

INTRODUCTION OF THE PROJECT : In this California Dataset.They have provided with Longitude,latitude,total incomes and household as per the standards.This sample data is provided to get the idea about how people live thier lifes according to their median incomes.They also provided the overall population in data set.

This Dataset has around 20,640 observation in it with 10 columns and it's a mixed dataset containing categorical & numerical values.

```
In [1]: #importing necessary libraries..
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

```
In [2]: df=pd.read_csv('housing.csv') #importing csv file
```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
0	-122.23	37.88	41	880	129.0	322	126	8.3252	452600	NEAR BAY
1	-122.22	37.86	21	7099	1106.0	2401	1138	8.3014	358500	NEAR BAY
2	-122.24	37.85	52	1467	190.0	496	177	7.2574	352100	NEAR BAY
3	-122.25	37.85	52	1274	235.0	558	219	5.6431	341300	NEAR BAY
4	-122.25	37.85	52	1627	280.0	565	259	3.8462	342200	NEAR BAY
...
20635	-121.09	39.48	25	1665	374.0	845	330	1.5603	78100	INLAND
20636	-121.21	39.49	18	697	150.0	356	114	2.5568	77100	INLAND
20637	-121.22	39.43	17	2254	485.0	1007	433	1.7000	92300	INLAND
20638	-121.32	39.43	18	1860	409.0	741	349	1.8672	84700	INLAND
20639	-121.24	39.37	16	2785	616.0	1387	530	2.3886	89400	INLAND

20640 rows × 10 columns

```
In [4]: df.head(10) #importing n number of starting data...
```

```
Out[4]:
```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
0	-122.23	37.88	41	880	129.0	322	126	8.3252	452600	NEAR BAY
1	-122.22	37.86	21	7099	1106.0	2401	1138	8.3014	358500	NEAR BAY

```
In [4]: df.head(10) #importing n number of starting data...
```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
0	-122.23	37.88	41	880	129.0	322	126	8.3252	452600	NEAR BAY
1	-122.22	37.86	21	7099	1106.0	2401	1138	8.3014	358500	NEAR BAY
2	-122.24	37.85	52	1467	190.0	496	177	7.2574	352100	NEAR BAY
3	-122.25	37.85	52	1274	235.0	558	219	5.6431	341300	NEAR BAY
4	-122.25	37.85	52	1627	280.0	565	259	3.8462	342200	NEAR BAY
5	-122.25	37.85	52	919	213.0	413	193	4.0368	269700	NEAR BAY
6	-122.25	37.84	52	2535	489.0	1094	514	3.6591	299200	NEAR BAY
7	-122.25	37.84	52	3104	687.0	1157	647	3.1200	241400	NEAR BAY
8	-122.26	37.84	42	2555	665.0	1206	595	2.0804	226700	NEAR BAY
9	-122.25	37.84	52	3549	707.0	1551	714	3.6912	261100	NEAR BAY

```
In [5]: df.shape #finding total number of rows and Columns...
```

```
Out[5]: (20640, 10)
```

```
In [6]: df.info() #Complete information about the data...
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20640 entries, 0 to 20639
Data columns (total 10 columns):
 #   Column              Non-Null Count  Dtype  
---  --
 0   longitude            20640 non-null  float64
 1   latitude             20640 non-null  float64
 2   housing_median_age   20640 non-null  int64  
 3   total_rooms          20640 non-null  int64  
 4   total_bedrooms       20433 non-null  float64
 5   population            20640 non-null  int64  
 6   households            20640 non-null  int64  
 7   median_income         20640 non-null  float64
 8   median_house_value   20640 non-null  int64  
 9   ocean_proximity      20640 non-null  object  
dtypes: float64(4), int64(5), object(1)
memory usage: 1.6+ MB
```

```
In [7]: df.describe() #describing statistics...
```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value
count	20640.000000	20640.000000	20640.000000	20640.000000	20433.000000	20640.000000	20640.000000	20640.000000	20640.000000
mean	-119.569704	35.631861	28.639486	2635.763061	537.870553	1425.476744	499.539660	3.870671	206855.816909
std	2.003532	2.135952	12.585558	2181.615252	421.385070	1132.462122	382.329753	1.899822	115395.615874
min	-124.350000	32.540000	1.000000	2.000000	1.000000	3.000000	1.000000	0.499900	14999.000000
25%	-121.800000	33.930000	18.000000	1447.750000	296.000000	787.000000	280.000000	2.563400	119600.000000
50%	-118.490000	34.260000	29.000000	2127.000000	435.000000	1166.000000	409.000000	3.548000	179700.000000
75%	-118.000000	37.710000	37.000000	3148.000000	647.000000	1725.000000	605.000000	4.743250	264725.000000
max	-114.310000	41.950000	52.000000	39320.000000	6445.000000	35682.000000	6082.000000	15.000100	500001.000000

```
In [8]: df.isna().sum() #By this function, we can find a Null values in a dataset...
```

longitude	0
latitude	0
housing_median_age	0
total_rooms	0
total_bedrooms	207
population	0
households	0
median_income	0
median_house_value	0
ocean_proximity	0
dtype: int64	

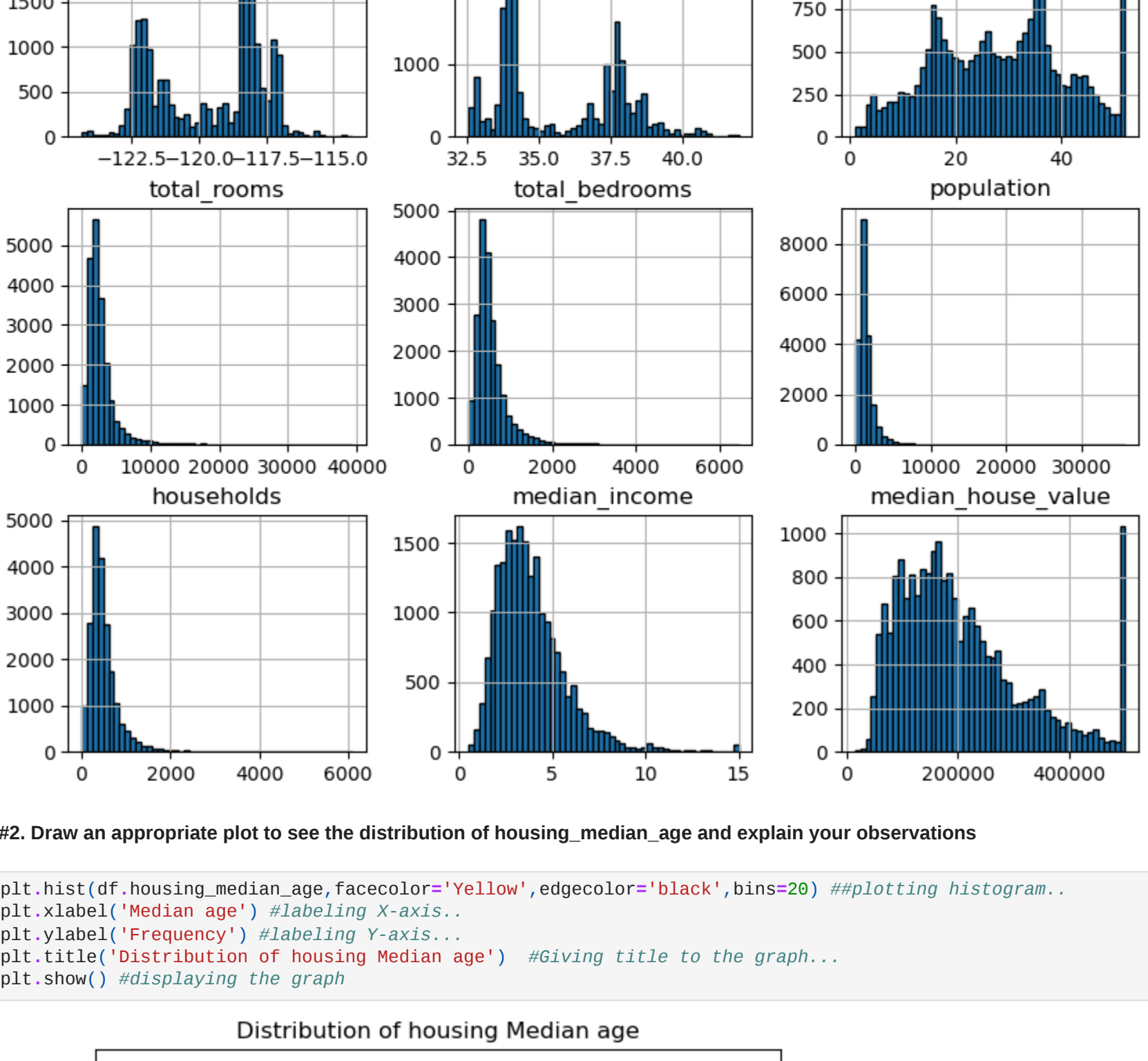
#1. What is the average median income of the data set and check the distribution of data using appropriate plots. Please explain the distribution of the plot

```
In [9]: df['median_income'].mean() #Finding average median income from the given Dataset...
```

```
Out[9]: 3.8706710029069766
```

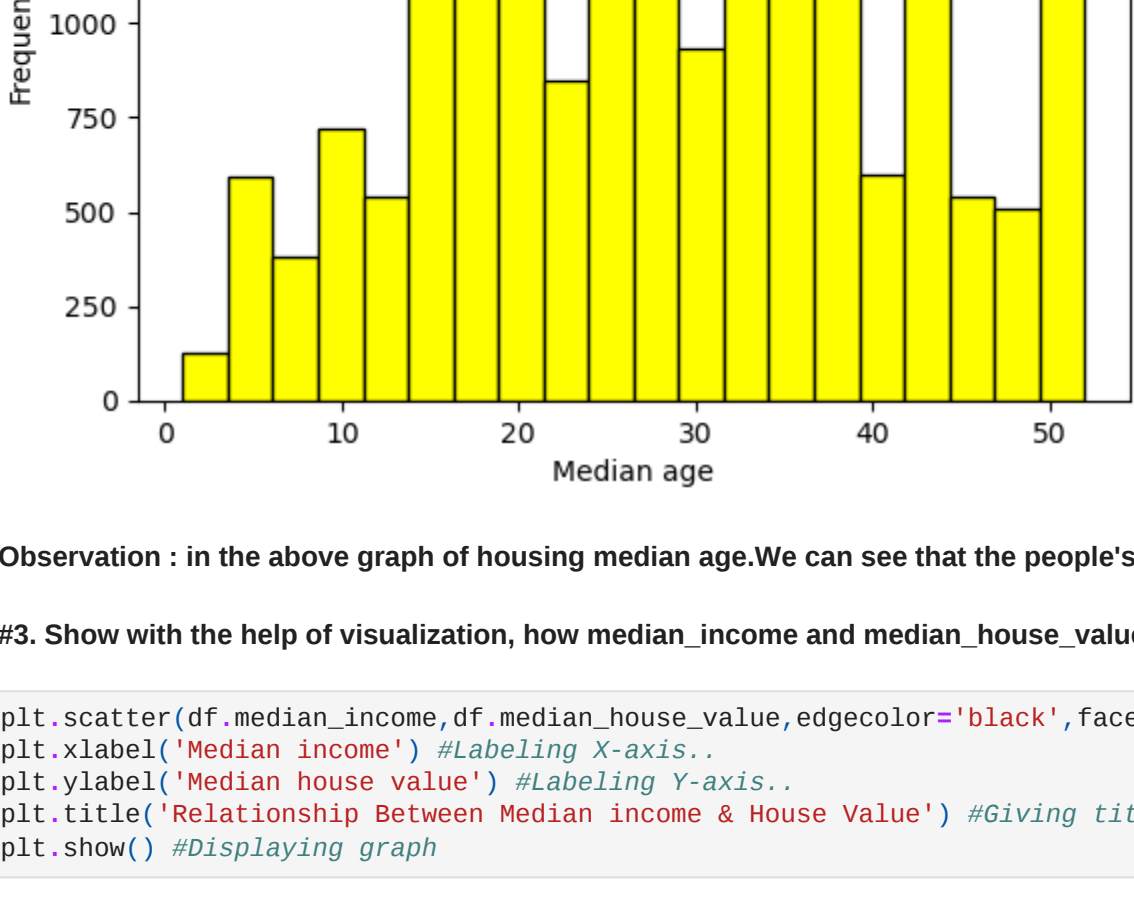
```
In [10]: df.hist(edgecolor='black',bins=50,figsize=(10,8)) #creating the histogram plots for the dataset
```

```
array([[<Axes: title='center': 'longitude'>],
       [<Axes: title='center': 'latitude'>],
       [<Axes: title='center': 'housing_median_age'>],
       [<Axes: title='center': 'total_rooms'>],
       [<Axes: title='center': 'total_bedrooms'>],
       [<Axes: title='center': 'population'>],
       [<Axes: title='center': 'households'>],
       [<Axes: title='center': 'median_income'>],
       [<Axes: title='center': 'median_house_value'>]], dtype=object)
```



#2. Draw an appropriate plot to see the distribution of housing_median_age and explain your observations

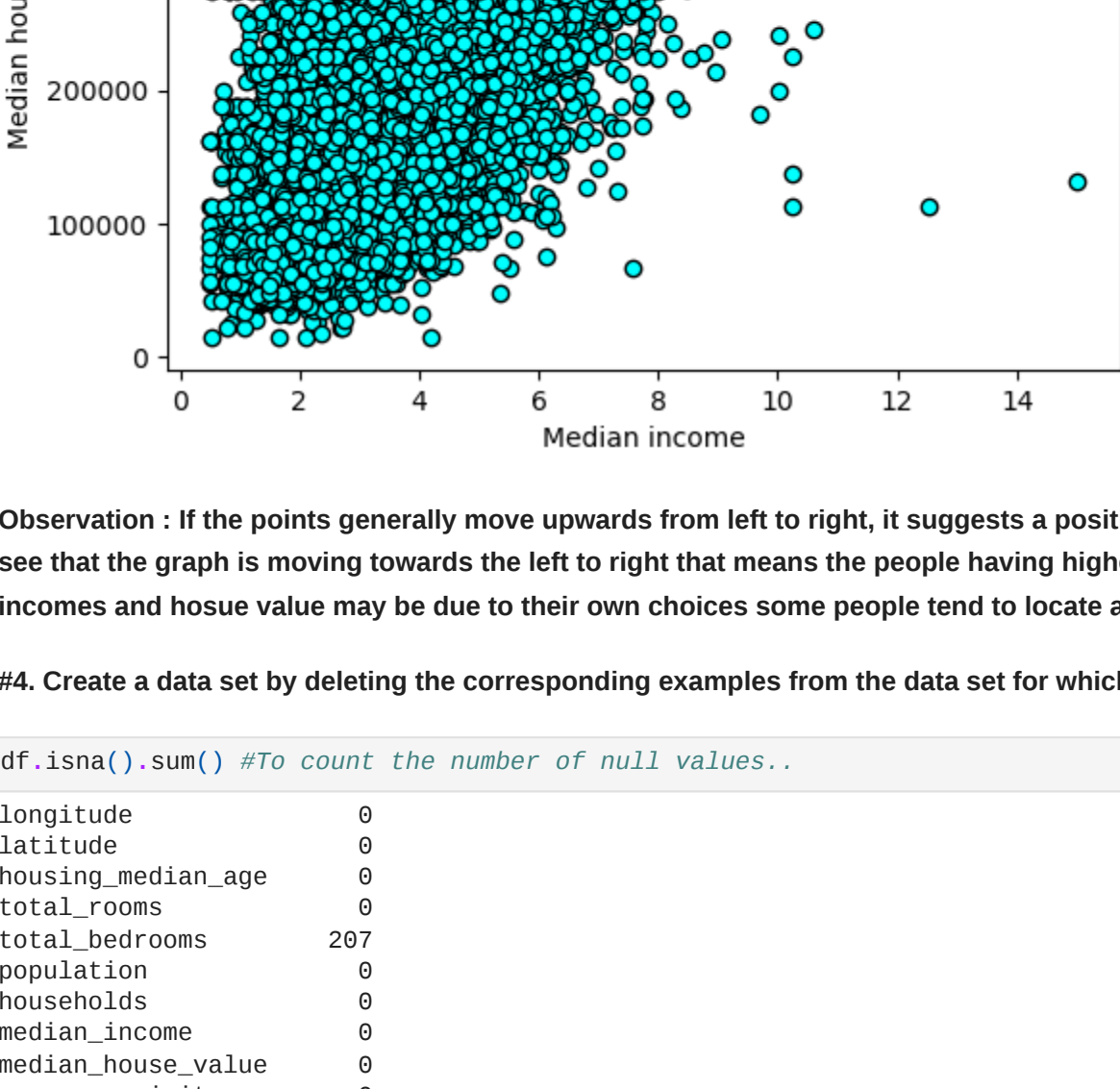
```
In [11]: plt.hist(df.housing_median_age,facecolor='Yellow',edgecolor='black',bins=20) ##plotting histogram...
plt.xlabel('Median age') #labeling x-axis..
plt.ylabel('Frequency') #labeling Y-axis...
plt.title('Distribution of housing Median age') #Giving title to the graph...
plt.show() #displaying the graph
```



Observation : in the above graph of housing median age.We can see that the people's who are in the age between 30-40 they have more income

#3. Show with the help of visualization, how median_income and median_house_values are related

```
In [12]: plt.scatter(df.median_income,df.median_house_value,edgecolor='black',facecolor='cyan') #Plotting Scatter plot
plt.xlabel('Median income') #labeling x-axis..
plt.ylabel('Median house value') #labeling Y-axis..
plt.title('Relationship Between Median income & House Value') #Giving title to graph
plt.show() #displaying graph
```



Observation : If the points generally move upwards from left to right, it suggests a positive correlation: areas with higher median incomes tend to have higher median house values.in our scenario as we can see that the graph is moving towards the left to right that means the people having higher incomes are moving towards to higher pricing houses.somewhat its showing a slight difference inbetween median incomes and hosue value may be due to their own choices some people tend to locate at the same place.overall there is a positive correlation we can say

#4. Create a data set by deleting the corresponding examples from the data set for which total_bedrooms are not available

```
In [13]: df.isna().sum() #To count the number of null values...
```

longitude	0
latitude	0
housing_median_age	0
total_rooms	0
total_bedrooms	207
population	0
households	0
median_income	0
median_house_value	0
ocean_proximity	0
dtype: int64	

```
In [14]: df1=df.dropna() #it will delete all the null rows...
```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
0	-122.23	37.88	41	880	129.0	322	126	8.3252	452600	NEAR BAY
1	-122.22	37.86	21	7099	1106.0	2401	1138	8.3014	358500	NEAR BAY
2	-122.24	37.85	52	1467	190.0	496	177	7.2574	352100	NEAR BAY
3	-122.25	37.85	52	1274	235.0	558	219	5.6431	341300	NEAR BAY
4	-122.25	37.85	52	1627	280.0	565	259	3.8462	342200	NEAR BAY
...
20635	-121.09	39.48	25	1665	374.0	845	330	1.5603	78100	INLAND
20636	-121.21	39.49	18	697	150.0	356	114	2.5568	77100	INLAND
20637	-121.22	39.43	17	2254	485.0	1007	433	1.7000	92300	INLAND
20638	-121.32	39.43	18	1860	409.0	741	349	1.8672	84700	INLAND
20639	-121.24	39.37	16	2785	616.0	1387	530	2.3886	89400	INLAND

20433 rows × 10 columns

```
In [15]: df.dropna(subset=['total_bedrooms']) #Creation of New Data after removing Null Values from the data..
df2 = New data without Null values...
```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
0	-122.23	37.88	41	880	129.0	322	126	8.3252	452600	NEAR BAY
1	-122.22	37.86	21	7099	1106.0	2401	1138	8.3014	358500	NEAR BAY
2	-122.24	37.85	52	1467	190.0	496	177	7.2574	352100	NEAR BAY
3	-122.25	37.85	52	1274	235.0	558	219	5.6431	341300	NEAR BAY
4	-122.25	37.85	52	1627	280.0	565	259	3.8462	342200	NEAR BAY
...
20635	-121.09	39.48	25	1665	374.0	845	330	1.5603	78100	INLAND
20636	-121.21	39.49	18	697	150.0	356	114	2.5568	77100	INLAND
20637	-121.22	39.43	17	2254	485.0	1007	433	1.7000	92300	INLAND
20638	-121.32	39.43	18	1860	409.0	741	349	1.8672	84700	INLAND
20639	-121.24	39.37	16	2785	616.0	1387	530	2.3886	89400	INLAND

20433 rows × 10 columns

#5. Create a data set by filling the missing data with the mean value of the total_bedrooms in the original data set

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
0	-122.23	37.88	41	880	129.0	322	126	8.3252	452600	NEAR BAY
1	-122.22	37.86	21	7099	1106.0	2401	1138	8.3014	358500	NEAR BAY
2	-122.24	37.85	52	1467	190.0	496	177	7.2574	352100	NEAR BAY
3	-122.25	37.85	52	1274	235.0	558	219	5.6431	341300	NEAR BAY
4	-122.25	37.85	52	1627	280.0	565	259	3.8462	342200	NEAR BAY
...
20635	-121.09	39.48	25	1665	374.0	845	330	1.5603	78100	INLAND
20636	-121.21	39.49	18	697	150.0	356	114	2.5568	77100	INLAND
20637	-121.22	39.43	17	2254	485.0	1007	433	1.7000	92300	INLAND
20638	-121.32	39.43	18	1860	409.0	741	349	1.8672	84700	INLAND
20639	-121.24	39.37	16	2785	616.0	1387	530	2.3886	89400	INLAND

20640 rows × 10 columns

```
In [17]: df2=df.fillna(value=df.mean()) #fillna function is used to fill the empty values of null..
df2
```

C:\Users\santo\AppData\Local\Temp\ipykernel_324\3772388735.py:1: FutureWarning: The median value of numeric_only in DataFrame.fillna is deprecated. In a future version, it will default to False. In addition, specifying 'numeric_only=None' is deprecated. Select only valid columns or specify the value of numeric_only to silence this warning.

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
0	-122.23	37.88	41	880	129.0	322	126	8.3252	452600	NEAR BAY
1	-122.22	37.86	21	7099	1106.0	2401	1138	8.3014	358500	NEAR BAY
2	-122.24	37.85	52	1467	190.0	496	177	7.2574	352100	NEAR BAY
3	-122.25	37.85	52	1274	235.0	558	219	5.6431	341300	NEAR BAY
4	-122.25	37.85	52	1627	280.0	565	259	3.8462	342200	NEAR BAY
...
20635	-121.09	39.48	25	1665	374.0	845	330	1.5603	78100	INLAND
20636	-121.21	39.49	18	697	150.0	356	114	2.5568	77100	INLAND
20637	-121.22	39.43	17	2254	485.0	1007	433	1.7000	92300	INLAND
20638	-121.32	39.43	18	1860	409.0	741	349	1.8672	84700	INLAND
20639	-121.24	39.37	16	2785	616.0	1387	530	2.3886	89400	INLAND

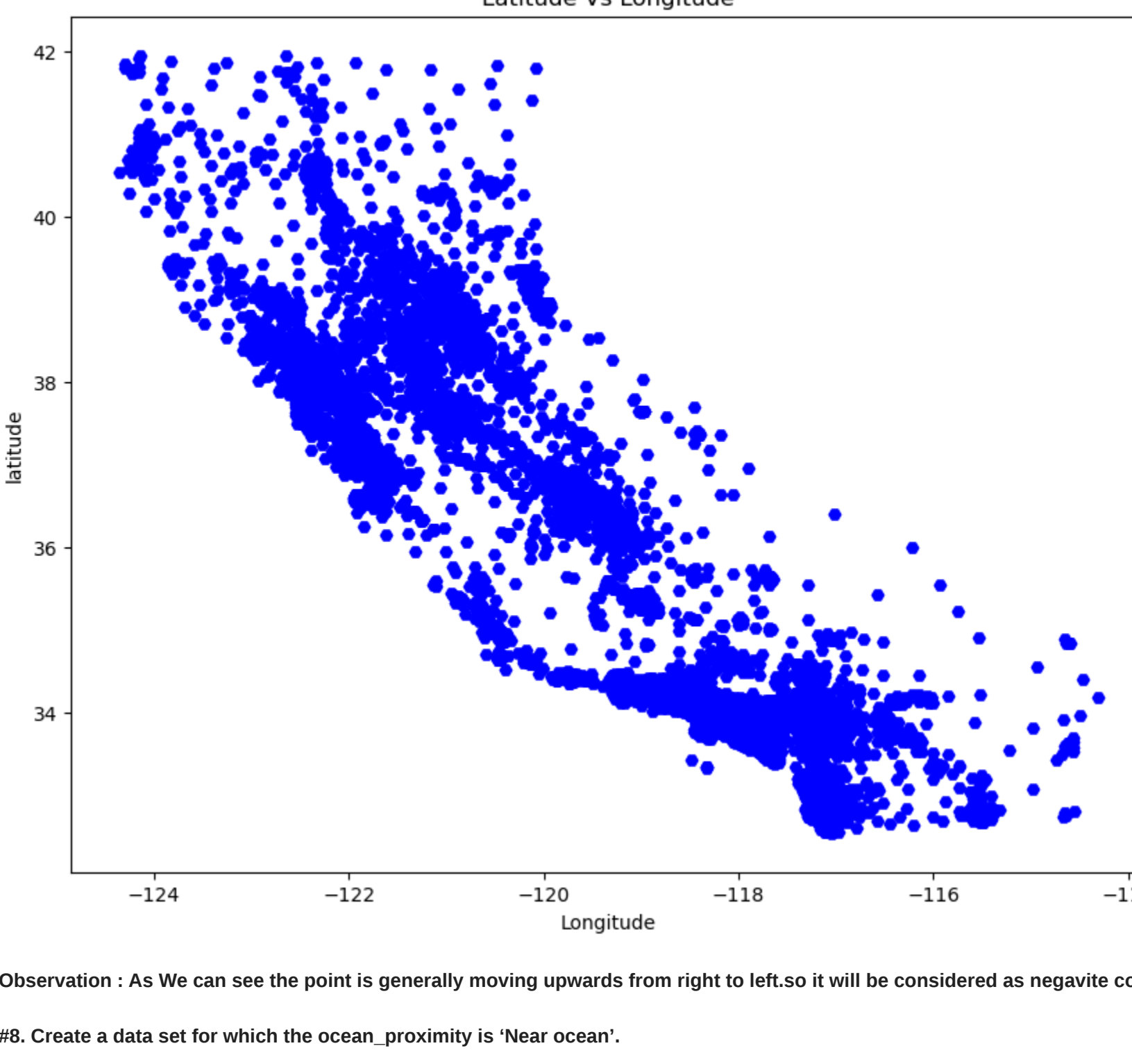
20640 rows × 10 columns

```
Total_Median_Bedrooms: 438.0
Total_median_population: 1266.0
Total median of house value: 179700.0
```

#7. Plot latitude versus longitude and explain your observations

```
In [23]: #Extracting latitude and longitude...
latitude=df2["latitude"]
longitude=df2["longitude"]
```

```
In [27]: #creating a scatter plot to see the relationship between both..
plt.figure(figsize=(10,8))
plt.scatter(longitude,latitude,facecolor='blue',marker='H')
plt.xlabel('longitude')
plt.ylabel('latitude')
plt.title('Latitude Vs Longitude')
plt.show()
```



Observation : As We can see the point is generally moving upwards from right to left.so it will be considered as negativ correlation.

#8. Create a data set for which the ocean_proximity is 'Near ocean'.

```
In [32]: filtered_data=df2[df2["ocean_proximity"]=="NEAR OCEAN"] #By using Filter function we can find the exact values of necessary data..
filtered_data
```

	longitude	
--	-----------	--