```
In [1]:   # Python ≥3.5 is required
          import sys
          assert sys.version_info >= (3, 5)

          # Scikit-Learn ≥0.20 is required
          import sklearn
          assert sklearn.__version__ >= "0.20"

          # Common imports
          import numpy as np
          import os

          # To plot pretty figures
          %matplotlib inline
          import matplotlib as mpl
          import matplotlib.pyplot as plt
          mpl.rc('axes', labelsize=14)
          mpl.rc('xtick', labelsize=12)
          mpl.rc('ytick', labelsize=12)

          # Where to save the figures
          PROJECT_ROOT_DIR = "."
          CHAPTER_ID = "end_to_end_project"
          IMAGES_PATH = os.path.join(PROJECT_ROOT_DIR, "images", CHAPTER_ID)
          os.makedirs(IMAGES_PATH, exist_ok=True)

          def save_fig(fig_id, tight_layout=True, fig_extension="png", resolution=300):
              path = os.path.join(IMAGES_PATH, fig_id + "." + fig_extension)
              print("Saving figure", fig_id)
              if tight_layout:
                  plt.tight_layout()
              plt.savefig(path, format=fig_extension, dpi=resolution)

          # Ignore useless warnings (see SciPy issue #5998)
          import warnings
          warnings.filterwarnings(action="ignore", message="^internal gelsd")
```

```python
In [2]: import os
        import tarfile
        import urllib

        DOWNLOAD_ROOT = "https://raw.githubusercontent.com/ageron/handson-ml2/master/"
        HOUSING_PATH = os.path.join("datasets", "housing")
        HOUSING_URL = DOWNLOAD_ROOT + "datasets/housing/housing.tgz"

        def fetch_housing_data(housing_url=HOUSING_URL, housing_path=HOUSING_PATH):
            if not os.path.isdir(housing_path):
                os.makedirs(housing_path)
            tgz_path = os.path.join(housing_path, "housing.tgz")
            urllib.request.urlretrieve(housing_url, tgz_path)
            housing_tgz = tarfile.open(tgz_path)
            housing_tgz.extractall(path=housing_path)
            housing_tgz.close()
```

```python
In [3]: fetch_housing_data()
```

```python
In [4]: import pandas as pd

        def load_housing_data(housing_path=HOUSING_PATH):
            csv_path = os.path.join(housing_path, "housing.csv")
            return pd.read_csv(csv_path)
```

```python
In [5]: housing = load_housing_data()
        housing.head()
```

Out[5]:

| | 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 |

```
In [6]: housing.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20640 entries, 0 to 20639
Data columns (total 10 columns):
longitude             20640 non-null float64
latitude              20640 non-null float64
housing_median_age    20640 non-null float64
total_rooms           20640 non-null float64
total_bedrooms        20433 non-null float64
population            20640 non-null float64
households            20640 non-null float64
median_income         20640 non-null float64
median_house_value    20640 non-null float64
ocean_proximity       20640 non-null object
dtypes: float64(9), object(1)
memory usage: 1.6+ MB
```

```
In [7]: housing["ocean_proximity"].value_counts()
```

```
Out[7]: <1H OCEAN      9136
        INLAND         6551
        NEAR OCEAN     2658
        NEAR BAY       2290
        ISLAND            5
        Name: ocean_proximity, dtype: int64
```
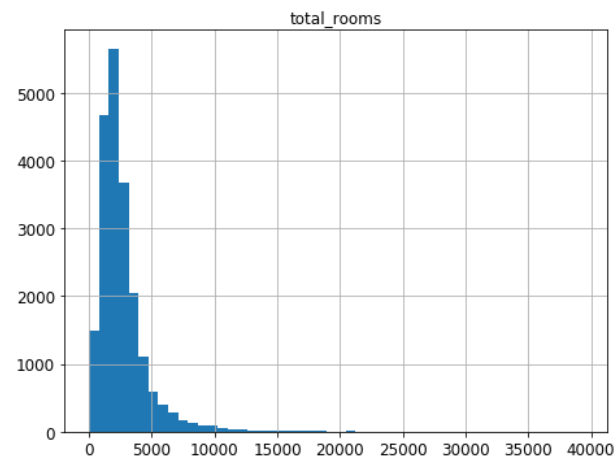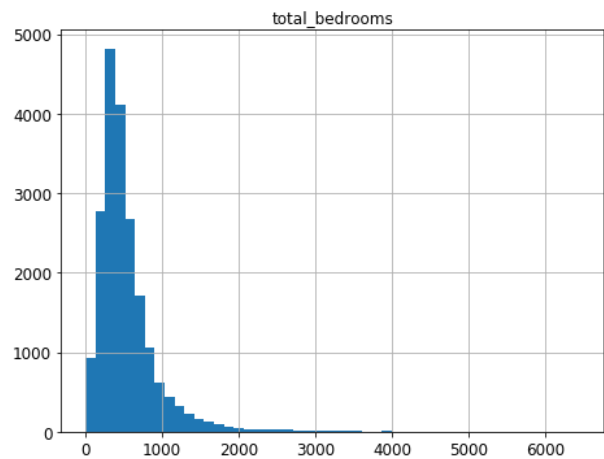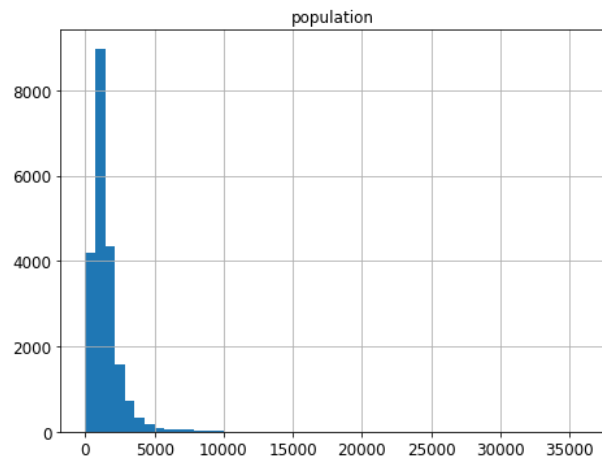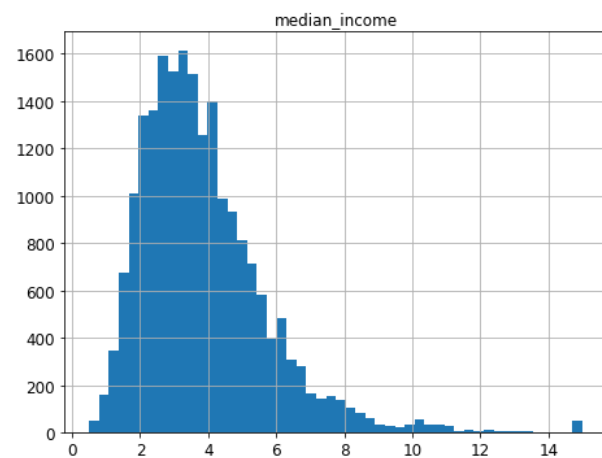
```
In [8]: housing.describe()
```

Out[8]:

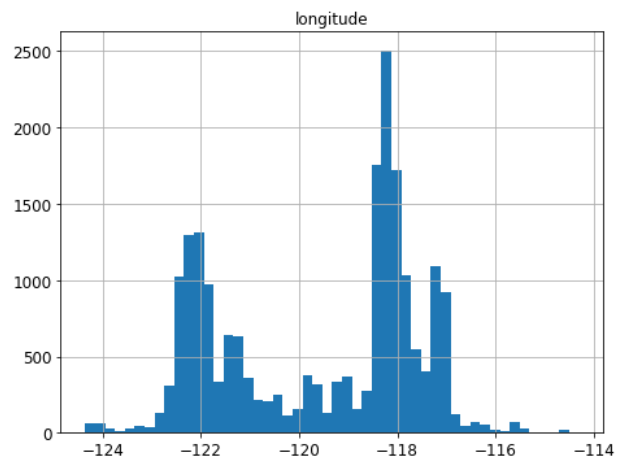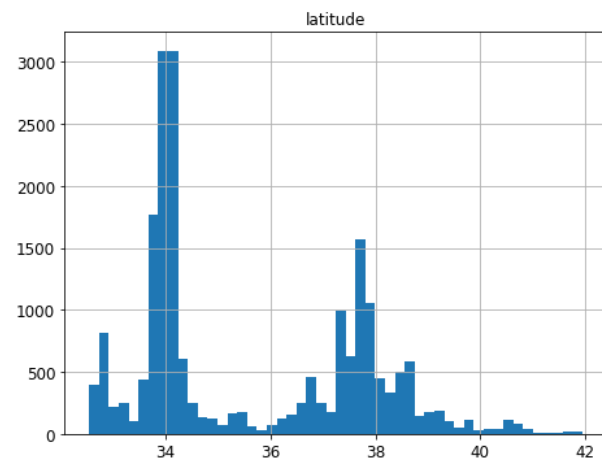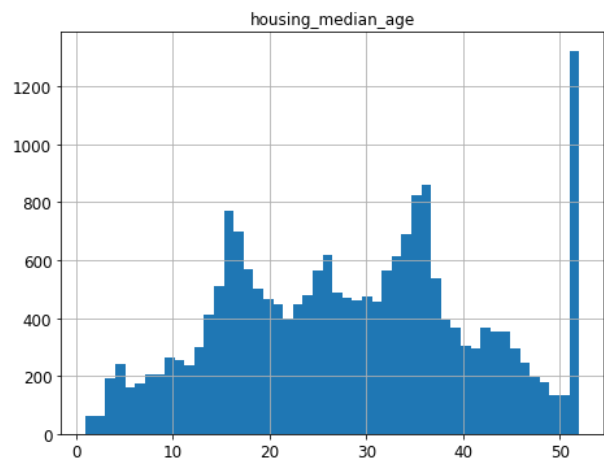| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | median_house_value |
|---|---|---|---|---|---|---|---|---|---|
| count | 20640.000000 | 20640.000000 | 20640.000000 | 20640.000000 | 20433.000000 | 20640.000000 | 20640.000000 | 20640.000000 | 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 |

```
In [9]: %matplotlib inline
        import matplotlib.pyplot as plt
        housing.hist(bins=50, figsize=(20,15))
        save_fig("attribute_histogram_plots")
        plt.show()
```

Saving figure attribute_histogram_plots

```
In [10]:   # to make this notebook's output identical at every run
           np.random.seed(42)
```

```
In [11]:   import numpy as np

           # For illustration only. Sklearn has train_test_split()
           def split_train_test(data, test_ratio):
               shuffled_indices = np.random.permutation(len(data))
               test_set_size = int(len(data) * test_ratio)
               test_indices = shuffled_indices[:test_set_size]
               train_indices = shuffled_indices[test_set_size:]
               return data.iloc[train_indices], data.iloc[test_indices]
```

```
In [12]:   train_set, test_set = split_train_test(housing, 0.2)
           len(train_set)
```

```
Out[12]:   16512
```

```
In [13]:   len(test_set)
```

```
Out[13]:   4128
```

```
In [14]:   from zlib import crc32

           def test_set_check(identifier, test_ratio):
               return crc32(np.int64(identifier)) & 0xffffffff < test_ratio * 2**32

           def split_train_test_by_id(data, test_ratio, id_column):
               ids = data[id_column]
               in_test_set = ids.apply(lambda id_: test_set_check(id_, test_ratio))
               return data.loc[~in_test_set], data.loc[in_test_set]
```

```
In [15]:   import hashlib

           def test_set_check(identifier, test_ratio, hash=hashlib.md5):
               return hash(np.int64(identifier)).digest()[-1] < 256 * test_ratio
```

```
In [16]:   def test_set_check(identifier, test_ratio, hash=hashlib.md5):
               return bytearray(hash(np.int64(identifier)).digest())[-1] < 256 * test_ratio
```

```
In [17]:   housing_with_id = housing.reset_index()   # adds an `index` column
           train_set, test_set = split_train_test_by_id(housing_with_id, 0.2, "index")
```

```
In [18]:   housing_with_id["id"] = housing["longitude"] * 1000 + housing["latitude"]
           train_set, test_set = split_train_test_by_id(housing_with_id, 0.2, "id")
```

```
In [19]: test_set.head()
```

Out[19]:

| | index | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | median_house_value | ocean_proximity | id |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **8** | 8 | -122.26 | 37.84 | 42.0 | 2555.0 | 665.0 | 1206.0 | 595.0 | 2.0804 | 226700.0 | NEAR BAY | -122222.16 |
| **10** | 10 | -122.26 | 37.85 | 52.0 | 2202.0 | 434.0 | 910.0 | 402.0 | 3.2031 | 281500.0 | NEAR BAY | -122222.15 |
| **11** | 11 | -122.26 | 37.85 | 52.0 | 3503.0 | 752.0 | 1504.0 | 734.0 | 3.2705 | 241800.0 | NEAR BAY | -122222.15 |
| **12** | 12 | -122.26 | 37.85 | 52.0 | 2491.0 | 474.0 | 1098.0 | 468.0 | 3.0750 | 213500.0 | NEAR BAY | -122222.15 |
| **13** | 13 | -122.26 | 37.84 | 52.0 | 696.0 | 191.0 | 345.0 | 174.0 | 2.6736 | 191300.0 | NEAR BAY | -122222.16 |

```
In [20]: from sklearn.model_selection import train_test_split

         train_set, test_set = train_test_split(housing, test_size=0.2, random_state=42)
```

```
In [21]: test_set.head()
```

Out[21]:

| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | median_house_value | ocean_proximity |
|---|---|---|---|---|---|---|---|---|---|---|
| **20046** | -119.01 | 36.06 | 25.0 | 1505.0 | NaN | 1392.0 | 359.0 | 1.6812 | 47700.0 | INLAND |
| **3024** | -119.46 | 35.14 | 30.0 | 2943.0 | NaN | 1565.0 | 584.0 | 2.5313 | 45800.0 | INLAND |
| **15663** | -122.44 | 37.80 | 52.0 | 3830.0 | NaN | 1310.0 | 963.0 | 3.4801 | 500001.0 | NEAR BAY |
| **20484** | -118.72 | 34.28 | 17.0 | 3051.0 | NaN | 1705.0 | 495.0 | 5.7376 | 218600.0 | <1H OCEAN |
| **9814** | -121.93 | 36.62 | 34.0 | 2351.0 | NaN | 1063.0 | 428.0 | 3.7250 | 278000.0 | NEAR OCEAN |

```
In [22]: housing["median_income"].hist()
```

Out[22]: <matplotlib.axes._subplots.AxesSubplot at 0x105702400>



```
In [23]: 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])
```

```
In [24]: housing["income_cat"].value_counts()
```

Out[24]: 3    7236
         2    6581
         4    3639
         5    2362
         1     822
         Name: income_cat, dtype: int64

```
In [25]:  housing["income_cat"].hist()
```

Out[25]:  <matplotlib.axes._subplots.AxesSubplot at 0x1057024e0>



```
In [26]:  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]
```
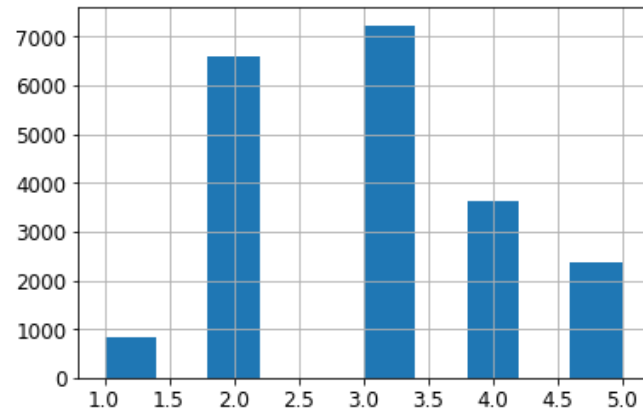
```
In [27]:  strat_test_set["income_cat"].value_counts() / len(strat_test_set)
```

Out[27]:  3    0.350533
          2    0.318798
          4    0.176357
          5    0.114583
          1    0.039729
          Name: income_cat, dtype: float64

```
In [28]:  housing["income_cat"].value_counts() / len(housing)
```

Out[28]:  3    0.350581
          2    0.318847
          4    0.176308
          5    0.114438
          1    0.039826
          Name: income_cat, dtype: float64

```
In [29]:  def income_cat_proportions(data):
              return data["income_cat"].value_counts() / len(data)

          train_set, test_set = train_test_split(housing, test_size=0.2, random_state=42)

          compare_props = pd.DataFrame({
              "Overall": income_cat_proportions(housing),
              "Stratified": income_cat_proportions(strat_test_set),
              "Random": income_cat_proportions(test_set),
          }).sort_index()
          compare_props["Rand. %error"] = 100 * compare_props["Random"] / compare_props["Overall"] - 100
          compare_props["Strat. %error"] = 100 * compare_props["Stratified"] / compare_props["Overall"] - 100
```

```
In [30]:  compare_props
```

Out[30]:

|   | Overall | Stratified | Random | Rand. %error | Strat. %error |
|---|---------|-----------|--------|--------------|---------------|
| 1 | 0.039826 | 0.039729 | 0.040213 | 0.973236 | -0.243309 |
| 2 | 0.318847 | 0.318798 | 0.324370 | 1.732260 | -0.015195 |
| 3 | 0.350581 | 0.350533 | 0.358527 | 2.266446 | -0.013820 |
| 4 | 0.176308 | 0.176357 | 0.167393 | -5.056334 | 0.027480 |
| 5 | 0.114438 | 0.114583 | 0.109496 | -4.318374 | 0.127011 |

```
In [31]:  for set_ in (strat_train_set, strat_test_set):
              set_.drop("income_cat", axis=1, inplace=True)
```

```
In [32]:  housing = strat_train_set.copy()
```

In [33]: 
```python
housing.plot(kind="scatter", x="longitude", y="latitude")
save_fig("bad_visualization_plot")
```

Saving figure bad_visualization_plot



In [34]: 
```python
housing.plot(kind="scatter", x="longitude", y="latitude", alpha=0.1)
save_fig("better_visualization_plot")
```

Saving figure better_visualization_plot

```
In [35]: 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,
             sharex=False)
         plt.legend()
         save_fig("housing_prices_scatterplot")
```

Saving figure housing_prices_scatterplot

```python
In [36]: # Download the California image
         images_path = os.path.join(PROJECT_ROOT_DIR, "images", "end_to_end_project")
         os.makedirs(images_path, exist_ok=True)
         DOWNLOAD_ROOT = "https://raw.githubusercontent.com/ageron/handson-ml2/master/"
         filename = "california.png"
         print("Downloading", filename)
         url = DOWNLOAD_ROOT + "images/end_to_end_project/" + filename
         urllib.request.urlretrieve(url, os.path.join(images_path, filename))
```

```
Downloading california.png
```

```
Out[36]: ('./images/end_to_end_project/california.png',
          <http.client.HTTPMessage at 0x1047fa2e8>)
```

```python
In [37]: import matplotlib.image as mpimg
         california_img=mpimg.imread(os.path.join(images_path, filename))
         ax = housing.plot(kind="scatter", x="longitude", y="latitude", figsize=(10,7),
                               s=housing['population']/100, label="Population",
                               c="median_house_value", cmap=plt.get_cmap("jet"),
                               colorbar=False, alpha=0.4,
                               )
         plt.imshow(california_img, extent=[-124.55, -113.80, 32.45, 42.05], alpha=0.5,
                     cmap=plt.get_cmap("jet"))
         plt.ylabel("Latitude", fontsize=14)
         plt.xlabel("Longitude", fontsize=14)

         prices = housing["median_house_value"]
         tick_values = np.linspace(prices.min(), prices.max(), 11)
         cbar = plt.colorbar()
         cbar.ax.set_yticklabels(["$%dk"%(round(v/1000)) for v in tick_values], fontsize=14)
         cbar.set_label('Median House Value', fontsize=16)

         plt.legend(fontsize=16)
         save_fig("california_housing_prices_plot")
         plt.show()
```

Saving figure california_housing_prices_plot



```
In [38]:  corr_matrix = housing.corr()
```

```
In [39]:  corr_matrix["median_house_value"].sort_values(ascending=False)
```

```
Out[39]:  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
```

```
In [40]:  from pandas.plotting import scatter_matrix

          attributes = ["median_house_value", "median_income", "total_rooms",
                        "housing_median_age"]
          scatter_matrix(housing[attributes], figsize=(12, 8))
          save_fig("scatter_matrix_plot")
```

Saving figure scatter_matrix_plot

```
In [41]: housing.plot(kind="scatter", x="median_income", y="median_house_value",
                       alpha=0.1)
         plt.axis([0, 16, 0, 550000])
         save_fig("income_vs_house_value_scatterplot")
```

Saving figure income_vs_house_value_scatterplot



```
In [42]: 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"]
```

```
In [43]: corr_matrix = housing.corr()
         corr_matrix["median_house_value"].sort_values(ascending=False)
```

```
Out[43]: 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
```

```
In [44]:  housing.plot(kind="scatter", x="rooms_per_household", y="median_house_value",
                        alpha=0.2)
          plt.axis([0, 5, 0, 520000])
          plt.show()
```



```
In [45]:  housing.describe()
```

Out[45]:

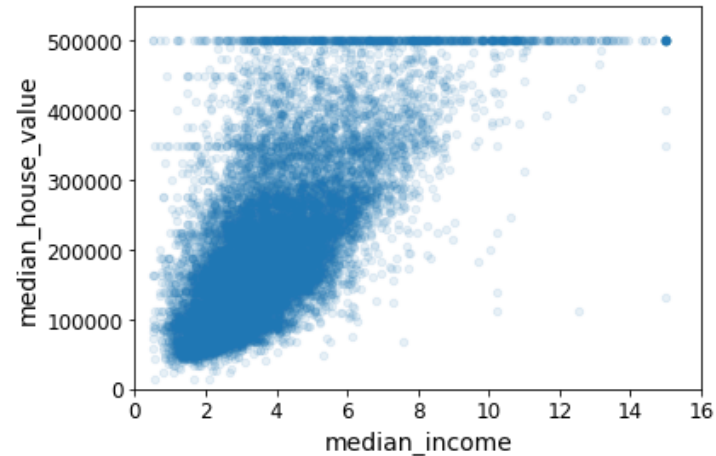| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | median_house_value | rooms_per_household |
|---|---|---|---|---|---|---|---|---|---|---|
| count | 16512.000000 | 16512.000000 | 16512.000000 | 16512.000000 | 16354.000000 | 16512.000000 | 16512.000000 | 16512.000000 | 16512.000000 | 16512.000000 |
| mean | -119.575834 | 35.639577 | 28.653101 | 2622.728319 | 534.973890 | 1419.790819 | 497.060380 | 3.875589 | 206990.920724 | 5.440341 |
| std | 2.001860 | 2.138058 | 12.574726 | 2138.458419 | 412.699041 | 1115.686241 | 375.720845 | 1.904950 | 115703.014830 | 2.611712 |
| min | -124.350000 | 32.540000 | 1.000000 | 6.000000 | 2.000000 | 3.000000 | 2.000000 | 0.499900 | 14999.000000 | 1.130435 |
| 25% | -121.800000 | 33.940000 | 18.000000 | 1443.000000 | 295.000000 | 784.000000 | 279.000000 | 2.566775 | 119800.000000 | 4.442040 |
| 50% | -118.510000 | 34.260000 | 29.000000 | 2119.500000 | 433.000000 | 1164.000000 | 408.000000 | 3.540900 | 179500.000000 | 5.232284 |
| 75% | -118.010000 | 37.720000 | 37.000000 | 3141.000000 | 644.000000 | 1719.250000 | 602.000000 | 4.744475 | 263900.000000 | 6.056361 |
| max | -114.310000 | 41.950000 | 52.000000 | 39320.000000 | 6210.000000 | 35682.000000 | 5358.000000 | 15.000100 | 500001.000000 | 141.909091 |

```
In [46]:  #Prepaing data for ML algorithm
          housing = strat_train_set.drop("median_house_value", axis=1) # drop labels for training set
          housing_labels = strat_train_set["median_house_value"].copy()
```

```
In [47]: sample_incomplete_rows = housing[housing.isnull().any(axis=1)].head()
         sample_incomplete_rows
```

Out[47]:

| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | ocean_proximity |
|---|---|---|---|---|---|---|---|---|---|
| **4629** | -118.30 | 34.07 | 18.0 | 3759.0 | NaN | 3296.0 | 1462.0 | 2.2708 | <1H OCEAN |
| **6068** | -117.86 | 34.01 | 16.0 | 4632.0 | NaN | 3038.0 | 727.0 | 5.1762 | <1H OCEAN |
| **17923** | -121.97 | 37.35 | 30.0 | 1955.0 | NaN | 999.0 | 386.0 | 4.6328 | <1H OCEAN |
| **13656** | -117.30 | 34.05 | 6.0 | 2155.0 | NaN | 1039.0 | 391.0 | 1.6675 | INLAND |
| **19252** | -122.79 | 38.48 | 7.0 | 6837.0 | NaN | 3468.0 | 1405.0 | 3.1662 | <1H OCEAN |

```
In [48]: sample_incomplete_rows.dropna(subset=["total_bedrooms"])     # option 1, ignoring rows with NA's
```

Out[48]:

| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | ocean_proximity |
|---|---|---|---|---|---|---|---|---|---|

```
In [49]: sample_incomplete_rows.drop("total_bedrooms", axis=1)     # option 2, dropping the whole column
```

Out[49]:

| | longitude | latitude | housing_median_age | total_rooms | population | households | median_income | ocean_proximity |
|---|---|---|---|---|---|---|---|---|
| **4629** | -118.30 | 34.07 | 18.0 | 3759.0 | 3296.0 | 1462.0 | 2.2708 | <1H OCEAN |
| **6068** | -117.86 | 34.01 | 16.0 | 4632.0 | 3038.0 | 727.0 | 5.1762 | <1H OCEAN |
| **17923** | -121.97 | 37.35 | 30.0 | 1955.0 | 999.0 | 386.0 | 4.6328 | <1H OCEAN |
| **13656** | -117.30 | 34.05 | 6.0 | 2155.0 | 1039.0 | 391.0 | 1.6675 | INLAND |
| **19252** | -122.79 | 38.48 | 7.0 | 6837.0 | 3468.0 | 1405.0 | 3.1662 | <1H OCEAN |

```
In [50]: median = housing["total_bedrooms"].median()
         sample_incomplete_rows["total_bedrooms"].fillna(median, inplace=True) # option 3, replacing NA's with mean value
```

```
In [51]: sample_incomplete_rows
```

Out[51]:

|  | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | ocean_proximity |
|---|---|---|---|---|---|---|---|---|---|
| **4629** | -118.30 | 34.07 | 18.0 | 3759.0 | 433.0 | 3296.0 | 1462.0 | 2.2708 | <1H OCEAN |
| **6068** | -117.86 | 34.01 | 16.0 | 4632.0 | 433.0 | 3038.0 | 727.0 | 5.1762 | <1H OCEAN |
| **17923** | -121.97 | 37.35 | 30.0 | 1955.0 | 433.0 | 999.0 | 386.0 | 4.6328 | <1H OCEAN |
| **13656** | -117.30 | 34.05 | 6.0 | 2155.0 | 433.0 | 1039.0 | 391.0 | 1.6675 | INLAND |
| **19252** | -122.79 | 38.48 | 7.0 | 6837.0 | 433.0 | 3468.0 | 1405.0 | 3.1662 | <1H OCEAN |

```
In [52]: from sklearn.impute import SimpleImputer    # Scikit's method to deal with NA's
         imputer = SimpleImputer(strategy="median")
```

```
In [53]: housing_num = housing.drop("ocean_proximity", axis=1)#Dropping the text column as above method works only on numeric data
         # alternatively: housing_num = housing.select_dtypes(include=[np.number])
```

```
In [54]: imputer.fit(housing_num)
```

Out[54]: SimpleImputer(copy=True, fill_value=None, missing_values=nan,
           strategy='median', verbose=0)

```
In [55]: imputer.statistics_
```

Out[55]: array([-118.51 ,   34.26 ,   29.   , 2119.5  ,  433.   , 1164.   ,
           408.   ,    3.5409])

```
In [56]: housing_num.median().values #chcking above method manually
```

Out[56]: array([-118.51 ,   34.26 ,   29.   , 2119.5  ,  433.   , 1164.   ,
           408.   ,    3.5409])

```
In [57]: X = imputer.transform(housing_num)
```

```
In [58]: housing_tr = pd.DataFrame(X, columns=housing_num.columns,
                         index=housing.index)
```

In [59]: `housing_tr.loc[sample_incomplete_rows.index.values]`

Out[59]:

| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income |
|---|---|---|---|---|---|---|---|---|
| **4629** | -118.30 | 34.07 | 18.0 | 3759.0 | 433.0 | 3296.0 | 1462.0 | 2.2708 |
| **6068** | -117.86 | 34.01 | 16.0 | 4632.0 | 433.0 | 3038.0 | 727.0 | 5.1762 |
| **17923** | -121.97 | 37.35 | 30.0 | 1955.0 | 433.0 | 999.0 | 386.0 | 4.6328 |
| **13656** | -117.30 | 34.05 | 6.0 | 2155.0 | 433.0 | 1039.0 | 391.0 | 1.6675 |
| **19252** | -122.79 | 38.48 | 7.0 | 6837.0 | 433.0 | 3468.0 | 1405.0 | 3.1662 |

In [60]: `imputer.strategy`

Out[60]: `'median'`

In [61]: `housing_tr = pd.DataFrame(X, columns=housing_num.columns,`
`                          index=housing_num.index)`

In [62]: `housing_tr.head()`

Out[62]:

| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income |
|---|---|---|---|---|---|---|---|---|
| **17606** | -121.89 | 37.29 | 38.0 | 1568.0 | 351.0 | 710.0 | 339.0 | 2.7042 |
| **18632** | -121.93 | 37.05 | 14.0 | 679.0 | 108.0 | 306.0 | 113.0 | 6.4214 |
| **14650** | -117.20 | 32.77 | 31.0 | 1952.0 | 471.0 | 936.0 | 462.0 | 2.8621 |
| **3230** | -119.61 | 36.31 | 25.0 | 1847.0 | 371.0 | 1460.0 | 353.0 | 1.8839 |
| **3555** | -118.59 | 34.23 | 17.0 | 6592.0 | 1525.0 | 4459.0 | 1463.0 | 3.0347 |

```
In [63]: housing_cat = housing[["ocean_proximity"]]
         housing_cat.head(10)
```

Out[63]:

|       | ocean_proximity |
|-------|-----------------|
| 17606 | <1H OCEAN       |
| 18632 | <1H OCEAN       |
| 14650 | NEAR OCEAN      |
| 3230  | INLAND          |
| 3555  | <1H OCEAN       |
| 19480 | INLAND          |
| 8879  | <1H OCEAN       |
| 13685 | INLAND          |
| 4937  | <1H OCEAN       |
| 4861  | <1H OCEAN       |

```
In [64]: from sklearn.preprocessing import OrdinalEncoder

         ordinal_encoder = OrdinalEncoder()
         housing_cat_encoded = ordinal_encoder.fit_transform(housing_cat)
         housing_cat_encoded[:10]
```

```
Out[64]: array([[0.],
                [0.],
                [4.],
                [1.],
                [0.],
                [1.],
                [0.],
                [1.],
                [0.],
                [0.]])
```

```
In [65]: ordinal_encoder.categories_
```

```
Out[65]: [array(['<1H OCEAN', 'INLAND', 'ISLAND', 'NEAR BAY', 'NEAR OCEAN'],
                dtype=object)]
```

```python
In [66]: from sklearn.preprocessing import OneHotEncoder

         cat_encoder = OneHotEncoder()
         housing_cat_1hot = cat_encoder.fit_transform(housing_cat)
         housing_cat_1hot
```

Out[66]: <16512x5 sparse matrix of type '<class 'numpy.float64'>'
             with 16512 stored elements in Compressed Sparse Row format>

```python
In [67]: housing_cat_1hot.toarray()
```

Out[67]: array([[1., 0., 0., 0., 0.],
                [1., 0., 0., 0., 0.],
                [0., 0., 0., 0., 1.],
                ...,
                [0., 1., 0., 0., 0.],
                [1., 0., 0., 0., 0.],
                [0., 0., 0., 1., 0.]])

```python
In [68]: cat_encoder = OneHotEncoder(sparse=False)
         housing_cat_1hot = cat_encoder.fit_transform(housing_cat)
         housing_cat_1hot
```

Out[68]: array([[1., 0., 0., 0., 0.],
                [1., 0., 0., 0., 0.],
                [0., 0., 0., 0., 1.],
                ...,
                [0., 1., 0., 0., 0.],
                [1., 0., 0., 0., 0.],
                [0., 0., 0., 1., 0.]])

```python
In [69]: cat_encoder.categories_
```

Out[69]: [array(['<1H OCEAN', 'INLAND', 'ISLAND', 'NEAR BAY', 'NEAR OCEAN'],
                dtype=object)]
```

```
In [70]: 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]

         attr_adder = CombinedAttributesAdder(add_bedrooms_per_room=False)
         housing_extra_attribs = attr_adder.transform(housing.values)
```

```
In [71]: housing_extra_attribs = pd.DataFrame(
             housing_extra_attribs,
             columns=list(housing.columns)+["rooms_per_household", "population_per_household"],
             index=housing.index)
         housing_extra_attribs.head()
```

Out[71]:

| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | ocean_proximity | rooms_per_household | population_per_ho |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 17606 | -121.89 | 37.29 | 38 | 1568 | 351 | 710 | 339 | 2.7042 | <1H OCEAN | 4.62537 | |
| 18632 | -121.93 | 37.05 | 14 | 679 | 108 | 306 | 113 | 6.4214 | <1H OCEAN | 6.00885 | |
| 14650 | -117.2 | 32.77 | 31 | 1952 | 471 | 936 | 462 | 2.8621 | NEAR OCEAN | 4.22511 | |
| 3230 | -119.61 | 36.31 | 25 | 1847 | 371 | 1460 | 353 | 1.8839 | INLAND | 5.23229 | |
| 3555 | -118.59 | 34.23 | 17 | 6592 | 1525 | 4459 | 1463 | 3.0347 | <1H OCEAN | 4.50581 | |

```python
In [72]:  from sklearn.pipeline import Pipeline
          from sklearn.preprocessing import StandardScaler

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

          housing_num_tr = num_pipeline.fit_transform(housing_num)
```

```python
In [73]:  housing_num_tr
```

```
Out[73]:  array([[-1.15604281,  0.77194962,  0.74333089, ..., -0.31205452,
                  -0.08649871,  0.15531753],
                 [-1.17602483,  0.6596948 , -1.1653172 , ...,  0.21768338,
                  -0.03353391, -0.83628902],
                 [ 1.18684903, -1.34218285,  0.18664186, ..., -0.46531516,
                  -0.09240499,  0.4222004 ],
                 ...,
                 [ 1.58648943, -0.72478134, -1.56295222, ...,  0.3469342 ,
                  -0.03055414, -0.52177644],
                 [ 0.78221312, -0.85106801,  0.18664186, ...,  0.02499488,
                   0.06150916, -0.30340741],
                 [-1.43579109,  0.99645926,  1.85670895, ..., -0.22852947,
                  -0.09586294,  0.10180567]])
```

```python
In [74]:  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)
```

```
In [75]: housing_prepared
```

```
Out[75]: array([[-1.15604281,  0.77194962,  0.74333089, ...,  0.        ,
                  0.        ,  0.        ],
                [-1.17602483,  0.6596948 , -1.1653172 , ...,  0.        ,
                  0.        ,  0.        ],
                [ 1.18684903, -1.34218285,  0.18664186, ...,  0.        ,
                  0.        ,  1.        ],
                ...,
                [ 1.58648943, -0.72478134, -1.56295222, ...,  0.        ,
                  0.        ,  0.        ],
                [ 0.78221312, -0.85106801,  0.18664186, ...,  0.        ,
                  0.        ,  0.        ],
                [-1.43579109,  0.99645926,  1.85670895, ...,  0.        ,
                  1.        ,  0.        ]])
```

```
In [76]: housing_prepared.shape
```

```
Out[76]: (16512, 16)
```

```
In [77]: from sklearn.linear_model import LinearRegression

         lin_reg = LinearRegression()
         lin_reg.fit(housing_prepared, housing_labels)
```

```
Out[77]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None,
                  normalize=False)
```

```
In [78]: # let's try the full preprocessing pipeline on a few training instances
         some_data = housing.iloc[:5]
         some_labels = housing_labels.iloc[:5]
         some_data_prepared = full_pipeline.transform(some_data)

         print("Predictions:", lin_reg.predict(some_data_prepared))
```

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

```
In [79]: print("Labels:", list(some_labels))
```

```
         Labels: [286600.0, 340600.0, 196900.0, 46300.0, 254500.0]
```

```
In [80]: some_data_prepared
```

```
Out[80]: array([[-1.15604281,  0.77194962,  0.74333089, -0.49323393, -0.44543821,
                  -0.63621141, -0.42069842, -0.61493744, -0.31205452, -0.08649871,
                   0.15531753,  1.        ,  0.        ,  0.        ,  0.        ,
                   0.        ],
                 [-1.17602483,  0.6596948 , -1.1653172 , -0.90896655, -1.0369278 ,
                  -0.99833135, -1.02222705,  1.33645936,  0.21768338, -0.03353391,
                  -0.83628902,  1.        ,  0.        ,  0.        ,  0.        ,
                   0.        ],
                 [ 1.18684903, -1.34218285,  0.18664186, -0.31365989, -0.15334458,
                  -0.43363936, -0.0933178 , -0.5320456 , -0.46531516, -0.09240499,
                   0.4222004 ,  0.        ,  0.        ,  0.        ,  0.        ,
                   1.        ],
                 [-0.01706767,  0.31357576, -0.29052016, -0.36276217, -0.39675594,
                   0.03604096, -0.38343559, -1.04556555, -0.07966124,  0.08973561,
                  -0.19645314,  0.        ,  1.        ,  0.        ,  0.        ,
                   0.        ],
                 [ 0.49247384, -0.65929936, -0.92673619,  1.85619316,  2.41221109,
                   2.72415407,  2.57097492, -0.44143679, -0.35783383, -0.00419445,
                   0.2699277 ,  1.        ,  0.        ,  0.        ,  0.        ,
                   0.        ]])
```

```
In [81]: from sklearn.metrics import mean_squared_error

         housing_predictions = lin_reg.predict(housing_prepared)
         lin_mse = mean_squared_error(housing_labels, housing_predictions)
         lin_rmse = np.sqrt(lin_mse)
         lin_rmse
```

```
Out[81]: 68628.19819848923
```

```
In [82]: from sklearn.metrics import mean_absolute_error

         lin_mae = mean_absolute_error(housing_labels, housing_predictions)
         lin_mae
```

```
Out[82]: 49439.89599001897
```

```
In [83]: housing_predictions
```

```
Out[83]: array([210644.60459286, 317768.80697211, 210956.43331178, ...,
                 95464.57062437, 214353.22541713, 276426.4692067 ])
```

```
In [84]:  from sklearn.tree import DecisionTreeRegressor

          tree_reg = DecisionTreeRegressor(random_state=42)
          tree_reg.fit(housing_prepared, housing_labels)
```

```
Out[84]:  DecisionTreeRegressor(criterion='mse', max_depth=None, max_features=None,
                     max_leaf_nodes=None, min_impurity_decrease=0.0,
                     min_impurity_split=None, min_samples_leaf=1,
                     min_samples_split=2, min_weight_fraction_leaf=0.0,
                     presort=False, random_state=42, splitter='best')
```

```
In [85]:  housing_predictions = tree_reg.predict(housing_prepared)
          tree_mse = mean_squared_error(housing_labels, housing_predictions)
          tree_rmse = np.sqrt(tree_mse)
          tree_rmse
```

```
Out[85]:  0.0
```

```
In [86]:  from sklearn.model_selection import cross_val_score

          scores = cross_val_score(tree_reg, housing_prepared, housing_labels,
                                scoring="neg_mean_squared_error", cv=10)
          tree_rmse_scores = np.sqrt(-scores)
```

```
In [87]:  def display_scores(scores):
              print("Scores:", scores)
              print("Mean:", scores.mean())
              print("Standard deviation:", scores.std())

          display_scores(tree_rmse_scores)
```

```
          Scores: [70194.33680785 66855.16363941 72432.58244769 70758.73896782
           71115.88230639 75585.14172901 70262.86139133 70273.6325285
           75366.87952553 71231.65726027]
          Mean: 71407.68766037929
          Standard deviation: 2439.4345041191004
```

```
In [88]: lin_scores = cross_val_score(lin_reg, housing_prepared, housing_labels,
                                       scoring="neg_mean_squared_error", cv=10)
         lin_rmse_scores = np.sqrt(-lin_scores)
         display_scores(lin_rmse_scores)
```

```
Scores: [66782.73843989 66960.118071   70347.95244419 74739.57052552
 68031.13388938 71193.84183426 64969.63056405 68281.61137997
 71552.91566558 67665.10082067]
Mean: 69052.46136345083
Standard deviation: 2731.674001798349
```

```
In [89]: from sklearn.ensemble import RandomForestRegressor

         forest_reg = RandomForestRegressor(n_estimators=100, random_state=42)
         forest_reg.fit(housing_prepared, housing_labels)
```

```
Out[89]: RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
                   max_features='auto', max_leaf_nodes=None,
                   min_impurity_decrease=0.0, min_impurity_split=None,
                   min_samples_leaf=1, min_samples_split=2,
                   min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=None,
                   oob_score=False, random_state=42, verbose=0, warm_start=False)
```

```
In [90]: housing_predictions = forest_reg.predict(housing_prepared)
         forest_mse = mean_squared_error(housing_labels, housing_predictions)
         forest_rmse = np.sqrt(forest_mse)
         forest_rmse
```

```
Out[90]: 18603.515021376355
```

```
In [94]: from sklearn.model_selection import cross_val_score

         forest_scores = cross_val_score(forest_reg, housing_prepared, housing_labels,
                                         scoring="neg_mean_squared_error", cv=10)
         forest_rmse_scores = np.sqrt(-forest_scores)
         display_scores(forest_rmse_scores)
```

```
Scores: [49519.80364233 47461.9115823  50029.02762854 52325.28068953
 49308.39426421 53446.37892622 48634.8036574  47585.73832311
 53490.10699751 50021.5852922 ]
Mean: 50182.303100336096
Standard deviation: 2097.0810550985693
```

```
In [95]: scores = cross_val_score(lin_reg, housing_prepared, housing_labels, scoring="neg_mean_squared_error", cv=10)

         pd.Series(np.sqrt(-scores)).describe()
```

```
Out[95]: count        10.000000
         mean      69052.461363
         std        2879.437224
         min       64969.630564
         25%       67136.363758
         50%       68156.372635
         75%       70982.369487
         max       74739.570526
         dtype: float64
```

```
In [96]: from sklearn.svm import SVR

         svm_reg = SVR(kernel="linear")
         svm_reg.fit(housing_prepared, housing_labels)
         housing_predictions = svm_reg.predict(housing_prepared)
         svm_mse = mean_squared_error(housing_labels, housing_predictions)
         svm_rmse = np.sqrt(svm_mse)
         svm_rmse
```

```
Out[96]: 111094.6308539982
```

```
In [97]:  from sklearn.model_selection import GridSearchCV

          param_grid = [
              # try 12 (3×4) combinations of hyperparameters
              {'n_estimators': [3, 10, 30], 'max_features': [2, 4, 6, 8]},
              # then try 6 (2×3) combinations with bootstrap set as False
              {'bootstrap': [False], 'n_estimators': [3, 10], 'max_features': [2, 3, 4]},
            ]

          forest_reg = RandomForestRegressor(random_state=42)
          # train across 5 folds, that's a total of (12+6)*5=90 rounds of training
          grid_search = GridSearchCV(forest_reg, param_grid, cv=5,
                                     scoring='neg_mean_squared_error',
                                     return_train_score=True)
          grid_search.fit(housing_prepared, housing_labels)

Out[97]: GridSearchCV(cv=5, error_score='raise-deprecating',
                 estimator=RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
                     max_features='auto', max_leaf_nodes=None,
                     min_impurity_decrease=0.0, min_impurity_split=None,
                     min_samples_leaf=1, min_samples_split=2,
                     min_weight_fraction_leaf=0.0, n_estimators='warn', n_jobs=None,
                     oob_score=False, random_state=42, verbose=0, warm_start=False),
                 fit_params=None, iid='warn', n_jobs=None,
                 param_grid=[{'n_estimators': [3, 10, 30], 'max_features': [2, 4, 6, 8]}, {'bootstrap': [False], 'n_estimators': [3, 10],
          'max_features': [2, 3, 4]}],
                 pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
                 scoring='neg_mean_squared_error', verbose=0)


In [98]:  grid_search.best_params_

Out[98]: {'max_features': 8, 'n_estimators': 30}


In [99]:  grid_search.best_estimator_

Out[99]: RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
                 max_features=8, max_leaf_nodes=None, min_impurity_decrease=0.0,
                 min_impurity_split=None, min_samples_leaf=1,
                 min_samples_split=2, min_weight_fraction_leaf=0.0,
                 n_estimators=30, n_jobs=None, oob_score=False, random_state=42,
                 verbose=0, warm_start=False)
```

```
In [100]: cvres = grid_search.cv_results_
          for mean_score, params in zip(cvres["mean_test_score"], cvres["params"]):
              print(np.sqrt(-mean_score), params)
```

```
63669.05791727153 {'max_features': 2, 'n_estimators': 3}
55627.16171305252 {'max_features': 2, 'n_estimators': 10}
53384.57867637289 {'max_features': 2, 'n_estimators': 30}
60965.99185930139 {'max_features': 4, 'n_estimators': 3}
52740.98248528835 {'max_features': 4, 'n_estimators': 10}
50377.344409590376 {'max_features': 4, 'n_estimators': 30}
58663.84733372485 {'max_features': 6, 'n_estimators': 3}
52006.15355973719 {'max_features': 6, 'n_estimators': 10}
50146.465964159885 {'max_features': 6, 'n_estimators': 30}
57869.25504027614 {'max_features': 8, 'n_estimators': 3}
51711.09443660957 {'max_features': 8, 'n_estimators': 10}
49682.25345942335 {'max_features': 8, 'n_estimators': 30}
62895.088889905004 {'bootstrap': False, 'max_features': 2, 'n_estimators': 3}
54658.14484390074 {'bootstrap': False, 'max_features': 2, 'n_estimators': 10}
59470.399594730654 {'bootstrap': False, 'max_features': 3, 'n_estimators': 3}
52725.01091081235 {'bootstrap': False, 'max_features': 3, 'n_estimators': 10}
57490.612956065226 {'bootstrap': False, 'max_features': 4, 'n_estimators': 3}
51009.51445842374 {'bootstrap': False, 'max_features': 4, 'n_estimators': 10}
```

```python
In [101]: pd.DataFrame(grid_search.cv_results_)
```

| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_max_features | param_n_estimators | param_bootstrap | params | split0_test_score | split1_test_score |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.067311 | 0.013484 | 0.003808 | 0.000232 | 2 | 3 | NaN | {'max_features': 2, 'n_estimators': 3} | -3.837622e+09 | -4.147108e+09 |
| 1 | 0.206839 | 0.007902 | 0.011261 | 0.001797 | 2 | 10 | NaN | {'max_features': 2, 'n_estimators': 10} | -3.047771e+09 | -3.254861e+09 |
| 2 | 0.614305 | 0.010618 | 0.034848 | 0.009672 | 2 | 30 | NaN | {'max_features': 2, 'n_estimators': 30} | -2.689185e+09 | -3.021086e+09 |
| 3 | 0.099256 | 0.001032 | 0.003355 | 0.000233 | 4 | 3 | NaN | {'max_features': 4, 'n_estimators': 3} | -3.730181e+09 | -3.786886e+09 |
| 4 | 0.323640 | 0.002283 | 0.008854 | 0.000534 | 4 | 10 | NaN | {'max_features': 4, 'n_estimators': 10} | -2.666283e+09 | -2.784511e+09 |
| 5 | 0.981263 | 0.007037 | 0.030582 | 0.002164 | 4 | 30 | NaN | {'max_features': 4, 'n_estimators': 30} | -2.387153e+09 | -2.588448e+09 |
| 6 | 0.136442 | 0.004478 | 0.003718 | 0.000302 | 6 | 3 | NaN | {'max_features': 6, 'n_estimators': 3} | -3.119657e+09 | -3.586319e+09 |
| 7 | 0.466119 | 0.011197 | 0.011113 | 0.001169 | 6 | 10 | NaN | {'max_features': 6, 'n_estimators': 10} | -2.549663e+09 | -2.782039e+09 |
| 8 | 1.376538 | 0.030661 | 0.029462 | 0.004335 | 6 | 30 | NaN | {'max_features': 6, 'n_estimators': 30} | -2.370010e+09 | -2.583638e+09 |
| 9 | 0.171699 | 0.001843 | 0.003471 | 0.000349 | 8 | 3 | NaN | {'max_features': 8, 'n_estimators': 3} | -3.353504e+09 | -3.348552e+09 |
| 10 | 0.586707 | 0.008500 | 0.011020 | 0.001052 | 8 | 10 | NaN | {'max_features': 8, 'n_estimators': 10} | -2.571970e+09 | -2.718994e+09 |

| | mean_fit_time | std_fit_time | mean_score_time | std_score_time | param_max_features | param_n_estimators | param_bootstrap | params | split0_test_score | split1_test_score |
|---|---|---|---|---|---|---|---|---|---|---|
| **11** | 1.757874 | 0.018076 | 0.029640 | 0.003076 | 8 | 30 | NaN | {'max_features': 8, 'n_estimators': 30} | -2.357390e+09 | -2.546640e+09 |
| **12** | 0.098745 | 0.004214 | 0.004431 | 0.000682 | 2 | 3 | False | {'bootstrap': False, 'max_features': 2, 'n_est... | -3.785816e+09 | -4.166012e+09 |
| **13** | 0.317325 | 0.003440 | 0.011101 | 0.000829 | 2 | 10 | False | {'bootstrap': False, 'max_features': 2, 'n_est... | -2.810721e+09 | -3.107789e+09 |
| **14** | 0.125696 | 0.001922 | 0.004290 | 0.000671 | 3 | 3 | False | {'bootstrap': False, 'max_features': 3, 'n_est... | -3.618324e+09 | -3.441527e+09 |
| **15** | 0.422643 | 0.004466 | 0.012567 | 0.000663 | 3 | 10 | False | {'bootstrap': False, 'max_features': 3, 'n_est... | -2.757999e+09 | -2.851737e+09 |
| **16** | 0.159518 | 0.006842 | 0.004415 | 0.000307 | 4 | 3 | False | {'bootstrap': False, 'max_features': 4, 'n_est... | -3.134040e+09 | -3.559375e+09 |
| **17** | 0.534765 | 0.011351 | 0.012249 | 0.001461 | 4 | 10 | False | {'bootstrap': False, 'max_features': 4, 'n_est... | -2.525578e+09 | -2.710011e+09 |

18 rows × 23 columns

```python
In [102]: from sklearn.model_selection import RandomizedSearchCV
          from scipy.stats import randint

          param_distribs = {
                  'n_estimators': randint(low=1, high=200),
                  'max_features': randint(low=1, high=8),
              }

          forest_reg = RandomForestRegressor(random_state=42)
          rnd_search = RandomizedSearchCV(forest_reg, param_distributions=param_distribs,
                                          n_iter=10, cv=5, scoring='neg_mean_squared_error', random_state=42)
          rnd_search.fit(housing_prepared, housing_labels)
```

```
Out[102]: RandomizedSearchCV(cv=5, error_score='raise-deprecating',
                    estimator=RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
                     max_features='auto', max_leaf_nodes=None,
                     min_impurity_decrease=0.0, min_impurity_split=None,
                     min_samples_leaf=1, min_samples_split=2,
                     min_weight_fraction_leaf=0.0, n_estimators='warn', n_jobs=None,
                     oob_score=False, random_state=42, verbose=0, warm_start=False),
                    fit_params=None, iid='warn', n_iter=10, n_jobs=None,
                    param_distributions={'n_estimators': <scipy.stats._distn_infrastructure.rv_frozen object at 0x105557b38>, 'max_feature
          s': <scipy.stats._distn_infrastructure.rv_frozen object at 0x105557da0>},
                    pre_dispatch='2*n_jobs', random_state=42, refit=True,
                    return_train_score='warn', scoring='neg_mean_squared_error',
                    verbose=0)
```

```python
In [103]: cvres = rnd_search.cv_results_
          for mean_score, params in zip(cvres["mean_test_score"], cvres["params"]):
              print(np.sqrt(-mean_score), params)
```

```
49150.657232934034 {'max_features': 7, 'n_estimators': 180}
51389.85295710133 {'max_features': 5, 'n_estimators': 15}
50796.12045980556 {'max_features': 3, 'n_estimators': 72}
50835.09932039744 {'max_features': 5, 'n_estimators': 21}
49280.90117886215 {'max_features': 7, 'n_estimators': 122}
50774.86679035961 {'max_features': 3, 'n_estimators': 75}
50682.75001237282 {'max_features': 3, 'n_estimators': 88}
49608.94061293652 {'max_features': 5, 'n_estimators': 100}
50473.57642831875 {'max_features': 3, 'n_estimators': 150}
64429.763804893395 {'max_features': 5, 'n_estimators': 2}
```

```
In [104]: feature_importances = grid_search.best_estimator_.feature_importances_
          feature_importances
```

Out[104]: array([7.33442355e-02, 6.29090705e-02, 4.11437985e-02, 1.46726854e-02,
                 1.41064835e-02, 1.48742809e-02, 1.42575993e-02, 3.66158981e-01,
                 5.64191792e-02, 1.08792957e-01, 5.33510773e-02, 1.03114883e-02,
                 1.64780994e-01, 6.02803867e-05, 1.96041560e-03, 2.85647464e-03])

```
In [105]: extra_attribs = ["rooms_per_hhold", "pop_per_hhold", "bedrooms_per_room"]
          #cat_encoder = cat_pipeline.named_steps["cat_encoder"] # old solution
          cat_encoder = full_pipeline.named_transformers_["cat"]
          cat_one_hot_attribs = list(cat_encoder.categories_[0])
          attributes = num_attribs + extra_attribs + cat_one_hot_attribs
          sorted(zip(feature_importances, attributes), reverse=True)
```

Out[105]: [(0.3661589806181342, 'median_income'),
           (0.1647809935615905, 'INLAND'),
           (0.10879295677551573, 'pop_per_hhold'),
           (0.07334423551601242, 'longitude'),
           (0.0629090704826203, 'latitude'),
           (0.05641917918195401, 'rooms_per_hhold'),
           (0.05335107734767581, 'bedrooms_per_room'),
           (0.041143798478729635, 'housing_median_age'),
           (0.014874280890402767, 'population'),
           (0.014672685420543237, 'total_rooms'),
           (0.014257599323407807, 'households'),
           (0.014106483453584102, 'total_bedrooms'),
           (0.010311488326303787, '<1H OCEAN'),
           (0.00256474637320158, 'NEAR OCEAN'),
           (0.00196041559947807, 'NEAR BAY'),
           (6.028038672736599e-05, 'ISLAND')]

```
In [106]: final_model = grid_search.best_estimator_

          X_test = strat_test_set.drop("median_house_value", axis=1)
          y_test = strat_test_set["median_house_value"].copy()

          X_test_prepared = full_pipeline.transform(X_test)
          final_predictions = final_model.predict(X_test_prepared)

          final_mse = mean_squared_error(y_test, final_predictions)
          final_rmse = np.sqrt(final_mse)
```

```
In [107]: final_rmse
```

Out[107]: 47730.22690385927

```
In [108]: from scipy import stats

          confidence = 0.95
          squared_errors = (final_predictions - y_test) ** 2
          np.sqrt(stats.t.interval(confidence, len(squared_errors) - 1,
                                   loc=squared_errors.mean(),
                                   scale=stats.sem(squared_errors)))
```

Out[108]: array([45685.10470776, 49691.25001878])

```
In [109]: m = len(squared_errors) # Computing t score manually
          mean = squared_errors.mean()
          tscore = stats.t.ppf((1 + confidence) / 2, df=m - 1)
          tmargin = tscore * squared_errors.std(ddof=1) / np.sqrt(m)
          np.sqrt(mean - tmargin), np.sqrt(mean + tmargin)
```

Out[109]: (45685.10470776, 49691.25001877858)

```
In [110]: # Alternatively we can use z score
          zscore = stats.norm.ppf((1 + confidence) / 2)
          zmargin = zscore * squared_errors.std(ddof=1) / np.sqrt(m)
          np.sqrt(mean - zmargin), np.sqrt(mean + zmargin)
```

Out[110]: (45685.717918136455, 49690.68623889413)
```