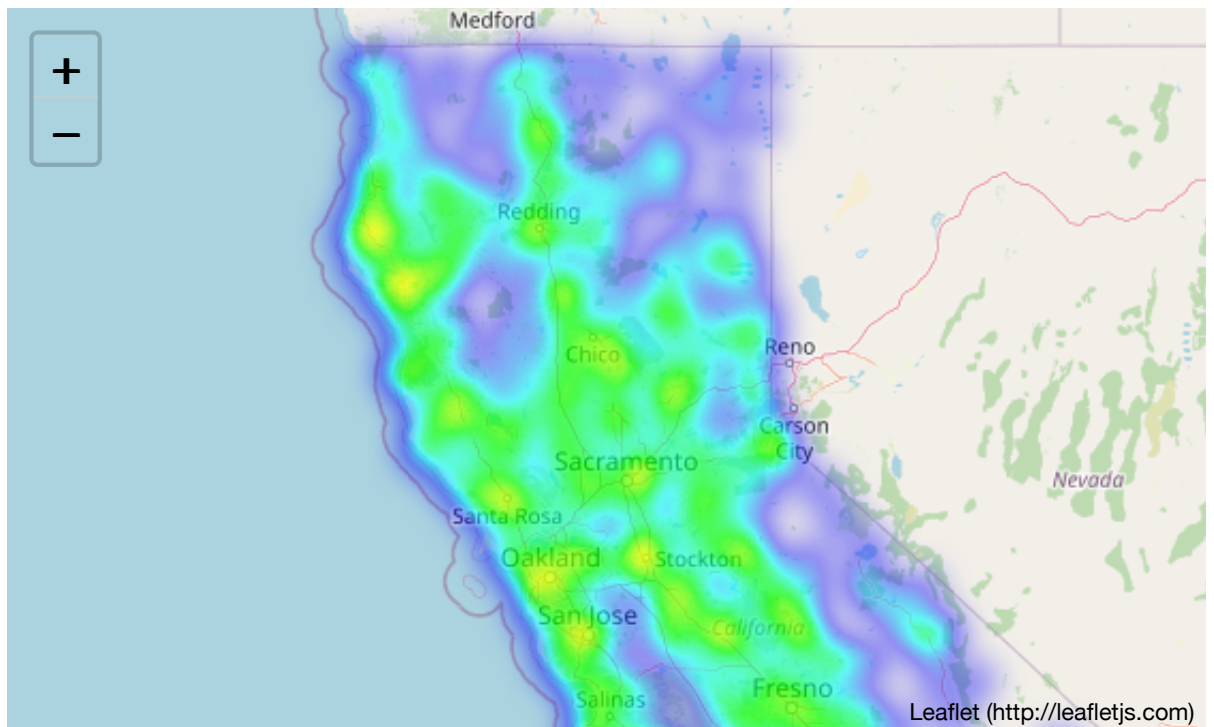


```

In [113]: 1 import folium
          2 from folium.plugins import HeatMap, MarkerCluster
          3 from folium import Choropleth, Circle, Marker
          4
          5
          6 cal_map = folium.Map(location=[36.7783, -119.4179], zoom_start =
          7 df_map = df[['latitude', 'longitude']]
          8 data = [[row['latitude'], row['longitude']] for index, row in df
          9 _ = HeatMap(data, radius=10).add_to(cal_map)
         10 cal_map

```

Out[113]:

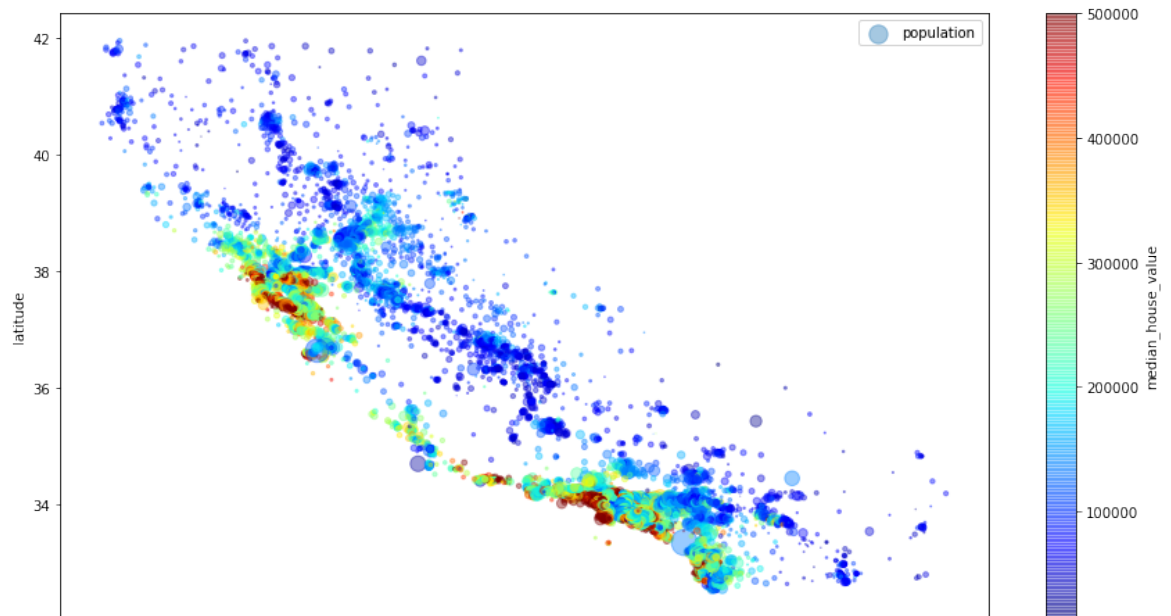


shape of price distribution by the latitude and longitude:

- high-density areas are:
 - the Bay Area
 - around Los Angeles and San Diego
 - plus a long line of fairly high density in the Central Valley, around Sacramento and Fresno.

```
In [152]: 1 df.plot(kind="scatter", x="longitude", y="latitude", alpha=0.4,
2           s=df["population"]/100, label="population", figsize=(15,8),
3           c="median_house_value", cmap=plt.get_cmap("jet"), colorbar=True
4           )
5           plt.legend()
```

Out[152]: <matplotlib.legend.Legend at 0x7fbf022ef850>



```
In [153]: 1 nearbay_df = df.loc[df['ocean_proximity'] == 'NEAR BAY']
2           nearbay_df.head()
```

Out[153]:

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	house_value
0	-122.23	37.88	41.000000	880	129.0	322.000000	
1	-122.22	37.86	21.000000	7099	1106.0	2401.000000	
2	-122.24	37.85	52.000000	1467	190.0	496.000000	
3	-122.25	37.85	52.000000	1274	235.0	558.000000	
4	-122.25	37.85	28.676283	1627	280.0	1424.928724	

```
In [ ]: 1
```



```
In [108]: 1 latitude = 36.7783
2 longitude = -119.4179
3 island_df = df.loc[df['ocean_proximity'] == 'ISLAND']
4
5 m = folium.Map(location = [latitude, longitude], tiles='cartodb
6 for idx, row in island_df.iterrows():
7     Marker([row['latitude'], row['longitude']]).add_to(m)
8 m
```

Out[108]:



Maximum Values around the 500001 Dollar and they are all near the ocean

```
In [132]: 1 print("House statistics:\n")
2 print(df['median_house_value'].describe())
```

House statistics:

```
count      20640.000000
mean       206855.816909
std        115395.615874
min         14999.000000
25%        119600.000000
50%        179700.000000
75%        264725.000000
max         500001.000000
Name: median_house_value, dtype: float64
```

Linear Regression and Training:

```
In [133]: 1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LinearRegression
```

```
1 df.columns
```

```
Index(['longitude', 'latitude', 'housing_median_age', 'total_rooms',
      'total_bedrooms', 'population', 'households', 'median_income',
      'median_house_value', 'ocean_proximity'],
      dtype='object')
```

```
1 Training = df[['longitude', 'latitude', 'housing_median_age', '
2               'total_bedrooms', 'population', 'households', 'median_in
3               'median_house_value']]
```

```
1 x = Training.drop(['median_house_value'],axis = 1).values
2 y = Training["median_house_value"].values
```

```
1 x_train, x_test, y_train, y_test = train_test_split(x,y,test_si
```

```
1 x_train.shape
```

 $(15480, 8)$

```
1 x_test.shape
```

(5160, 8)

```
1 y_train.shape
```

(15480,)

```
1 y_test.shape
```

(5160,)

```
1 100*x_test.shape[0]/x_train.shape[0]
```

33.333333333333336

```
1 lr = LinearRegression()
2 lr.fit(x_train, y_train)
3 y_predict = lr.predict(x_test)
```

```
1 y_predict = lr.predict(x_test)
```

```
1 lr.score(x_train, y_train)
```

0.5868339060852498

```
In [147]: 1 lr.score(x_test,y_test)
```

```
Out[147]: 0.581149194438561
```

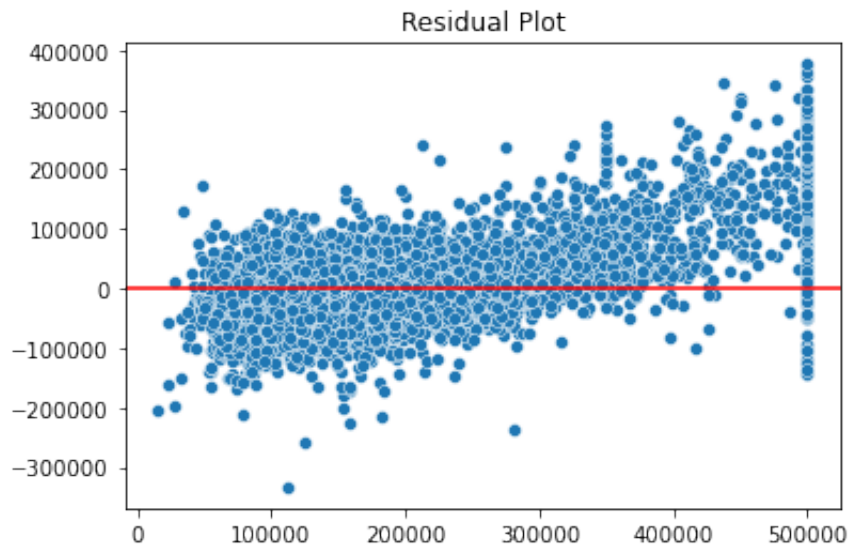
Model Visualization:

Errors are close to normal distribution and also the mean is pretty close to zero

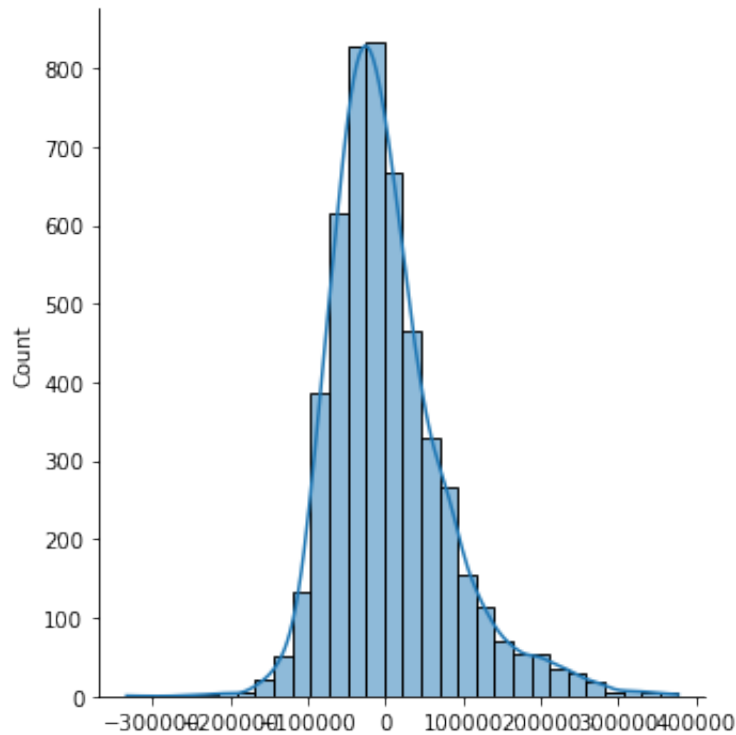
```
In [116]: 1 test_residuals = y_test - y_predict  
2 test_residuals  
3
```

```
Out[116]: array([ 11185.38093889, -61306.43844511, 13715.36551086, ...,  
-49592.36875296, -73804.00073398, -84994.9557901 ])
```

```
In [117]: 1 sns.scatterplot(x=y_test,y=test_residuals)  
2 plt.axhline(y=0,color='red')  
3 plt.title('Residual Plot');
```



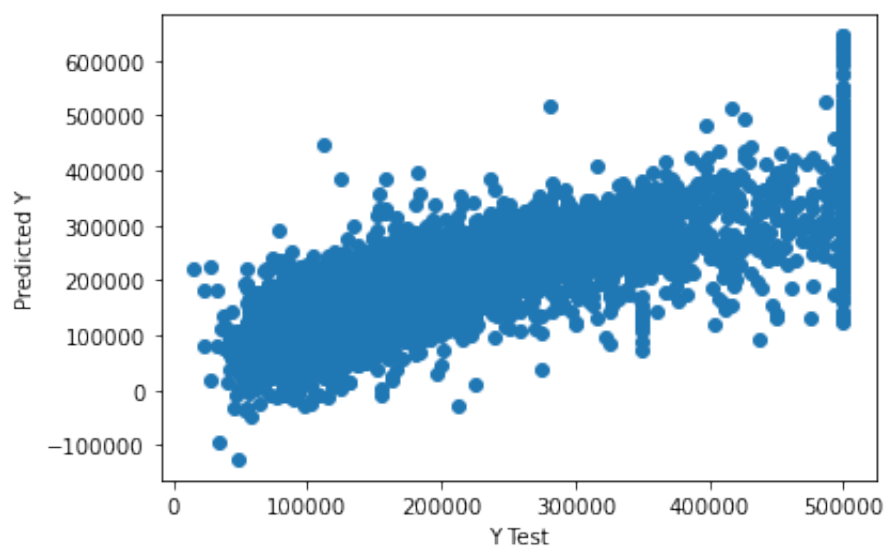
```
In [148]: 1 sns.displot(test_residuals,bins=30,kde=True);
```



```
In [149]: 1 predictions = lr.predict(x_test)
```

```
In [150]: 1 plt.scatter(x = y_test, y = predictions)
          2 plt.xlabel('Y Test')
          3 plt.ylabel('Predicted Y')
```

```
Out[150]: Text(0, 0.5, 'Predicted Y')
```



```
In [ ]: 1
```

