# Lab 2-Demonstrate the steps to build a machine-learning model that predicts the median housing price using the California housing price dataset.

Download the dataset :
https://media.geeksforgeeks.org/wp-content/uploads/20240319120216/housing.csv
1. Perform the describe and info steps

```python
#1.Import and head
import pandas as pd
url = "https://media.geeksforgeeks.org/wp-content/uploads/20240319120216/housing.csv"
housing = pd.read_csv(url)

housing.head()
```

| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | median_house_value | ocean_proximity |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | -122.23 | 37.88 | 41.0 | 880.0 | 129.0 | 322.0 | 126.0 | 8.3252 | 452600.0 | NEAR BAY |
| 1 | -122.22 | 37.86 | 21.0 | 7099.0 | 1106.0 | 2401.0 | 1138.0 | 8.3014 | 358500.0 | NEAR BAY |
| 2 | -122.24 | 37.85 | 52.0 | 1467.0 | 190.0 | 496.0 | 177.0 | 7.2574 | 352100.0 | NEAR BAY |
| 3 | -122.25 | 37.85 | 52.0 | 1274.0 | 235.0 | 558.0 | 219.0 | 5.6431 | 341300.0 | NEAR BAY |
| 4 | -122.25 | 37.85 | 52.0 | 1627.0 | 280.0 | 565.0 | 259.0 | 3.8462 | 342200.0 | NEAR BAY |

```python
#Info
housing.info()
```

```
<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  float64
 3   total_rooms         20640 non-null  float64
 4   total_bedrooms      20433 non-null  float64
 5   population          20640 non-null  float64
 6   households          20640 non-null  float64
 7   median_income       20640 non-null  float64
 8   median_house_value  20640 non-null  float64
 9   ocean_proximity     20640 non-null  object
dtypes: float64(9), object(1)
memory usage: 1.6+ MB
```

```
#Describe
housing.describe()
```
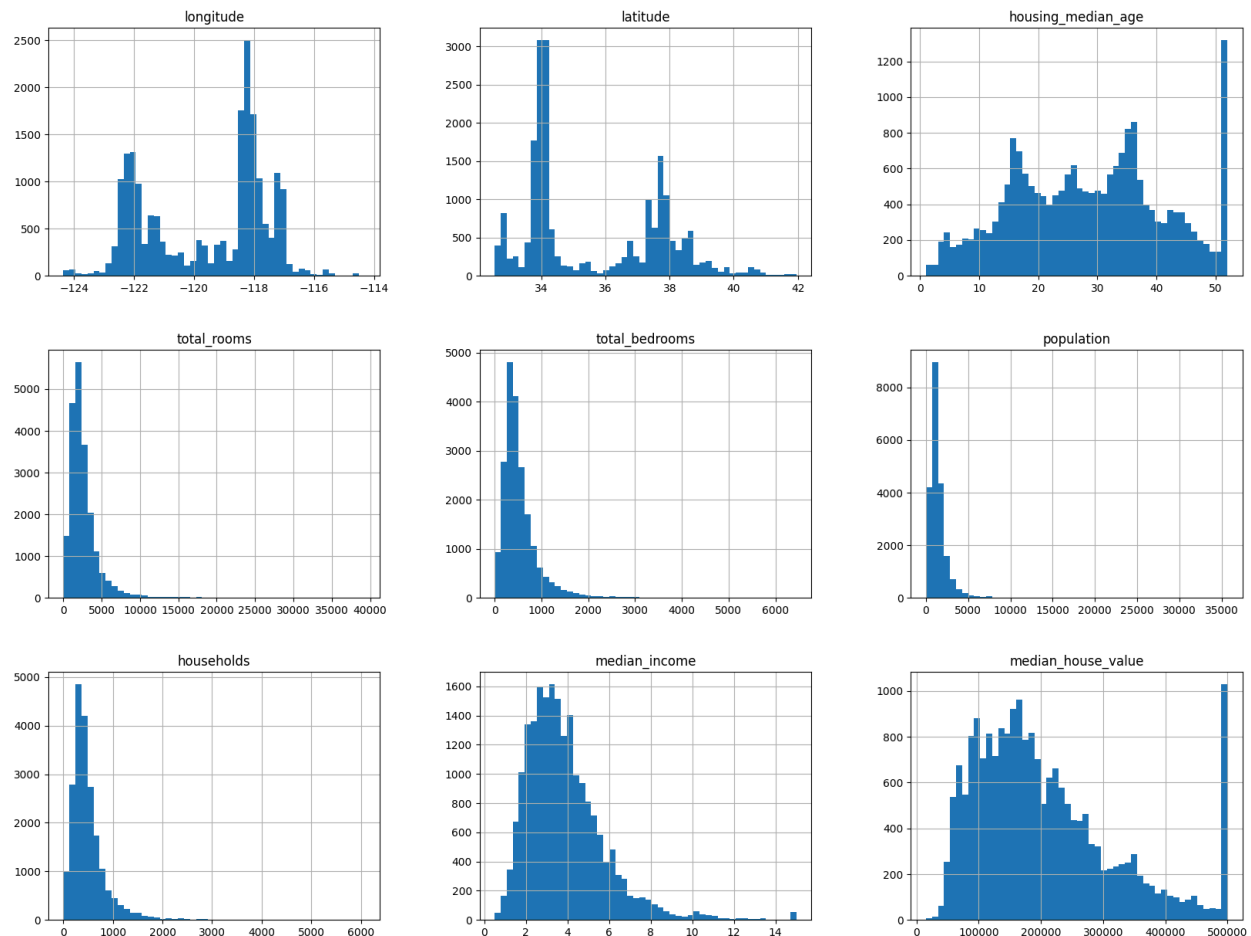
| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | median_house_value |
|---|---|---|---|---|---|---|---|---|---|
| count | 20640.000000 | 20640.000000 | 20640.000000 | 20640.000000 | 20433.000000 | 20640.000000 | 20640.000000 | 20640.000000 | 20640.000000 |
| mean | -119.569704 | 35.631861 | 28.639486 | 2635.763081 | 537.870553 | 1425.476744 | 499.539680 | 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 | 3.534800 | 179700.000000 |
| 75% | -118.010000 | 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 |

## 2. Plot the histogram of each feature( Indicate what does histogram indicate on median_income and house_median_age)

```
housing.hist(bins=50, figsize=(20,15))
plt.show()
```

## 3. Demonstrate the process of creating a test set( write the difference between random and stratified test set)

```python
#Random test set
from sklearn.model_selection import train_test_split
train_set, test_set = train_test_split(housing, test_size=0.2,
random_state=42)


#Stratified test set based on income category
import numpy as np
housing["income_cat"] = pd.cut(
    housing["median_income"],
    bins=[0., 1.5, 3.0, 4.5, 6., np.inf],
    labels=[1, 2, 3, 4, 5]
)
from sklearn.model_selection import StratifiedShuffleSplit
split = StratifiedShuffleSplit(n_splits=1, test_size=0.2, random_state=42)
for train_index, test_index in split.split(housing,
housing["income_cat"]):
    strat_train_set = housing.loc[train_index]
    strat_test_set = housing.loc[test_index]
```
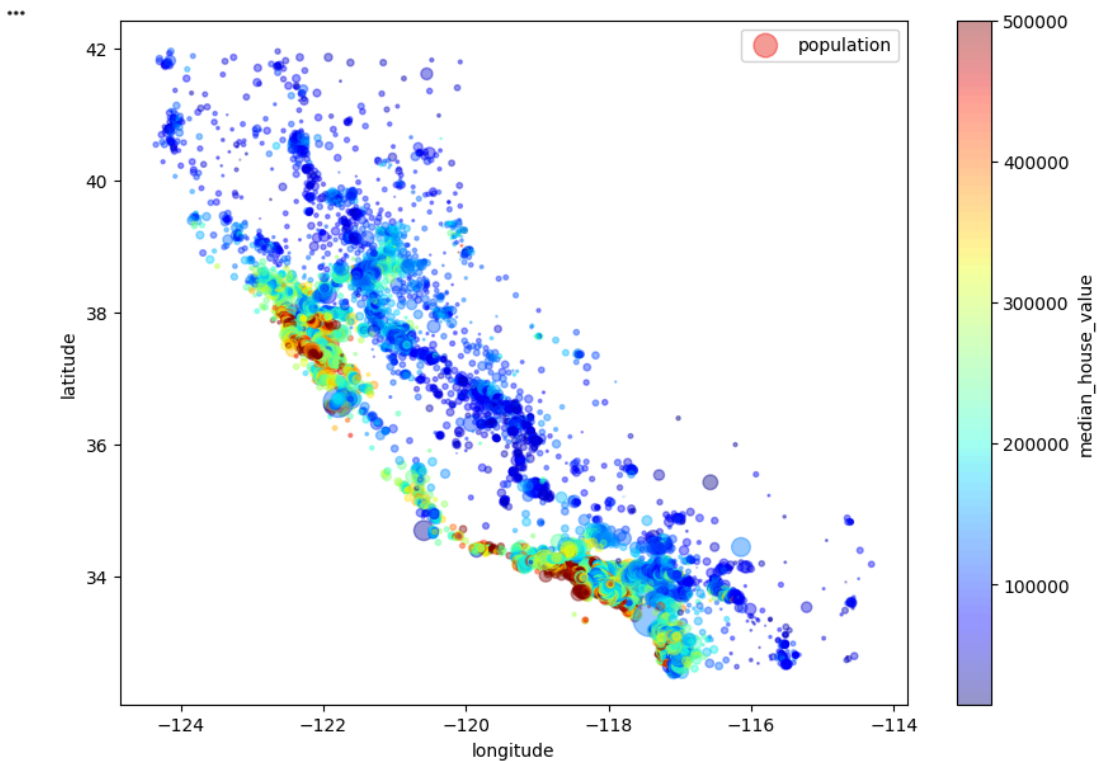
## 4. List the geographical features from the dataset and plot a graph to Visualize Geographical Data(what does the graph indicate w.r.t housing prices and location)

```python
housing.plot(kind="scatter", x="longitude", y="latitude",
            alpha=0.4,
            s=housing["population"]/100,
            label="population",
            figsize=(10,7),
            c="median_house_value",
            cmap=plt.get_cmap("jet"),
            colorbar=True)

plt.legend()
plt.show()
```

5. Plot a graph to show features correlation with housing price. Which feature corelates to the maximum. Plot the graph for that with housing price and analyze what the graph indicate.

```
#Correlation with Housing Price
corr_matrix = housing.corr(numeric_only=True)


corr_matrix["median_house_value"].sort_values(ascending=False)
```

| | median_house_value |
|---|---|
| median_house_value | 1.000000 |
| median_income | 0.688075 |
| total_rooms | 0.134153 |
| housing_median_age | 0.105623 |
| households | 0.065843 |
| total_bedrooms | 0.049686 |
| population | -0.024650 |
| longitude | -0.045967 |
| latitude | -0.144160 |

**dtype:** float64
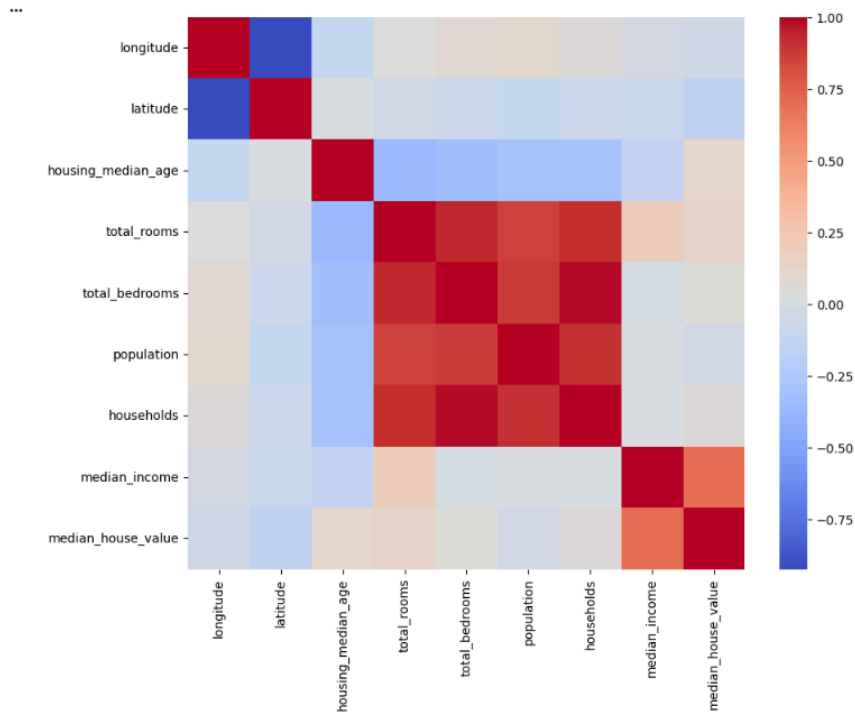
```
#Correlation graph
import seaborn as sns

plt.figure(figsize=(10,8))
sns.heatmap(corr_matrix, annot=False, cmap="coolwarm")
plt.show()
```
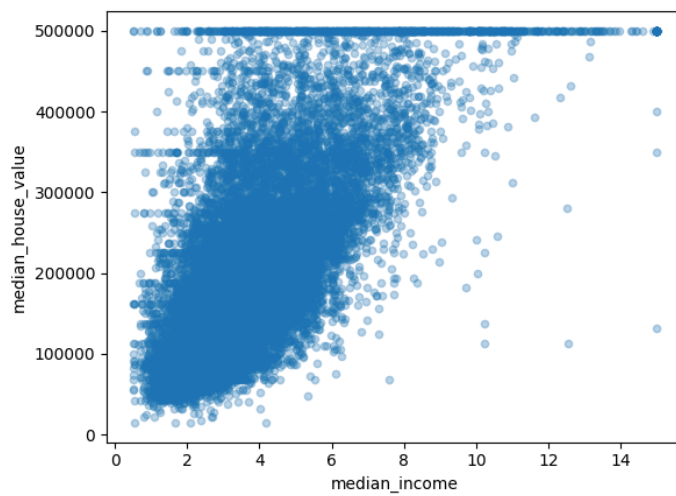


```
#Plot median_income vs price
housing.plot(kind="scatter", x="median_income", y="median_house_value",
alpha=0.3)
plt.show()
```

## 6. List the features that could be combined to improve correlation and plot again to see if correlation has improved.

```python
#6. List the features that could be combined to improve correlation and
plot again to see if correlation has improved
housing["rooms_per_household"] = housing["total_rooms"] /
housing["households"]
housing["bedrooms_per_room"] = housing["total_bedrooms"] /
housing["total_rooms"]
housing["population_per_household"] = housing["population"] /
housing["households"]
corr_matrix = housing.corr(numeric_only=True)
corr_matrix["median_house_value"].sort_values(ascending=False)
```

|  | median_house_value |
|---|---|
| median_house_value | 1.000000 |
| median_income | 0.688075 |
| rooms_per_household | 0.151948 |
| total_rooms | 0.134153 |
| housing_median_age | 0.105623 |
| households | 0.065843 |
| total_bedrooms | 0.049686 |
| population_per_household | -0.023737 |
| population | -0.024650 |
| longitude | -0.045967 |
| latitude | -0.144160 |
| bedrooms_per_room | -0.255880 |

dtype: float64

## 7. List the features that needs to be cleaned and demonstrate the process of cleaning.

```python
housing.dropna(subset=["total_bedrooms"])
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy="median")
housing_num = housing.drop("ocean_proximity", axis=1)
imputer.fit(housing_num)
housing_num_imputed = imputer.transform(housing_num)
```

## 8. Is there any categorical data that needs to be converted to numerical? If so explain the method used to convert and code the same and show the output.

```python
from sklearn.preprocessing import OneHotEncoder
housing_cat = housing[["ocean_proximity"]]
encoder = OneHotEncoder()
housing_cat_1hot = encoder.fit_transform(housing_cat)
housing_cat_1hot.toarray()
```

```
...    array([[0., 0., 0., 1., 0.],
              [0., 0., 0., 1., 0.],
              [0., 0., 0., 1., 0.],
              ...,
              [0., 1., 0., 0., 0.],
              [0., 1., 0., 0., 0.],
              [0., 1., 0., 0., 0.]])
```

## 9. Discuss the importance of feature scaling.

```python
#Scaling
from sklearn.preprocessing import StandardScaler


scaler = StandardScaler()
housing_scaled = scaler.fit_transform(housing_num)
```

## 10. Design a pipeline inculcating (Custom transform, feature scaling and encoding). Explain how it works

```python
#Custom Transformer
from sklearn.base import BaseEstimator, TransformerMixin

class CombinedAttributesAdder(BaseEstimator, TransformerMixin):
    def __init__(self, add_bedrooms_per_room=True):
        self.add_bedrooms_per_room = add_bedrooms_per_room

    def fit(self, X, y=None):
        return self

    def transform(self, X):
        rooms_per_household = X[:, 3] / X[:, 6]
        population_per_household = X[:, 5] / X[:, 6]
```

```python
        if self.add_bedrooms_per_room:
            bedrooms_per_room = X[:, 4] / X[:, 3]
            return np.c_[X, rooms_per_household,
                         population_per_household,
                         bedrooms_per_room]
        else:
            return np.c_[X, rooms_per_household,
                         population_per_household]
#Full Pipeline
from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer

num_attribs = list(housing.drop("ocean_proximity", axis=1))
cat_attribs = ["ocean_proximity"]

num_pipeline = Pipeline([
    ('imputer', SimpleImputer(strategy="median")),
    ('attribs_adder', CombinedAttributesAdder()),
    ('scaler', StandardScaler()),
])

full_pipeline = ColumnTransformer([
    ("num", num_pipeline, num_attribs),
    ("cat", OneHotEncoder(), cat_attribs),
])

housing_prepared = full_pipeline.fit_transform(housing)
```