# 18IT140_PRACTICAL_12_House_Price_Prediction_with_Python

October 29, 2021

## 1 Machine Learning Project on House Price Prediction with Python

```python
[1]: import pandas as pd
     housing = pd.read_csv("/content/sample_data/housing.csv")
     housing.head()
```

```
[1]:    longitude  latitude  ...  median_house_value  ocean_proximity
     0   -122.23     37.88   ...            452600.0         NEAR BAY
     1   -122.22     37.86   ...            358500.0         NEAR BAY
     2   -122.24     37.85   ...            352100.0         NEAR BAY
     3   -122.25     37.85   ...            341300.0         NEAR BAY
     4   -122.25     37.85   ...            342200.0         NEAR BAY

     [5 rows x 10 columns]
```
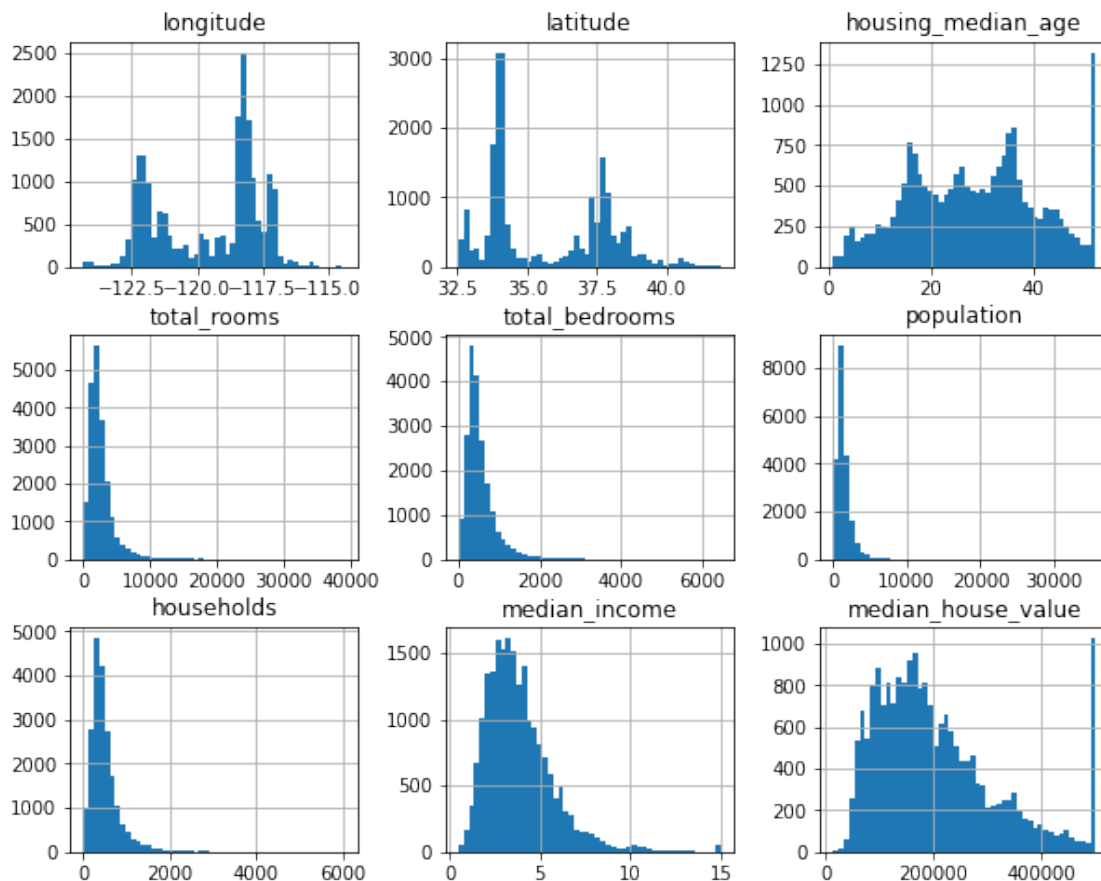
```python
[2]: housing.info()
```

```
<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
```

```
[3]: housing.ocean_proximity.value_counts()
```
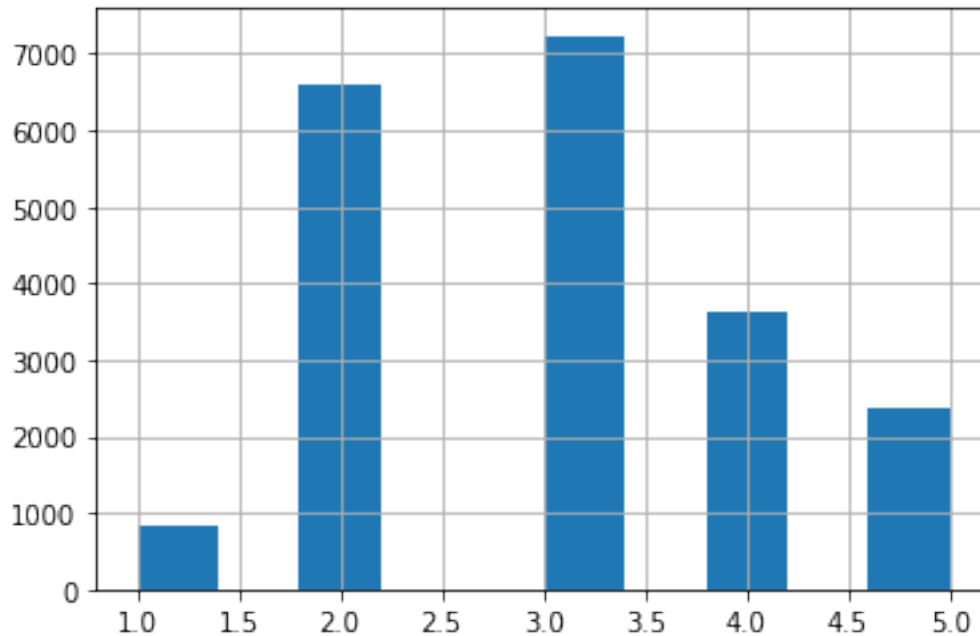
```
[3]: <1H OCEAN      9136
     INLAND        6551
     NEAR OCEAN    2658
     NEAR BAY      2290
     ISLAND           5
     Name: ocean_proximity, dtype: int64
```

```
[4]: import matplotlib.pyplot as plt
     housing.hist(bins=50, figsize=(10, 8))
     plt.show()
```



```
[5]: from sklearn.model_selection import train_test_split
     train_set, test_set = train_test_split(housing, test_size=0.2, random_state=42)
```

```
[6]: 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])
     housing['income_cat'].hist()
     plt.show()
```
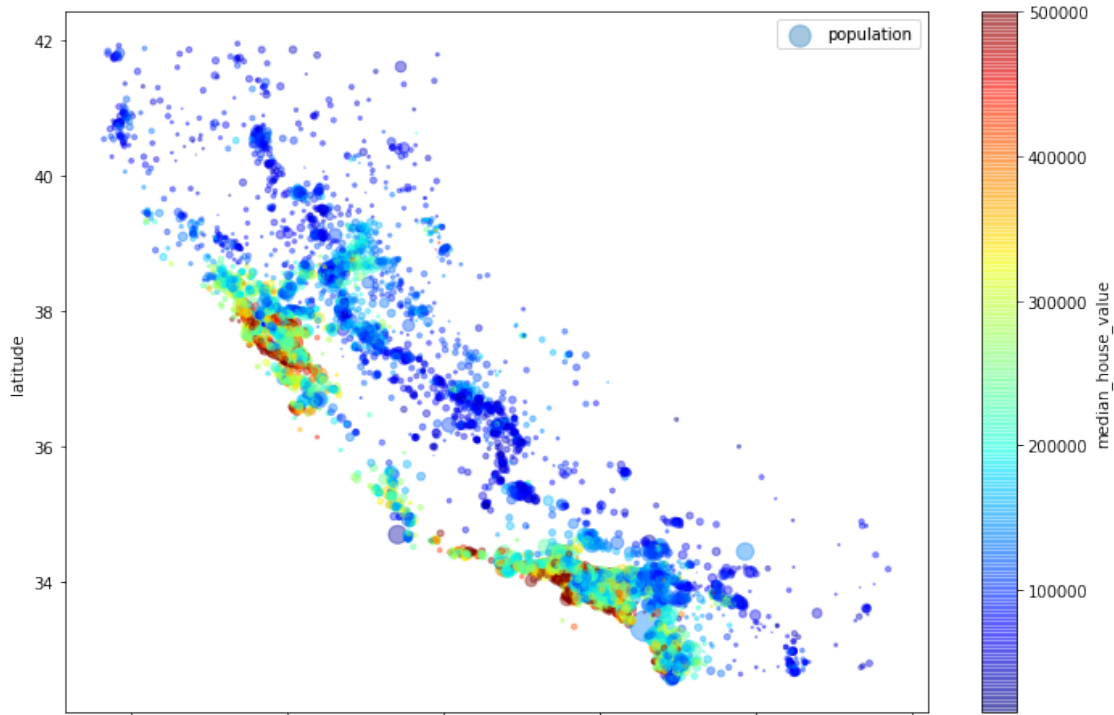
```
[7]: from sklearn.model_selection import StratifiedShuffleSplit
     split = StratifiedShuffleSplit(n_splits=1, test_size=0.2, random_state=42)
     for train_index, test_index in split.split(housing, housing["income_cat"]):
         strat_train_set = housing.loc[train_index]
         strat_test_set = housing.loc[test_index]
     print(strat_test_set['income_cat'].value_counts() / len(strat_test_set))
```

```
3    0.350533
2    0.318798
4    0.176357
5    0.114583
1    0.039729
Name: income_cat, dtype: float64
```

```
[8]: for set_ in (strat_train_set, strat_test_set):
         set_.drop('income_cat', axis=1, inplace=True)
     housing = strat_train_set.copy()
```

```
[9]: housing.plot(kind='scatter', x='longitude', y='latitude', alpha=0.4,␣
       ↪s=housing['population']/100, label='population',
     figsize=(12, 8), c='median_house_value', cmap=plt.get_cmap('jet'),␣
       ↪colorbar=True)
     plt.legend()
     plt.show()
```

```
[10]: corr_matrix = housing.corr()
      print(corr_matrix.median_house_value.sort_values(ascending=False))
```

```
median_house_value      1.000000
median_income           0.687160
total_rooms             0.135097
housing_median_age      0.114110
households              0.064506
total_bedrooms          0.047689
population              -0.026920
longitude               -0.047432
latitude                -0.142724
Name: median_house_value, dtype: float64
```

```
[11]: 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()
      print(corr_matrix["median_house_value"].sort_values(ascending=False))
```

```
median_house_value        1.000000
median_income             0.687160
```

```
rooms_per_household          0.146285
total_rooms                  0.135097
housing_median_age           0.114110
households                   0.064506
total_bedrooms               0.047689
population_per_household     -0.021985
population                   -0.026920
longitude                    -0.047432
latitude                     -0.142724
bedrooms_per_room            -0.259984
Name: median_house_value, dtype: float64
```

[12]:
```python
# Data Preparation
housing = strat_train_set.drop("median_house_value", axis=1)
housing_labels = strat_train_set["median_house_value"].copy()

median = housing["total_bedrooms"].median()
housing["total_bedrooms"].fillna(median, inplace=True)

housing_num = housing.drop("ocean_proximity", axis=1)

from sklearn.base import BaseEstimator, TransformerMixin

# column index
rooms_ix, bedrooms_ix, population_ix, households_ix = 3, 4, 5, 6

class CombinedAttributesAdder(BaseEstimator, TransformerMixin):
    def __init__(self, add_bedrooms_per_room=True): # no *args or **kargs
        self.add_bedrooms_per_room = add_bedrooms_per_room
    def fit(self, X, y=None):
        return self  # nothing else to do
    def transform(self, X):
        rooms_per_household = X[:, rooms_ix] / X[:, households_ix]
        population_per_household = X[:, population_ix] / X[:, households_ix]
        if self.add_bedrooms_per_room:
            bedrooms_per_room = X[:, bedrooms_ix] / X[:, rooms_ix]
            return np.c_[X, rooms_per_household, population_per_household,
                         bedrooms_per_room]
        else:
            return np.c_[X, rooms_per_household, population_per_household]
```

[13]:
```python
from sklearn.preprocessing import OneHotEncoder
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.impute import SimpleImputer
num_pipeline = Pipeline([
    ('imputer',SimpleImputer(strategy="median")),
    ('attribs_adder', CombinedAttributesAdder()),
```

```python
    ('std_scaler', StandardScaler()),
])
housing_num_tr = num_pipeline.fit_transform(housing_num)

from sklearn.compose import ColumnTransformer
num_attribs = list(housing_num)
cat_attribs = ["ocean_proximity"]
full_pipeline = ColumnTransformer([
    ("num", num_pipeline, num_attribs),
    ("cat", OneHotEncoder(), cat_attribs),
])
housing_prepared = full_pipeline.fit_transform(housing)
```

## 2 Linear Regression for House Price Prediction with Python

```python
[14]: from sklearn.linear_model import LinearRegression
lin_reg = LinearRegression()
lin_reg.fit(housing_prepared, housing_labels)

data = housing.iloc[:5]
labels = housing_labels.iloc[:5]
data_preparation = full_pipeline.transform(data)
print("Predictions: ", lin_reg.predict(data_preparation))
```

```
Predictions:  [210644.60459286 317768.80697211 210956.43331178  59218.98886849
 189747.55849879]
```

I hope you liked this article on Machine Learning project on House Price Prediction with Python.

```python
[ ]: !wget -nc https://raw.githubusercontent.com/brpy/colab-pdf/master/colab_pdf.py
from colab_pdf import colab_pdf
colab_pdf('18IT140_PRACTICAL_12_House_Price_Prediction_with_Python.ipynb')
```

```
--2021-10-29 15:10:39--  https://raw.githubusercontent.com/brpy/colab-
pdf/master/colab_pdf.py
Resolving raw.githubusercontent.com (raw.githubusercontent.com)...
185.199.108.133, 185.199.110.133, 185.199.111.133, ...
Connecting to raw.githubusercontent.com
(raw.githubusercontent.com)|185.199.108.133|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 1864 (1.8K) [text/plain]
Saving to: colab_pdf.py

colab_pdf.py          100%[===================>]   1.82K  --.-KB/s    in 0s
```

2021-10-29 15:10:39 (40.1 MB/s) - colab_pdf.py saved [1864/1864]


WARNING: apt does not have a stable CLI interface. Use with caution in scripts.


WARNING: apt does not have a stable CLI interface. Use with caution in scripts.

Extracting templates from packages: 100%