**Chapter 2 – End-to-end Machine Learning project**

*Welcome to Machine Learning Housing Corp.! Your task is to predict median house values in Californian districts, given a number of features from these districts.*

*This notebook contains all the sample code and solutions to the exercices in chapter 2.*

**CO** Run in Google Colab

## ▾ Setup

First, let's make sure this notebook works well in both python 2 and 3, import a few common modules, ensure MatplotLib plots figures inline and prepare a function to save the figures:

```
1  # To support both python 2 and python 3
2  from __future__ import division, print_function, unicode_literals
3
4  # Common imports
5  import numpy as np
6  import os
7
8  # to make this notebook's output stable across runs
9  np.random.seed(42)
10
11 # To plot pretty figures
12 %matplotlib inline
13 import matplotlib as mpl
14 import matplotlib.pyplot as plt
15 mpl.rc('axes', labelsize=14)
16 mpl.rc('xtick', labelsize=12)
17 mpl.rc('ytick', labelsize=12)
18
19 # Where to save the figures
20 PROJECT_ROOT_DIR = "."
21 CHAPTER_ID = "end_to_end_project"
22 IMAGES_PATH = os.path.join(PROJECT_ROOT_DIR, "images", CHAPTER_ID)
23 os.makedirs(IMAGES_PATH, exist_ok=True)
24
25 def save_fig(fig_id, tight_layout=True, fig_extension="png", resolution=300):
26     path = os.path.join(IMAGES_PATH, fig_id + "." + fig_extension)
27     print("Saving figure", fig_id)
28     if tight_layout:
29         plt.tight_layout()
30     plt.savefig(path, format=fig_extension, dpi=resolution)
```

## ▾ Get the data

```
1  import os
2  import tarfile
3  import urllib.request
4
5  DOWNLOAD_ROOT = "https://raw.githubusercontent.com/ageron/handson-ml/master/"
6  HOUSING_PATH = os.path.join("datasets", "housing")
7  HOUSING_URL = DOWNLOAD_ROOT + "datasets/housing/housing.tgz"
8
9  def fetch_housing_data(housing_url=HOUSING_URL, housing_path=HOUSING_PATH):
10     os.makedirs(housing_path, exist_ok=True)
11     tgz_path = os.path.join(housing_path, "housing.tgz")
12     urllib.request.urlretrieve(housing_url, tgz_path)
13     housing_tgz = tarfile.open(tgz_path)
14     housing_tgz.extractall(path=housing_path)
15     housing_tgz.close()
```

```
1  fetch_housing_data()
```

```
1 import pandas as pd
2
3 def load_housing_data(housing_path=HOUSING_PATH):
4     csv_path = os.path.join(housing_path, "housing.csv")
5     return pd.read_csv(csv_path)
```

```
1 housing = load_housing_data()
2 housing.head()
```

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

```
1 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
```

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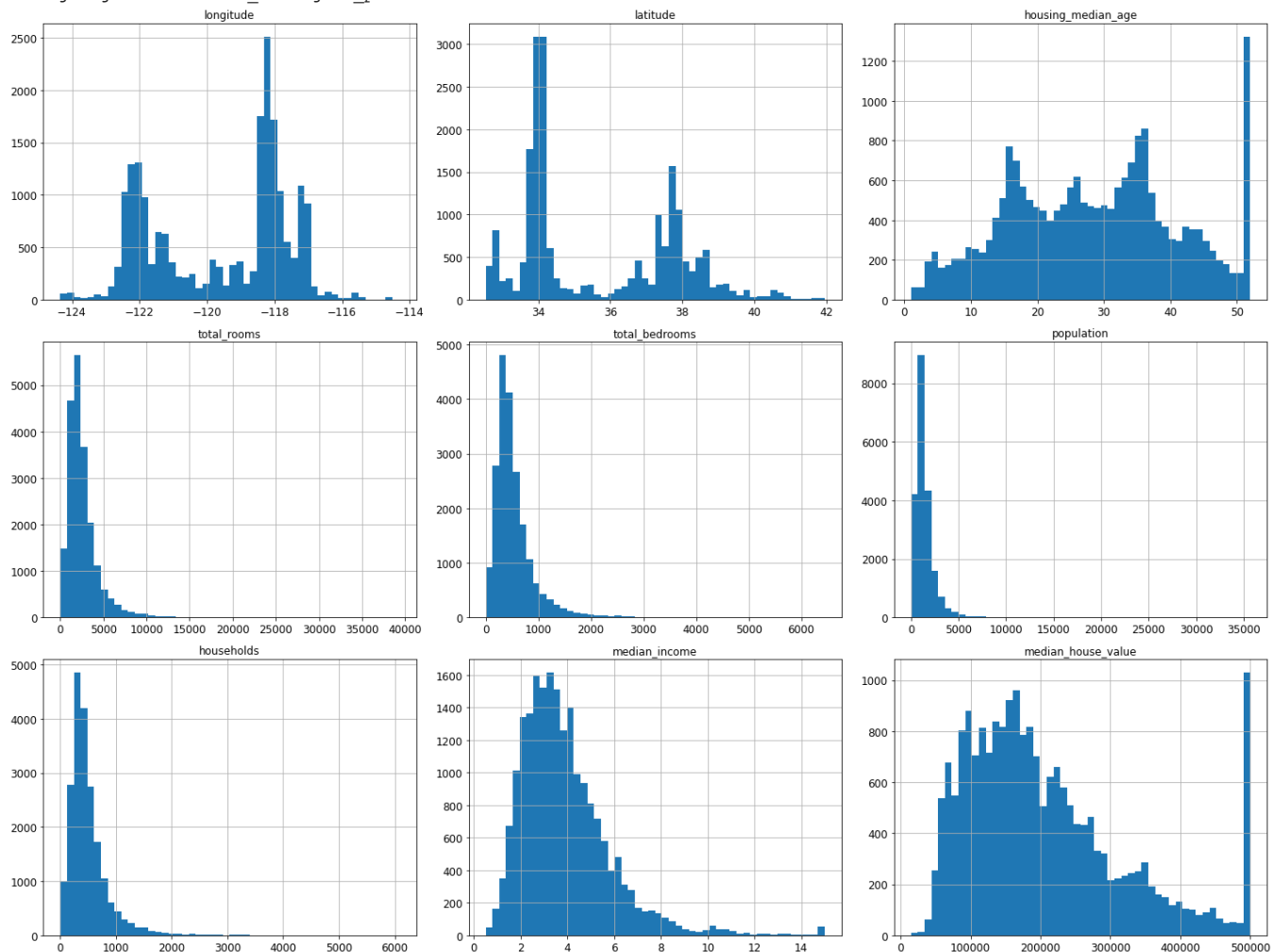
```
<1H OCEAN     9136
INLAND        6551
NEAR OCEAN    2658
NEAR BAY      2290
ISLAND           5
Name: ocean_proximity, dtype: int64
```

```
1 housing.describe()
```

| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | m |
|---|---|---|---|---|---|---|---|---|---|
| count | 20640.000000 | 20640.000000 | 20640.000000 | 20640.000000 | 20433.000000 | 20640.000000 | 20640.000000 | 20640.000000 | |
| mean | -119.569704 | 35.631861 | 28.639486 | 2635.763081 | 537.870553 | 1425.476744 | 499.539680 | 3.870671 | |
| std | 2.003532 | 2.135952 | 12.585558 | 2181.615252 | 421.385070 | 1132.462122 | 382.329753 | 1.899822 | |
| min | -124.350000 | 32.540000 | 1.000000 | 2.000000 | 1.000000 | 3.000000 | 1.000000 | 0.499900 | |
| 25% | -121.800000 | 33.930000 | 18.000000 | 1447.750000 | 296.000000 | 787.000000 | 280.000000 | 2.563400 | |
| 50% | -118.490000 | 34.260000 | 29.000000 | 2127.000000 | 435.000000 | 1166.000000 | 409.000000 | 3.534800 | |
| 75% | -118.010000 | 37.710000 | 37.000000 | 3148.000000 | 647.000000 | 1725.000000 | 605.000000 | 4.743250 | |
| max | -114.310000 | 41.950000 | 52.000000 | 39320.000000 | 6445.000000 | 35682.000000 | 6082.000000 | 15.000100 | |

```
1 %matplotlib inline
2 import matplotlib.pyplot as plt
3 housing.hist(bins=50, figsize=(20,15))
4 save_fig("attribute_histogram_plots")
5 plt.show()
```

```
Saving figure attribute_histogram_plots
```



```
1 # to make this notebook's output identical at every run
2 np.random.seed(42)
```

```
1 import numpy as np
2
3 # For illustration only. Sklearn has train_test_split()
4 def split_train_test(data, test_ratio):
5     shuffled_indices = np.random.permutation(len(data))
6     test_set_size = int(len(data) * test_ratio)
7     test_indices = shuffled_indices[:test_set_size]
8     train_indices = shuffled_indices[test_set_size:]
9     return data.iloc[train_indices], data.iloc[test_indices]
```

```
1 train_set, test_set = split_train_test(housing, 0.2)
2 print(len(train_set), "train +", len(test_set), "test")
```

```
16512 train + 4128 test
```

```
1 from zlib import crc32
2
3 def test_set_check(identifier, test_ratio):
4     return crc32(np.int64(identifier)) & 0xffffffff < test_ratio * 2**32
5
6 def split_train_test_by_id(data, test_ratio, id_column):
7     ids = data[id_column]
8     in_test_set = ids.apply(lambda id_: test_set_check(id_, test_ratio))
9     return data.loc[~in_test_set], data.loc[in_test_set]
```

The implementation of `test_set_check()` above works fine in both Python 2 and Python 3. In earlier releases, the following implementation was proposed, which supported any hash function, but was much slower and did not support Python 2:

```
1 import hashlib
2
3 def test_set_check(identifier, test_ratio, hash=hashlib.md5):
4     return hash(np.int64(identifier)).digest()[-1] < 256 * test_ratio
```

If you want an implementation that supports any hash function and is compatible with both Python 2 and Python 3, here is one:

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

```
1 housing_with_id = housing.reset_index()   # adds an `index` column
2 train_set, test_set = split_train_test_by_id(housing_with_id, 0.2, "index")
```

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

```
1 test_set.head()
```

|    | index | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | med |
|----|-------|-----------|----------|--------------------|-------------|----------------|------------|------------|---------------|-----|
| 8  | 8     | -122.26   | 37.84    | 42.0               | 2555.0      | 665.0          | 1206.0     | 595.0      | 2.0804        |     |
| 10 | 10    | -122.26   | 37.85    | 52.0               | 2202.0      | 434.0          | 910.0      | 402.0      | 3.2031        |     |
| 11 | 11    | -122.26   | 37.85    | 52.0               | 3503.0      | 752.0          | 1504.0     | 734.0      | 3.2705        |     |
| 12 | 12    | -122.26   | 37.85    | 52.0               | 2491.0      | 474.0          | 1098.0     | 468.0      | 3.0750        |     |
| 13 | 13    | -122.26   | 37.84    | 52.0               | 696.0       | 191.0          | 345.0      | 174.0      | 2.6736        |     |

```
1 from sklearn.model_selection import train_test_split
2
3 train_set, test_set = train_test_split(housing, test_size=0.2, random_state=42)
```
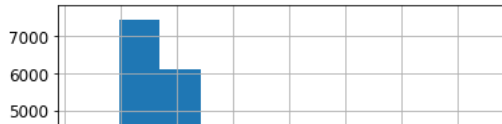
```
1 test_set.head()
```

|       | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | median_ |
|-------|-----------|----------|--------------------|-------------|----------------|------------|------------|---------------|---------|
| 20046 | -119.01   | 36.06    | 25.0               | 1505.0      | NaN            | 1392.0     | 359.0      | 1.6812        |         |
| 3024  | -119.46   | 35.14    | 30.0               | 2943.0      | NaN            | 1565.0     | 584.0      | 2.5313        |         |
| 15663 | -122.44   | 37.80    | 52.0               | 3830.0      | NaN            | 1310.0     | 963.0      | 3.4801        |         |
| 20484 | -118.72   | 34.28    | 17.0               | 3051.0      | NaN            | 1705.0     | 495.0      | 5.7376        |         |
| 9814  | -121.93   | 36.62    | 34.0               | 2351.0      | NaN            | 1063.0     | 428.0      | 3.7250        |         |

```
1 housing["median_income"].hist()
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f86c8c962b0>
```



**Warning**: in the book, I did not use `pd.cut()`, instead I used the code below. The `pd.cut()` solution gives the same result (except the labels are integers instead of floats), but it is simpler to understand:

```
# Divide by 1.5 to limit the number of income categories
housing["income_cat"] = np.ceil(housing["median_income"] / 1.5)
# Label those above 5 as 5
housing["income_cat"].where(housing["income_cat"] < 5, 5.0, inplace=True)
```
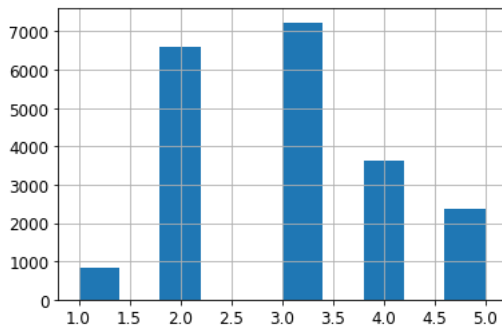
```
1 housing["income_cat"] = pd.cut(housing["median_income"],
2                                 bins=[0., 1.5, 3.0, 4.5, 6., np.inf],
3                                 labels=[1, 2, 3, 4, 5])
```

```
1 housing["income_cat"].value_counts()
```

```
3    7236
2    6581
4    3639
5    2362
1     822
Name: income_cat, dtype: int64
```

```
1 housing["income_cat"].hist()
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f86bcfa15e0>
```



```
1 from sklearn.model_selection import StratifiedShuffleSplit
2
3 split = StratifiedShuffleSplit(n_splits=1, test_size=0.2, random_state=42)
4 for train_index, test_index in split.split(housing, housing["income_cat"]):
5     strat_train_set = housing.loc[train_index]
6     strat_test_set = housing.loc[test_index]
```

```
1 strat_test_set["income_cat"].value_counts() / len(strat_test_set)
```

```
3    0.350533
2    0.318798
4    0.176357
5    0.114341
1    0.039971
Name: income_cat, dtype: float64
```

```
1 housing["income_cat"].value_counts() / len(housing)
```

```
3    0.350581
2    0.318847
4    0.176308
5    0.114438
1    0.039826
Name: income_cat, dtype: float64
```

```
1 def income_cat_proportions(data):
2     return data["income_cat"].value_counts() / len(data)
3
4 train_set, test_set = train_test_split(housing, test_size=0.2, random_state=42)
5
6 compare_props = pd.DataFrame({
7     "Overall": income_cat_proportions(housing),
8     "Stratified": income_cat_proportions(strat_test_set),
9     "Random": income_cat_proportions(test_set),
10 }).sort_index()
11 compare_props["Rand. %error"] = 100 * compare_props["Random"] / compare_props["Overall"] - 100
12 compare_props["Strat. %error"] = 100 * compare_props["Stratified"] / compare_props["Overall"] - 100
```

```
1 compare_props
```

|   | Overall | Stratified | Random | Rand. %error | Strat. %error |
|---|---------|------------|--------|--------------|---------------|
| 1 | 0.039826 | 0.039971 | 0.040213 | 0.973236 | 0.364964 |
| 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 | 0.114341 | 0.109496 | -4.318374 | -0.084674 |

```
1 for set_ in (strat_train_set, strat_test_set):
2     set_.drop("income_cat", axis=1, inplace=True)
```
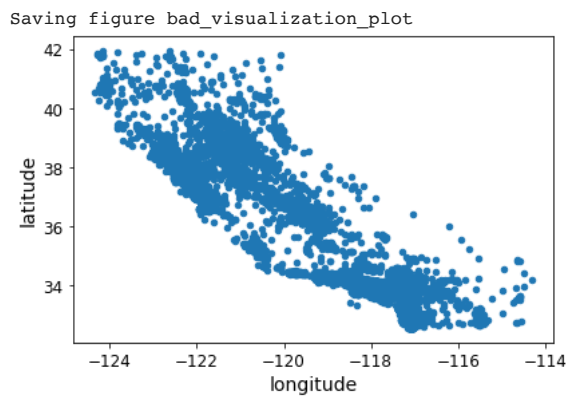
## Discover and visualize the data to gain insights

```
1 housing = strat_train_set.copy()
```

```
1 housing.plot(kind="scatter", x="longitude", y="latitude")
2 save_fig("bad_visualization_plot")
```

```
Saving figure bad_visualization_plot
```
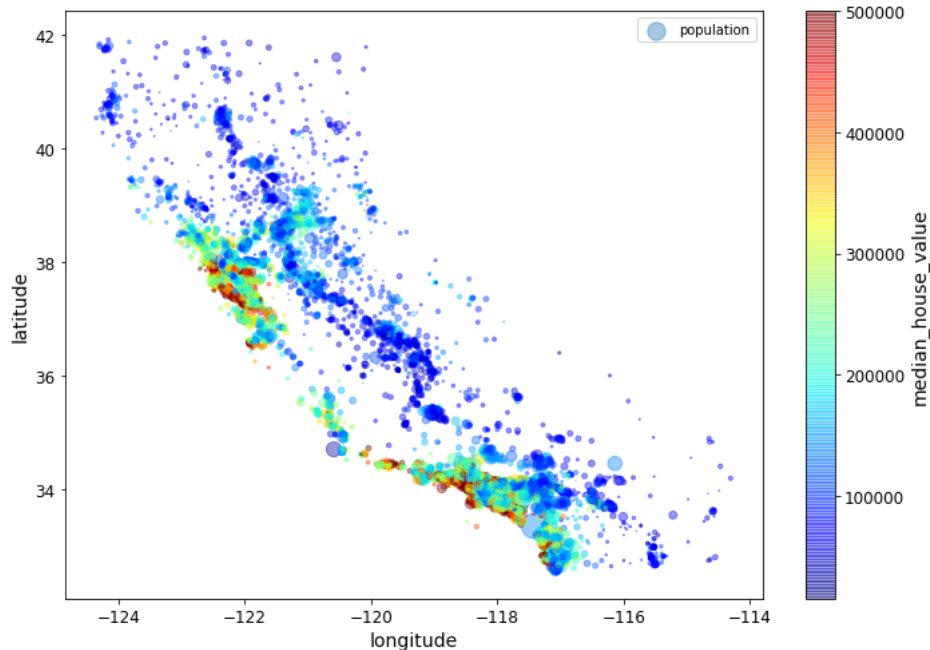


```
1 housing.plot(kind="scatter", x="longitude", y="latitude", alpha=0.1)
2 save_fig("better_visualization_plot")
```

The argument `sharex=False` fixes a display bug (the x-axis values and legend were not displayed). This is a temporary fix (see:
https://github.com/pandas-dev/pandas/issues/10611). Thanks to Wilmer Arellano for pointing it out.

```
1 housing.plot(kind="scatter", x="longitude", y="latitude", alpha=0.4,
2     s=housing["population"]/100, label="population", figsize=(10,7),
3     c="median_house_value", cmap=plt.get_cmap("jet"), colorbar=True,
4     sharex=False)
5 plt.legend()
6 save_fig("housing_prices_scatterplot")
```

```
Saving figure housing_prices_scatterplot
```
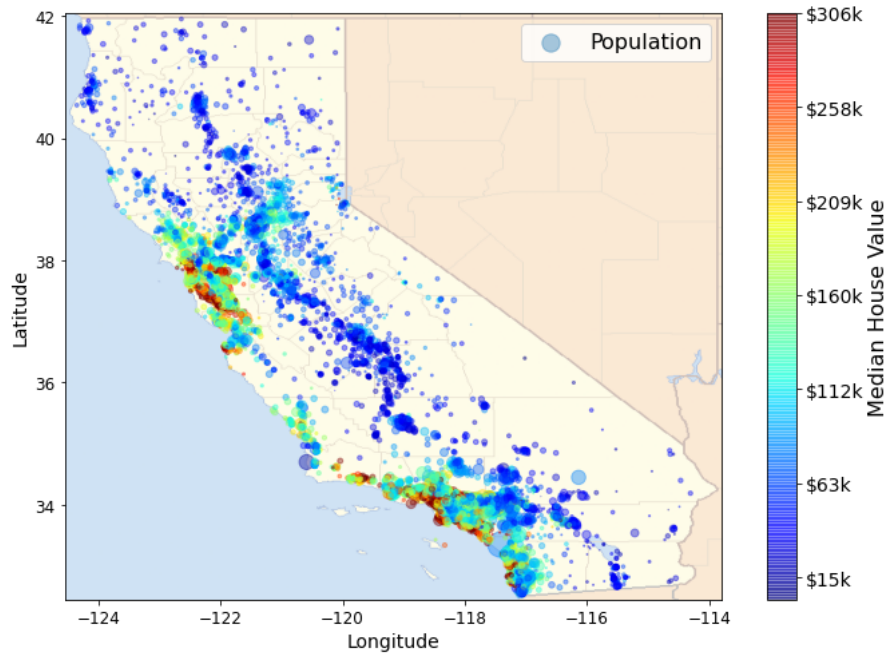


```
1 # Download the California image
2 images_path = os.path.join(PROJECT_ROOT_DIR, "images", "end_to_end_project")
3 os.makedirs(images_path, exist_ok=True)
4 DOWNLOAD_ROOT = "https://raw.githubusercontent.com/ageron/handson-ml/master/"
5 filename = "california.png"
6 print("Downloading", filename)
7 url = DOWNLOAD_ROOT + "images/end_to_end_project/" + filename
8 urllib.request.urlretrieve(url, os.path.join(images_path, filename))
```

```
Downloading california.png
('./images/end_to_end_project/california.png',
 <http.client.HTTPMessage at 0x7f86bce35670>)
```

```
1 import matplotlib.image as mpimg
2 california_img=mpimg.imread(PROJECT_ROOT_DIR + '/images/end_to_end_project/california.png')
3 ax = housing.plot(kind="scatter", x="longitude", y="latitude", figsize=(10,7),
4                       s=housing['population']/100, label="Population",
5                       c="median_house_value", cmap=plt.get_cmap("jet"),
6                       colorbar=False, alpha=0.4,
7                       )
8 plt.imshow(california_img, extent=[-124.55, -113.80, 32.45, 42.05], alpha=0.5,
9           cmap=plt.get_cmap("jet"))
10 plt.ylabel("Latitude", fontsize=14)
11 plt.xlabel("Longitude", fontsize=14)
12
13 prices = housing["median_house_value"]
14 tick_values = np.linspace(prices.min(), prices.max(), 11)
15 cbar = plt.colorbar()
16 cbar.ax.set_yticklabels(["$%dk"%(round(v/1000)) for v in tick_values], fontsize=14)
17 cbar.set_label('Median House Value', fontsize=16)
18
19 plt.legend(fontsize=16)
20 save_fig("california_housing_prices_plot")
21 plt.show()
```

```
Saving figure california_housing_prices_plot
```



```
1 corr_matrix = housing.corr()
```

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

```
median_house_value    1.000000
median_income         0.687151
total_rooms           0.135140
housing_median_age    0.114146
households            0.064590
total_bedrooms        0.047781
population            -0.026882
longitude             -0.047466
latitude              -0.142673
Name: median_house_value, dtype: float64
```

```
1 # from pandas.tools.plotting import scatter_matrix # For older versions of Pandas
2 from pandas.plotting import scatter_matrix
3
4 attributes = ["median_house_value", "median_income", "total_rooms",
5               "housing_median_age"]
6 scatter_matrix(housing[attributes], figsize=(12, 8))
7 save_fig("scatter_matrix_plot")
```
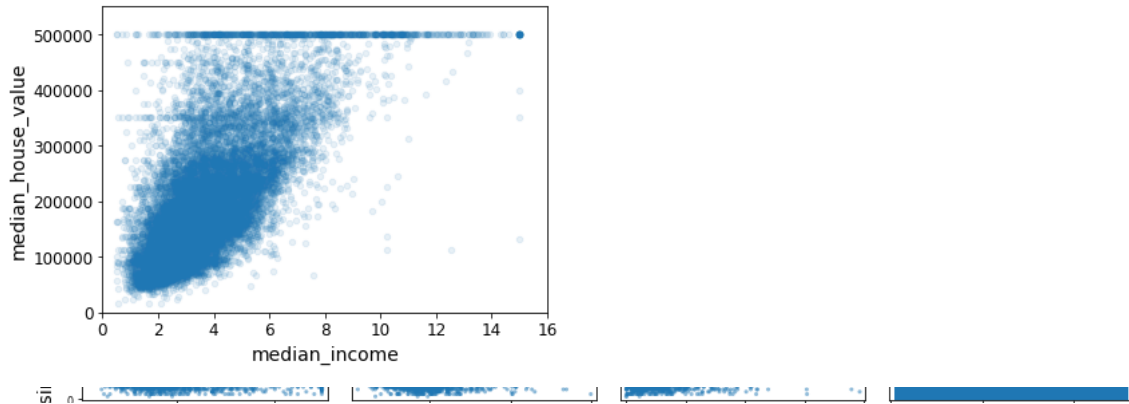
Saving figure scatter_matrix_plot



```
1 housing.plot(kind="scatter", x="median_income", y="median_house_value",
2                 alpha=0.1)
3 plt.axis([0, 16, 0, 550000])
4 save_fig("income_vs_house_value_scatterplot")
```

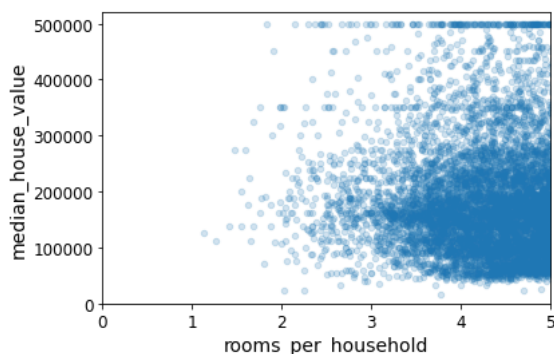Saving figure income_vs_house_value_scatterplot



```
1 housing["rooms_per_household"] = housing["total_rooms"]/housing["households"]
2 housing["bedrooms_per_room"] = housing["total_bedrooms"]/housing["total_rooms"]
3 housing["population_per_household"]=housing["population"]/housing["households"]
```

Note: there was a bug in the previous cell, in the definition of the `rooms_per_household` attribute. This explains why the correlation value below differs slightly from the value in the book (unless you are reading the latest version).

```
1 corr_matrix = housing.corr()
2 corr_matrix["median_house_value"].sort_values(ascending=False)
```

```
median_house_value          1.000000
median_income               0.687151
rooms_per_household         0.146255
total_rooms                 0.135140
housing_median_age          0.114146
households                  0.064590
total_bedrooms              0.047781
population_per_household    -0.021991
population                  -0.026882
longitude                   -0.047466
latitude                    -0.142673
bedrooms_per_room           -0.259952
Name: median_house_value, dtype: float64
```

```
1 housing.plot(kind="scatter", x="rooms_per_household", y="median_house_value",
2                 alpha=0.2)
3 plt.axis([0, 5, 0, 520000])
4 plt.show()
```



```
1 housing.describe()
```

|       | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | m |
|-------|-----------|----------|--------------------|-------------|----------------|------------|------------|---------------|---|
| count | 16512.000000 | 16512.000000 | 16512.000000 | 16512.000000 | 16354.000000 | 16512.000000 | 16512.000000 | 16512.000000 | |
| mean | -119.575635 | 35.639314 | 28.653404 | 2622.539789 | 534.914639 | 1419.687379 | 497.011810 | 3.875884 | |
| std | 2.001828 | 2.137963 | 12.574819 | 2138.417080 | 412.665649 | 1115.663036 | 375.696156 | 1.904931 | |
| min | -124.350000 | 32.540000 | 1.000000 | 6.000000 | 2.000000 | 3.000000 | 2.000000 | 0.499900 | |
| 25% | -121.800000 | 33.940000 | 18.000000 | 1443.000000 | 295.000000 | 784.000000 | 279.000000 | 2.566950 | |
| 50% | -118.510000 | 34.260000 | 29.000000 | 2119.000000 | 433.000000 | 1164.000000 | 408.000000 | 3.541550 | |
| 75% | -118.010000 | 37.720000 | 37.000000 | 3141.000000 | 644.000000 | 1719.000000 | 602.000000 | 4.745325 | |
| max | -114.310000 | 41.950000 | 52.000000 | 39320.000000 | 6210.000000 | 35682.000000 | 5358.000000 | 15.000100 | |

## Prepare the data for Machine Learning algorithms

```
1 housing = strat_train_set.drop("median_house_value", axis=1) # drop labels for training set
2 housing_labels = strat_train_set["median_house_value"].copy()
```

```
1 sample_incomplete_rows = housing[housing.isnull().any(axis=1)].head()
2 sample_incomplete_rows
```

|       | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | ocean_p |
|-------|-----------|----------|--------------------|-------------|----------------|------------|------------|---------------|---------|
| 1606 | -122.08 | 37.88 | 26.0 | 2947.0 | NaN | 825.0 | 626.0 | 2.9330 | |
| 10915 | -117.87 | 33.73 | 45.0 | 2264.0 | NaN | 1970.0 | 499.0 | 3.4193 | < |
| 19150 | -122.70 | 38.35 | 14.0 | 2313.0 | NaN | 954.0 | 397.0 | 3.7813 | < |
| 4186 | -118.23 | 34.13 | 48.0 | 1308.0 | NaN | 835.0 | 294.0 | 4.2891 | < |
| 16885 | -122.40 | 37.58 | 26.0 | 3281.0 | NaN | 1145.0 | 480.0 | 6.3580 | NE/ |

```
1 sample_incomplete_rows.dropna(subset=["total_bedrooms"])    # option 1
```

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

```
1 sample_incomplete_rows.drop("total_bedrooms", axis=1)       # option 2
```

|       | longitude | latitude | housing_median_age | total_rooms | population | households | median_income | ocean_proximity |
|-------|-----------|----------|--------------------|-------------|------------|------------|---------------|-----------------|
| 1606 | -122.08 | 37.88 | 26.0 | 2947.0 | 825.0 | 626.0 | 2.9330 | NEAR BAY |
| 10915 | -117.87 | 33.73 | 45.0 | 2264.0 | 1970.0 | 499.0 | 3.4193 | <1H OCEAN |
| 19150 | -122.70 | 38.35 | 14.0 | 2313.0 | 954.0 | 397.0 | 3.7813 | <1H OCEAN |
| 4186 | -118.23 | 34.13 | 48.0 | 1308.0 | 835.0 | 294.0 | 4.2891 | <1H OCEAN |
| 16885 | -122.40 | 37.58 | 26.0 | 3281.0 | 1145.0 | 480.0 | 6.3580 | NEAR OCEAN |

```
1 median = housing["total_bedrooms"].median()
2 sample_incomplete_rows["total_bedrooms"].fillna(median, inplace=True) # option 3
3 sample_incomplete_rows
```

|       | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | ocean_p |
|-------|-----------|----------|--------------------|-------------|----------------|------------|------------|---------------|---------|
| 1606 | -122.08 | 37.88 | 26.0 | 2947.0 | 433.0 | 825.0 | 626.0 | 2.9330 | |
| 10915 | -117.87 | 33.73 | 45.0 | 2264.0 | 433.0 | 1970.0 | 499.0 | 3.4193 | < |
| 19150 | -122.70 | 38.35 | 14.0 | 2313.0 | 433.0 | 954.0 | 397.0 | 3.7813 | < |
| 4186 | -118.23 | 34.13 | 48.0 | 1308.0 | 433.0 | 835.0 | 294.0 | 4.2891 | < |
| 16885 | -122.40 | 37.58 | 26.0 | 3281.0 | 433.0 | 1145.0 | 480.0 | 6.3580 | NE/ |

**Warning**: Since Scikit-Learn 0.20, the `sklearn.preprocessing.Imputer` class was replaced by the `sklearn.impute.SimpleImputer` class.

```
1 try:
2     from sklearn.impute import SimpleImputer # Scikit-Learn 0.20+
3 except ImportError:
4     from sklearn.preprocessing import Imputer as SimpleImputer
5
6 imputer = SimpleImputer(strategy="median")
```

Remove the text attribute because median can only be calculated on numerical attributes:

```
1 housing_num = housing.drop('ocean_proximity', axis=1)
2 # alternatively: housing_num = housing.select_dtypes(include=[np.number])
```

```
1 imputer.fit(housing_num)
```

```
   SimpleImputer(strategy='median')
```

```
1 imputer.statistics_
```

```
   array([-118.51  ,   34.26  ,   29.    , 2119.    ,  433.    ,
          1164.    ,  408.    ,    3.54155])
```

Check that this is the same as manually computing the median of each attribute:

```
1 housing_num.median().values
```

```
   array([-118.51  ,   34.26  ,   29.    , 2119.    ,  433.    ,
          1164.    ,  408.    ,    3.54155])
```

Transform the training set:

```
1 X = imputer.transform(housing_num)
```

```
1 housing_tr = pd.DataFrame(X, columns=housing_num.columns,
2                           index=housing.index)
```

```
1 housing_tr.loc[sample_incomplete_rows.index.values]
```

|       | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income |
|-------|-----------|----------|--------------------|-------------|----------------|------------|------------|---------------|
| 1606  | -122.08   | 37.88    | 26.0               | 2947.0      | 433.0          | 825.0      | 626.0      | 2.9330        |
| 10915 | -117.87   | 33.73    | 45.0               | 2264.0      | 433.0          | 1970.0     | 499.0      | 3.4193        |
| 19150 | -122.70   | 38.35    | 14.0               | 2313.0      | 433.0          | 954.0      | 397.0      | 3.7813        |
| 4186  | -118.23   | 34.13    | 48.0               | 1308.0      | 433.0          | 835.0      | 294.0      | 4.2891        |
| 16885 | -122.40   | 37.58    | 26.0               | 3281.0      | 433.0          | 1145.0     | 480.0      | 6.3580        |

```
1 imputer.strategy
```

```
   'median'
```

```
1 housing_tr = pd.DataFrame(X, columns=housing_num.columns,
2                           index=housing_num.index)
3 housing_tr.head()
```

Now let's preprocess the categorical input feature, `ocean_proximity`:

```
1 housing_cat = housing[['ocean_proximity']]
2 housing_cat.head(10)
```

|       | ocean_proximity |
|-------|-----------------|
| 12655 | INLAND          |
| 15502 | NEAR OCEAN      |
| 2908  | INLAND          |
| 14053 | NEAR OCEAN      |
| 20496 | <1H OCEAN       |
| 1481  | NEAR BAY        |
| 18125 | <1H OCEAN       |
| 5830  | <1H OCEAN       |
| 17989 | <1H OCEAN       |
| 4861  | <1H OCEAN       |

**Warning**: earlier versions of the book used the `LabelEncoder` class or Pandas' `Series.factorize()` method to encode string categorical attributes as integers. However, the `OrdinalEncoder` class that was introduced in Scikit-Learn 0.20 (see PR #10521) is preferable since it is designed for input features ( `x` instead of labels `y` ) and it plays well with pipelines (introduced later in this notebook). If you are using an older version of Scikit-Learn (<0.20), then you can import it from `future_encoders.py` instead.

```
1 try:
2     from sklearn.preprocessing import OrdinalEncoder
3 except ImportError:
4     from future_encoders import OrdinalEncoder # Scikit-Learn < 0.20
```

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

```
    array([[1.],
           [4.],
           [1.],
           [4.],
           [0.],
           [3.],
           [0.],
           [0.],
           [0.],
           [0.]])
```

```
1 ordinal_encoder.categories_
```

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

**Warning**: earlier versions of the book used the `LabelBinarizer` or `CategoricalEncoder` classes to convert each categorical value to a one-hot vector. It is now preferable to use the `OneHotEncoder` class. Since Scikit-Learn 0.20 it can handle string categorical inputs (see PR #10521), not just integer categorical inputs. If you are using an older version of Scikit-Learn, you can import the new version from `future_encoders.py`:

```
1 try:
2     from sklearn.preprocessing import OrdinalEncoder # just to raise an ImportError if Scikit-Learn < 0.20
3     from sklearn.preprocessing import OneHotEncoder
4 except ImportError:
5     from future_encoders import OneHotEncoder # Scikit-Learn < 0.20
6
7 cat_encoder = OneHotEncoder()
8 housing_cat_1hot = cat_encoder.fit_transform(housing_cat)
9 housing_cat_1hot
```

```
<16512x5 sparse matrix of type '<class 'numpy.float64'>'
        with 16512 stored elements in Compressed Sparse Row format>
```

By default, the `OneHotEncoder` class returns a sparse array, but we can convert it to a dense array if needed by calling the `toarray()` method:

```
1 housing_cat_1hot.toarray()
```

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

Alternatively, you can set `sparse=False` when creating the `OneHotEncoder`:

```
1 cat_encoder = OneHotEncoder(sparse=False)
2 housing_cat_1hot = cat_encoder.fit_transform(housing_cat)
3 housing_cat_1hot
```

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

```
1 cat_encoder.categories_
```

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

Let's create a custom transformer to add extra attributes:

```
1 housing.columns
```

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

```
 1 from sklearn.base import BaseEstimator, TransformerMixin
 2
 3 # get the right column indices: safer than hard-coding indices 3, 4, 5, 6
 4 rooms_ix, bedrooms_ix, population_ix, household_ix = [
 5     list(housing.columns).index(col)
 6     for col in ("total_rooms", "total_bedrooms", "population", "households")]
 7
 8 class CombinedAttributesAdder(BaseEstimator, TransformerMixin):
 9     def __init__(self, add_bedrooms_per_room = True): # no *args or **kwargs
10         self.add_bedrooms_per_room = add_bedrooms_per_room
11     def fit(self, X, y=None):
12         return self  # nothing else to do
13     def transform(self, X, y=None):
14         rooms_per_household = X[:, rooms_ix] / X[:, household_ix]
15         population_per_household = X[:, population_ix] / X[:, household_ix]
16         if self.add_bedrooms_per_room:
17             bedrooms_per_room = X[:, bedrooms_ix] / X[:, rooms_ix]
18             return np.c_[X, rooms_per_household, population_per_household,
19                          bedrooms_per_room]
20         else:
21             return np.c_[X, rooms_per_household, population_per_household]
22
23 attr_adder = CombinedAttributesAdder(add_bedrooms_per_room=False)
24 housing_extra_attribs = attr_adder.transform(housing.values)
```

Alternatively, you can use Scikit-Learn's `FunctionTransformer` class that lets you easily create a transformer based on a transformation function (thanks to [Hanmin Qin](#) for suggesting this code). Note that we need to set `validate=False` because the data contains non-float values (`validate` will default to `False` in Scikit-Learn 0.22).

```
1 from sklearn.preprocessing import FunctionTransformer
2
3 def add_extra_features(X, add_bedrooms_per_room=True):
4     rooms_per_household = X[:, rooms_ix] / X[:, household_ix]
5     population_per_household = X[:, population_ix] / X[:, household_ix]
6     if add_bedrooms_per_room:
7         bedrooms_per_room = X[:, bedrooms_ix] / X[:, rooms_ix]
8         return np.c_[X, rooms_per_household, population_per_household,
9                      bedrooms_per_room]
10    else:
11        return np.c_[X, rooms_per_household, population_per_household]
12
13 attr_adder = FunctionTransformer(add_extra_features, validate=False,
14                                  kw_args={"add_bedrooms_per_room": False})
15 housing_extra_attribs = attr_adder.fit_transform(housing.values)
```

```
1 housing_extra_attribs = pd.DataFrame(
2     housing_extra_attribs,
3     columns=list(housing.columns)+["rooms_per_household", "population_per_household"],
4     index=housing.index)
5 housing_extra_attribs.head()
```

| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population | households | median_income | ocean_p |
|---|---|---|---|---|---|---|---|---|---|
| 12655 | -121.46 | 38.52 | 29.0 | 3873.0 | 797.0 | 2237.0 | 706.0 | 2.1736 | |
| 15502 | -117.23 | 33.09 | 7.0 | 5320.0 | 855.0 | 2015.0 | 768.0 | 6.3373 | NE/ |
| 2908 | -119.04 | 35.37 | 44.0 | 1618.0 | 310.0 | 667.0 | 300.0 | 2.875 | |
| 14053 | -117.13 | 32.75 | 24.0 | 1877.0 | 519.0 | 898.0 | 483.0 | 2.2264 | NE/ |
| 20496 | -118.7 | 34.28 | 27.0 | 3536.0 | 646.0 | 1837.0 | 580.0 | 4.4964 | < |

Now let's build a pipeline for preprocessing the numerical attributes (note that we could use `CombinedAttributesAdder()` instead of `FunctionTransformer(...)` if we preferred):

```
1 from sklearn.pipeline import Pipeline
2 from sklearn.preprocessing import StandardScaler
3
4 num_pipeline = Pipeline([
5         ('imputer', SimpleImputer(strategy="median")),
6         ('attribs_adder', FunctionTransformer(add_extra_features, validate=False)),
7         ('std_scaler', StandardScaler()),
8     ])
9
10 housing_num_tr = num_pipeline.fit_transform(housing_num)
```

```
1 housing_num_tr
```

```
array([[-0.94135046,  1.34743822,  0.02756357, ...,  0.01739526,
         0.00622264, -0.12112176],
       [ 1.17178212, -1.19243966, -1.72201763, ...,  0.56925554,
        -0.04081077, -0.81086696],
       [ 0.26758118, -0.1259716 ,  1.22045984, ..., -0.01802432,
        -0.07537122, -0.33827252],
       ...,
       [-1.5707942 ,  1.31001828,  1.53856552, ..., -0.5092404 ,
        -0.03743619,  0.32286937],
       [-1.56080303,  1.2492109 , -1.1653327 , ...,  0.32814891,
        -0.05915604, -0.45702273],
       [-1.28105026,  2.02567448, -0.13148926, ...,  0.01407228,
         0.00657083, -0.12169672]])
```

**Warning**: earlier versions of the book applied different transformations to different columns using a solution based on a `DataFrameSelector` transformer and a `FeatureUnion` (see below). It is now preferable to use the `ColumnTransformer` class that was introduced in Scikit-Learn 0.20. If you are using an older version of Scikit-Learn, you can import it from `future_encoders.py`:

```
1 try:
2     from sklearn.compose import ColumnTransformer
3 except ImportError:
4     from future_encoders import ColumnTransformer # Scikit-Learn < 0.20
```

```
1 num_attribs = list(housing_num)
2 cat_attribs = ["ocean_proximity"]
3
4 full_pipeline = ColumnTransformer([
5         ("num", num_pipeline, num_attribs),
6         ("cat", OneHotEncoder(), cat_attribs),
7     ])
8
9 housing_prepared = full_pipeline.fit_transform(housing)
```

```
1 housing_prepared
```

```
   array([[-0.94135046,  1.34743822,  0.02756357, ...,  0.         ,
            0.         ,  0.         ],
          [ 1.17178212, -1.19243966, -1.72201763, ...,  0.         ,
            0.         ,  1.         ],
          [ 0.26758118, -0.1259716 ,  1.22045984, ...,  0.         ,
            0.         ,  0.         ],
          ...,
          [-1.5707942 ,  1.31001828,  1.53856552, ...,  0.         ,
            0.         ,  0.         ],
          [-1.56080303,  1.2492109 , -1.1653327 , ...,  0.         ,
            0.         ,  0.         ],
          [-1.28105026,  2.02567448, -0.13148926, ...,  0.         ,
            0.         ,  0.         ]])
```

```
1 housing_prepared.shape
```

```
   (16512, 16)
```

For reference, here is the old solution based on a `DataFrameSelector` transformer (to just select a subset of the Pandas `DataFrame` columns), and a `FeatureUnion`:

```
1 from sklearn.base import BaseEstimator, TransformerMixin
2
3 # Create a class to select numerical or categorical columns
4 class OldDataFrameSelector(BaseEstimator, TransformerMixin):
5     def __init__(self, attribute_names):
6         self.attribute_names = attribute_names
7     def fit(self, X, y=None):
8         return self
9     def transform(self, X):
10        return X[self.attribute_names].values
```

Now let's join all these components into a big pipeline that will preprocess both the numerical and the categorical features (again, we could use `CombinedAttributesAdder()` instead of `FunctionTransformer(...)` if we preferred):

```
1 num_attribs = list(housing_num)
2 cat_attribs = ["ocean_proximity"]
3
4 old_num_pipeline = Pipeline([
5         ('selector', OldDataFrameSelector(num_attribs)),
6         ('imputer', SimpleImputer(strategy="median")),
7         ('attribs_adder', FunctionTransformer(add_extra_features, validate=False)),
8         ('std_scaler', StandardScaler()),
9     ])
10
11 old_cat_pipeline = Pipeline([
12        ('selector', OldDataFrameSelector(cat_attribs)),
13        ('cat_encoder', OneHotEncoder(sparse=False)),
14    ])
```

```
1 from sklearn.pipeline import FeatureUnion
2
3 old_full_pipeline = FeatureUnion(transformer_list=[
4         ("num_pipeline", old_num_pipeline),
```

```
5         ("cat_pipeline", old_cat_pipeline),
```

```
1 old_housing_prepared = old_full_pipeline.fit_transform(housing)
2 old_housing_prepared
```

```
array([[-0.94135046,  1.34743822,  0.02756357, ...,  0.         ,
         0.        ,  0.        ],
       [ 1.17178212, -1.19243966, -1.72201763, ...,  0.         ,
         0.        ,  1.        ],
       [ 0.26758118, -0.1259716 ,  1.22045984, ...,  0.         ,
         0.        ,  0.        ],
       ...,
       [-1.5707942 ,  1.31001828,  1.53856552, ...,  0.         ,
         0.        ,  0.        ],
       [-1.56080303,  1.2492109 , -1.1653327 , ...,  0.         ,
         0.        ,  0.        ],
       [-1.28105026,  2.02567448, -0.13148926, ...,  0.         ,
         0.        ,  0.        ]])
```

The result is the same as with the `ColumnTransformer`:

```
1 np.allclose(housing_prepared, old_housing_prepared)
```

```
True
```

# Select and train a model

```
1 from sklearn.linear_model import LinearRegression
2
3 lin_reg = LinearRegression()
4 lin_reg.fit(housing_prepared, housing_labels)
```

```
LinearRegression()
```

```
1 # let's try the full preprocessing pipeline on a few training instances
2 some_data = housing.iloc[:5]
3 some_labels = housing_labels.iloc[:5]
4 some_data_prepared = full_pipeline.transform(some_data)
5
6 print("Predictions:", lin_reg.predict(some_data_prepared))
```

```
Predictions: [ 85657.90192014 305492.60737488 152056.46122456 186095.70946094
  244550.67966089]
```

Compare against the actual values:

```
1 print("Labels:", list(some_labels))
```

```
Labels: [72100.0, 279600.0, 82700.0, 112500.0, 238300.0]
```

```
1 some_data_prepared
```

```
array([[-0.94135046,  1.34743822,  0.02756357,  0.58477745,  0.64037127,
         0.73260236,  0.55628602, -0.8936472 ,  0.01739526,  0.00622264,
        -0.12112176,  0.        ,  1.        ,  0.        ,  0.        ,
         0.        ],
       [ 1.17178212, -1.19243966, -1.72201763,  1.26146668,  0.78156132,
         0.53361152,  0.72131799,  1.292168  ,  0.56925554, -0.04081077,
        -0.81086696,  0.        ,  0.        ,  0.        ,  0.        ,
         1.        ],
       [ 0.26758118, -0.1259716 ,  1.22045984, -0.46977281, -0.54513828,
        -0.67467519, -0.52440722, -0.52543365, -0.01802432, -0.07537122,
        -0.33827252,  0.        ,  1.        ,  0.        ,  0.        ,
         0.        ],
       [ 1.22173797, -1.35147437, -0.37006852, -0.34865152, -0.03636724,
        -0.46761716, -0.03729672, -0.86592882, -0.59513997, -0.10680295,
         0.96120521,  0.        ,  0.        ,  0.        ,  0.        ,
         1.        ],
       [ 0.43743108, -0.63581817, -0.13148926,  0.42717947,  0.27279028,
         0.37406031,  0.22089846,  0.32575178,  0.2512412 ,  0.00610923,
        -0.47451338,  1.        ,  0.        ,  0.        ,  0.        ,
         0.        ]])
```

```
1 from sklearn.metrics import mean_squared_error
2
3 housing_predictions = lin_reg.predict(housing_prepared)
4 lin_mse = mean_squared_error(housing_labels, housing_predictions)
5 lin_rmse = np.sqrt(lin_mse)
6 lin_rmse
```

```
68627.87390018745
```

```
1 from sklearn.metrics import mean_absolute_error
2
3 lin_mae = mean_absolute_error(housing_labels, housing_predictions)
4 lin_mae
```

```
49438.66860915802
```

```
1 from sklearn.tree import DecisionTreeRegressor
2
3 tree_reg = DecisionTreeRegressor(random_state=42)
4 tree_reg.fit(housing_prepared, housing_labels)
```

```
DecisionTreeRegressor(random_state=42)
```

```
1 housing_predictions = tree_reg.predict(housing_prepared)
2 tree_mse = mean_squared_error(housing_labels, housing_predictions)
3 tree_rmse = np.sqrt(tree_mse)
4 tree_rmse
```

```
0.0
```

## ▾ Fine-tune your model

```
1 from sklearn.model_selection import cross_val_score
2
3 scores = cross_val_score(tree_reg, housing_prepared, housing_labels,
4                          scoring="neg_mean_squared_error", cv=10)
5 tree_rmse_scores = np.sqrt(-scores)
```

```
1 def display_scores(scores):
2     print("Scores:", scores)
3     print("Mean:", scores.mean())
4     print("Standard deviation:", scores.std())
5
6 display_scores(tree_rmse_scores)
```

```
1 lin_scores = cross_val_score(lin_reg, housing_prepared, housing_labels,
2                              scoring="neg_mean_squared_error", cv=10)
3 lin_rmse_scores = np.sqrt(-lin_scores)
4 display_scores(lin_rmse_scores)
```

```
Scores: [71762.76364394 64114.99166359 67771.17124356 68635.19072082
 66846.14089488 72528.03725385 73997.08050233 68802.33629334
 66443.28836884 70139.79923956]
Mean: 69104.07998247063
Standard deviation: 2880.3282098180634
```

**Note**: we specify `n_estimators=10` to avoid a warning about the fact that the default value is going to change to 100 in Scikit-Learn 0.22.

```
1 from sklearn.ensemble import RandomForestRegressor
2
3 forest_reg = RandomForestRegressor(n_estimators=10, random_state=42)
4 forest_reg.fit(housing_prepared, housing_labels)
```

```
RandomForestRegressor(n_estimators=10, random_state=42)
```

```
1 housing_predictions = forest_reg.predict(housing_prepared)
2 forest_mse = mean_squared_error(housing_labels, housing_predictions)
3 forest_rmse = np.sqrt(forest_mse)
4 forest_rmse
```

```
   22413.454658589766
```

```
1 from sklearn.model_selection import cross_val_score
2
3 forest_scores = cross_val_score(forest_reg, housing_prepared, housing_labels,
4                                 scoring="neg_mean_squared_error", cv=10)
5 forest_rmse_scores = np.sqrt(-forest_scores)
6 display_scores(forest_rmse_scores)
```

```
   Scores: [53519.05518628 50467.33817051 48924.16513902 53771.72056856
    50810.90996358 54876.09682033 56012.79985518 52256.88927227
    51527.73185039 55762.56008531]
   Mean: 52792.92669114079
   Standard deviation: 2262.8151900582
```

```
1 scores = cross_val_score(lin_reg, housing_prepared, housing_labels, scoring="neg_mean_squared_error", cv=10)
2 pd.Series(np.sqrt(-scores)).describe()
```

```
   count       10.000000
   mean     69104.079982
   std       3036.132517
   min      64114.991664
   25%      67077.398482
   50%      68718.763507
   75%      71357.022543
   max      73997.080502
   dtype: float64
```

```
1 from sklearn.svm import SVR
2
3 svm_reg = SVR(kernel="linear")
4 svm_reg.fit(housing_prepared, housing_labels)
5 housing_predictions = svm_reg.predict(housing_prepared)
6 svm_mse = mean_squared_error(housing_labels, housing_predictions)
7 svm_rmse = np.sqrt(svm_mse)
8 svm_rmse
```

```
   111095.06635291968
```

```
1 from sklearn.model_selection import GridSearchCV
2
3 param_grid = [
4     # try 12 (3×4) combinations of hyperparameters
5     {'n_estimators': [3, 10, 30], 'max_features': [2, 4, 6, 8]},
6     # then try 6 (2×3) combinations with bootstrap set as False
7     {'bootstrap': [False], 'n_estimators': [3, 10], 'max_features': [2, 3, 4]},
8   ]
9
10 forest_reg = RandomForestRegressor(random_state=42)
11 # train across 5 folds, that's a total of (12+6)*5=90 rounds of training
12 grid_search = GridSearchCV(forest_reg, param_grid, cv=5,
13                            scoring='neg_mean_squared_error', return_train_score=True)
14 grid_search.fit(housing_prepared, housing_labels)
```

```
   GridSearchCV(cv=5, estimator=RandomForestRegressor(random_state=42),
                param_grid=[{'max_features': [2, 4, 6, 8],
                             'n_estimators': [3, 10, 30]},
                            {'bootstrap': [False], 'max_features': [2, 3, 4],
                             'n_estimators': [3, 10]}],
                return_train_score=True, scoring='neg_mean_squared_error')
```

The best hyperparameter combination found:

```
1 grid_search.best_params_
```

```
   {'max_features': 8, 'n_estimators': 30}
```

```
1 grid_search.best_estimator_
```

```
RandomForestRegressor(max_features=8, n_estimators=30, random_state=42)
```

Let's look at the score of each hyperparameter combination tested during the grid search:

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

```
63895.161577951665 {'max_features': 2, 'n_estimators': 3}
54916.32386349543 {'max_features': 2, 'n_estimators': 10}
52885.86715332332 {'max_features': 2, 'n_estimators': 30}
60075.3680329983 {'max_features': 4, 'n_estimators': 3}
52495.01284985185 {'max_features': 4, 'n_estimators': 10}
50187.24324926565 {'max_features': 4, 'n_estimators': 30}
58064.73529982314 {'max_features': 6, 'n_estimators': 3}
51519.32062366315 {'max_features': 6, 'n_estimators': 10}
49969.80441627874 {'max_features': 6, 'n_estimators': 30}
58895.824998155826 {'max_features': 8, 'n_estimators': 3}
52459.79624724529 {'max_features': 8, 'n_estimators': 10}
49898.98913455217 {'max_features': 8, 'n_estimators': 30}
62381.765106921855 {'bootstrap': False, 'max_features': 2, 'n_estimators': 3}
54476.57050944266 {'bootstrap': False, 'max_features': 2, 'n_estimators': 10}
59974.60028085155 {'bootstrap': False, 'max_features': 3, 'n_estimators': 3}
52754.5632813202 {'bootstrap': False, 'max_features': 3, 'n_estimators': 10}
57831.136061214274 {'bootstrap': False, 'max_features': 4, 'n_estimators': 3}
51278.37877140253 {'bootstrap': False, 'max_features': 4, 'n_estimators': 10}
```

```
1 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 |
|---|---|---|---|---|---|---|---|
| 0 | 0.070521 | 0.001859 | 0.005390 | 0.000674 | 2 | 3 | NaN |
| 1 | 0.231319 | 0.012147 | 0.013777 | 0.000302 | 2 | 10 | NaN |
| 2 | 0.676404 | 0.012335 | 0.039669 | 0.003132 | 2 | 30 | NaN |
| 3 | 0.112226 | 0.001885 | 0.004928 | 0.000053 | 4 | 3 | NaN |
| 4 | 0.369924 | 0.008944 | 0.015067 | 0.002453 | 4 | 10 | NaN |
| 5 | 1.255543 | 0.214540 | 0.038651 | 0.001850 | 4 | 30 | NaN |
| 6 | 0.147990 | 0.001848 | 0.005060 | 0.000122 | 6 | 3 | NaN |
| 7 | 0.500964 | 0.003588 | 0.013439 | 0.000229 | 6 | 10 | NaN |

```
1 from sklearn.model_selection import RandomizedSearchCV
2 from scipy.stats import randint
3
4 param_distribs = {
5         'n_estimators': randint(low=1, high=200),
6         'max_features': randint(low=1, high=8),
7     }
8
9 forest_reg = RandomForestRegressor(random_state=42)
10 rnd_search = RandomizedSearchCV(forest_reg, param_distributions=param_distribs,
11                                 n_iter=10, cv=5, scoring='neg_mean_squared_error', random_state=42)
12 rnd_search.fit(housing_prepared, housing_labels)

    RandomizedSearchCV(cv=5, estimator=RandomForestRegressor(random_state=42),
                    param_distributions={'max_features': <scipy.stats._distn_infrastructure.rv_frozen object at
        0x7f86bca20fd0>,
                                        'n_estimators': <scipy.stats._distn_infrastructure.rv_frozen object at
        0x7f86bce50130>},
                    random_state=42, scoring='neg_mean_squared_error')
```

```
1 cvres = rnd_search.cv_results_
2 for mean_score, params in zip(cvres["mean_test_score"], cvres["params"]):
3     print(np.sqrt(-mean_score), params)

    49117.55344336652 {'max_features': 7, 'n_estimators': 180}
    51450.63202856348 {'max_features': 5, 'n_estimators': 15}
    50692.53588182537 {'max_features': 3, 'n_estimators': 72}
    50783.614493515 {'max_features': 5, 'n_estimators': 21}
    49162.89877456354 {'max_features': 7, 'n_estimators': 122}
    50655.798471042704 {'max_features': 3, 'n_estimators': 75}
    50513.856319990606 {'max_features': 3, 'n_estimators': 88}
    49521.17201976928 {'max_features': 5, 'n_estimators': 100}
    50302.90440763418 {'max_features': 3, 'n_estimators': 150}
    65167.02018649492 {'max_features': 5, 'n_estimators': 2}
```

```
1 feature_importances = grid_search.best_estimator_.feature_importances_
2 feature_importances
```

```
array([6.96542523e-02, 6.04213840e-02, 4.21882202e-02, 1.52450557e-02,
       1.55545295e-02, 1.58491147e-02, 1.49346552e-02, 3.79009225e-01,
       5.47789150e-02, 1.07031322e-01, 4.82031213e-02, 6.79266007e-03,
       1.65706303e-01, 7.83480660e-05, 1.52473276e-03, 3.02816106e-03])
```

```
1 extra_attribs = ["rooms_per_hhold", "pop_per_hhold", "bedrooms_per_room"]
2 #cat_encoder = cat_pipeline.named_steps["cat_encoder"] # old solution
3 cat_encoder = full_pipeline.named_transformers_["cat"]
4 cat_one_hot_attribs = list(cat_encoder.categories_[0])
5 attributes = num_attribs + extra_attribs + cat_one_hot_attribs
6 sorted(zip(feature_importances, attributes), reverse=True)
```

```
[(0.3790092248170967, 'median_income'),
 (0.16570630316895876, 'INLAND'),
 (0.10703132208204354, 'pop_per_hhold'),
 (0.06965425227942929, 'longitude'),
 (0.0604213440080722, 'latitude'),
 (0.054778915018283726, 'rooms_per_hhold'),
 (0.048203121338269206, 'bedrooms_per_room'),
 (0.04218822024391753, 'housing_median_age'),
 (0.015849114744428634, 'population'),
 (0.015554529490469328, 'total_bedrooms'),
 (0.01524505568840977, 'total_rooms'),
 (0.014934655161887776, 'households'),
 (0.006792660074259966, '<1H OCEAN'),
 (0.0030281610628962747, 'NEAR OCEAN'),
 (0.0015247327555504937, 'NEAR BAY'),
 (7.834806602687504e-05, 'ISLAND')]
```

```
 1 final_model = grid_search.best_estimator_
 2
 3 X_test = strat_test_set.drop("median_house_value", axis=1)
 4 y_test = strat_test_set["median_house_value"].copy()
 5
 6 X_test_prepared = full_pipeline.transform(X_test)
 7 final_predictions = final_model.predict(X_test_prepared)
 8
 9 final_mse = mean_squared_error(y_test, final_predictions)
10 final_rmse = np.sqrt(final_mse)
```

```
1 final_rmse
```

```
47873.26095812988
```

We can compute a 95% confidence interval for the test RMSE:

```
1 from scipy import stats
```

```
1 confidence = 0.95
2 squared_errors = (final_predictions - y_test) ** 2
3 mean = squared_errors.mean()
4 m = len(squared_errors)
5
6 np.sqrt(stats.t.interval(confidence, m - 1,
7                          loc=np.mean(squared_errors),
8                          scale=stats.sem(squared_errors)))
```

```
array([45893.36082829, 49774.46796717])
```

We could compute the interval manually like this:

```
1 tscore = stats.t.ppf((1 + confidence) / 2, df=m - 1)
2 tmargin = tscore * squared_errors.std(ddof=1) / np.sqrt(m)
3 np.sqrt(mean - tmargin), np.sqrt(mean + tmargin)
```

```
(45893.360828285535, 49774.46796717361)
```

Alternatively, we could use a z-scores rather than t-scores:

```
1 zscore = stats.norm.ppf((1 + confidence) / 2)
2 zmargin = zscore * squared_errors.std(ddof=1) / np.sqrt(m)
```

```
3 np.sqrt(mean - zmargin), np.sqrt(mean + zmargin)
   (45893.9540110131, 49773.921030650374)
```

## ▾ Extra material

## ▾ A full pipeline with both preparation and prediction

```
1 full_pipeline_with_predictor = Pipeline([
2        ("preparation", full_pipeline),
3        ("linear", LinearRegression())
4    ])
5
6 full_pipeline_with_predictor.fit(housing, housing_labels)
7 full_pipeline_with_predictor.predict(some_data)
```

```
   array([ 85657.90192014, 305492.60737488, 152056.46122456, 186095.70946094,
           244550.67966089])
```
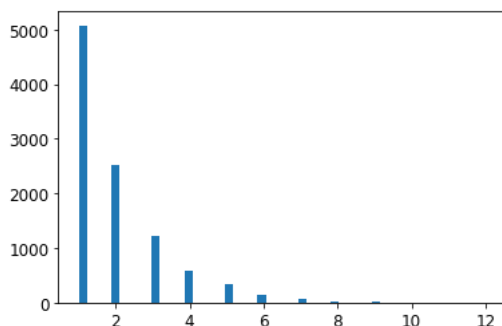
## ▾ Model persistence using joblib

```
1 my_model = full_pipeline_with_predictor
```

```
1 #from sklearn.externals import joblib # deprecated, use import joblib instead
2 import joblib
3
4 joblib.dump(my_model, "my_model.pkl") # DIFF
5 #...
6 my_model_loaded = joblib.load("my_model.pkl") # DIFF
```

## ▾ Example SciPy distributions for `RandomizedSearchCV`

```
1 from scipy.stats import geom, expon
2 geom_distrib=geom(0.5).rvs(10000, random_state=42)
3 expon_distrib=expon(scale=1).rvs(10000, random_state=42)
4 plt.hist(geom_distrib, bins=50)
5 plt.show()
6 plt.hist(expon_distrib, bins=50)
7 plt.show()
```

# ▾ Exercise solutions

## ▾ 1.

Question: Try a Support Vector Machine regressor (`sklearn.svm.SVR`), with various hyperparameters such as `kernel="linear"` (with various values for the `C` hyperparameter) or `kernel="rbf"` (with various values for the `C` and `gamma` hyperparameters). Don't worry about what these hyperparameters mean for now. How does the best `SVR` predictor perform?

```
1 from sklearn.model_selection import GridSearchCV
2
3 param_grid = [
4         {'kernel': ['linear'], 'C': [10., 30., 100., 300., 1000., 3000., 10000., 30000.0]},
5         {'kernel': ['rbf'], 'C': [1.0, 3.0, 10., 30., 100., 300., 1000.0],
6          'gamma': [0.01, 0.03, 0.1, 0.3, 1.0, 3.0]},
7     ]
8
9 svm_reg = SVR()
10 grid_search = GridSearchCV(svm_reg, param_grid, cv=5, scoring='neg_mean_squared_error', verbose=2, n_jobs=4)
11 grid_search.fit(housing_prepared, housing_labels)
```

```
    Fitting 5 folds for each of 50 candidates, totalling 250 fits
    GridSearchCV(cv=5, estimator=SVR(), n_jobs=4,
                 param_grid=[{'C': [10.0, 30.0, 100.0, 300.0, 1000.0, 3000.0,
                                    10000.0, 30000.0],
                              'kernel': ['linear']},
                             {'C': [1.0, 3.0, 10.0, 30.0, 100.0, 300.0, 1000.0],
                              'gamma': [0.01, 0.03, 0.1, 0.3, 1.0, 3.0],
                              'kernel': ['rbf']}],
                 scoring='neg_mean_squared_error', verbose=2)
```

The best model achieves the following score (evaluated using 5-fold cross validation):

```
1 negative_mse = grid_search.best_score_
2 rmse = np.sqrt(-negative_mse)
3 rmse
```

```
    70286.61835383571
```

That's much worse than the `RandomForestRegressor`. Let's check the best hyperparameters found:

```
1 grid_search.best_params_
```

```
    {'C': 30000.0, 'kernel': 'linear'}
```

The linear kernel seems better than the RBF kernel. Notice that the value of `C` is the maximum tested value. When this happens you definitely want to launch the grid search again with higher values for `C` (removing the smallest values), because it is likely that higher values of `C` will be better.

## ▾ 2.

Question: Try replacing `GridSearchCV` with `RandomizedSearchCV`.

```
1 from sklearn.model_selection import RandomizedSearchCV
2 from scipy.stats import expon, reciprocal
3
4 # see https://docs.scipy.org/doc/scipy/reference/stats.html
5 # for `expon()` and `reciprocal()` documentation and more probability distribution functions.
6
7 # Note: gamma is ignored when kernel is "linear"
8 param_distribs = {
9         'kernel': ['linear', 'rbf'],
10        'C': reciprocal(20, 200000),
11        'gamma': expon(scale=1.0),
```

```
12      }
13
14 svm_reg = SVR()
15 rnd_search = RandomizedSearchCV(svm_reg, param_distributions=param_distribs,
16                                 n_iter=50, cv=5, scoring='neg_mean_squared_error',
17                                 verbose=2, n_jobs=4, random_state=42)
18 rnd_search.fit(housing_prepared, housing_labels)
```

```
Fitting 5 folds for each of 50 candidates, totalling 250 fits
RandomizedSearchCV(cv=5, estimator=SVR(), n_iter=50, n_jobs=4,
                   param_distributions={'C': <scipy.stats._distn_infrastructure.rv_frozen object at 0x7f86b5b45970>,
                                        'gamma': <scipy.stats._distn_infrastructure.rv_frozen object at 0x7f86b5b45be0>,
                                        'kernel': ['linear', 'rbf']},
                   random_state=42, scoring='neg_mean_squared_error',
                   verbose=2)
```

The best model achieves the following score (evaluated using 5-fold cross validation):

```
1 negative_mse = rnd_search.best_score_
2 rmse = np.sqrt(-negative_mse)
3 rmse
```

```
54751.69009488048
```

Now this is much closer to the performance of the `RandomForestRegressor` (but not quite there yet). Let's check the best hyperparameters found:

```
1 rnd_search.best_params_
```

```
{'C': 157055.10989448498, 'gamma': 0.26497040005002437, 'kernel': 'rbf'}
```
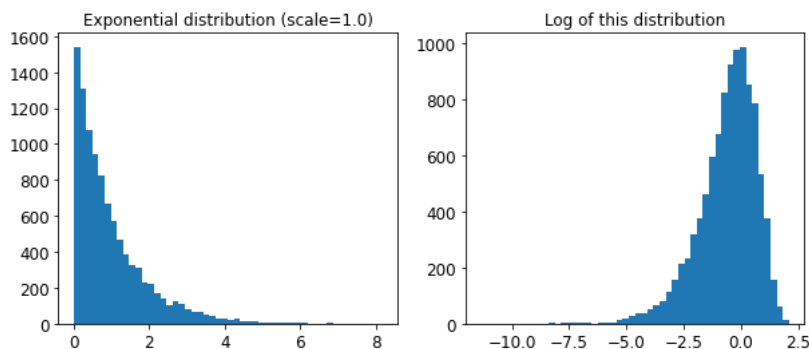
This time the search found a good set of hyperparameters for the RBF kernel. Randomized search tends to find better hyperparameters than grid search in the same amount of time.

Let's look at the exponential distribution we used, with `scale=1.0`. Note that some samples are much larger or smaller than 1.0, but when you look at the log of the distribution, you can see that most values are actually concentrated roughly in the range of exp(-2) to exp(+2), which is about 0.1 to 7.4.

```
 1 expon_distrib = expon(scale=1.)
 2 samples = expon_distrib.rvs(10000, random_state=42)
 3 plt.figure(figsize=(10, 4))
 4 plt.subplot(121)
 5 plt.title("Exponential distribution (scale=1.0)")
 6 plt.hist(samples, bins=50)
 7 plt.subplot(122)
 8 plt.title("Log of this distribution")
 9 plt.hist(np.log(samples), bins=50)
10 plt.show()
```
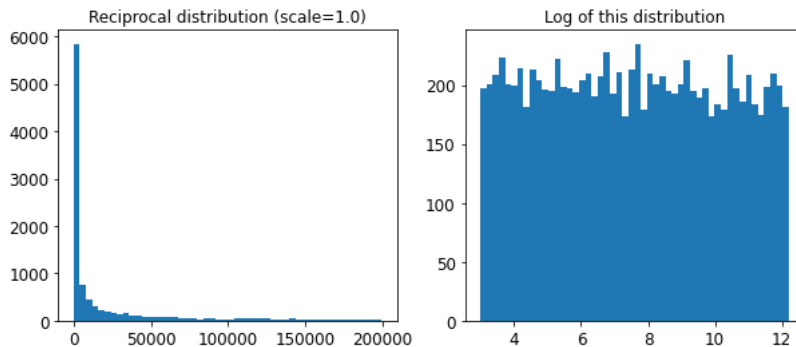


The distribution we used for `c` looks quite different: the scale of the samples is picked from a uniform distribution within a given range, which is why the right graph, which represents the log of the samples, looks roughly constant. This distribution is useful when you don't have a clue of what the target scale is:

```
 1 reciprocal_distrib = reciprocal(20, 200000)
 2 samples = reciprocal_distrib.rvs(10000, random_state=42)
 3 plt.figure(figsize=(10, 4))
 4 plt.subplot(121)
 5 plt.title("Reciprocal distribution (scale=1.0)")
 6 plt.hist(samples, bins=50)
 7 plt.subplot(122)
 8 plt.title("Log of this distribution")
 9 plt.hist(np.log(samples), bins=50)
10 plt.show()
```

The reciprocal distribution is useful when you have no idea what the scale of the hyperparameter should be (indeed, as you can see on the figure on the right, all scales are equally likely, within the given range), whereas the exponential distribution is best when you know (more or less) what the scale of the hyperparameter should be.

## ▾ 3.

Question: Try adding a transformer in the preparation pipeline to select only the most important attributes.

```
 1 from sklearn.base import BaseEstimator, TransformerMixin
 2
 3 def indices_of_top_k(arr, k):
 4     return np.sort(np.argpartition(np.array(arr), -k)[-k:])
 5
 6 class TopFeatureSelector(BaseEstimator, TransformerMixin):
 7     def __init__(self, feature_importances, k):
 8         self.feature_importances = feature_importances
 9         self.k = k
10     def fit(self, X, y=None):
11         self.feature_indices_ = indices_of_top_k(self.feature_importances, self.k)
12         return self
13     def transform(self, X):
14         return X[:, self.feature_indices_]
```

Note: this feature selector assumes that you have already computed the feature importances somehow (for example using a `RandomForestRegressor`). You may be tempted to compute them directly in the `TopFeatureSelector`'s `fit()` method, however this would likely slow down grid/randomized search since the feature importances would have to be computed for every hyperparameter combination (unless you implement some sort of cache).

Let's define the number of top features we want to keep:

```
 1 k = 5
```

Now let's look for the indices of the top k features:

```
 1 top_k_feature_indices = indices_of_top_k(feature_importances, k)
 2 top_k_feature_indices
```

```
   array([ 0,  1,  7,  9, 12])
```

```
1 np.array(attributes)[top_k_feature_indices]
```

```
array(['longitude', 'latitude', 'median_income', 'pop_per_hhold',
       'INLAND'], dtype='<U18')
```

Let's double check that these are indeed the top k features:

```
1 sorted(zip(feature_importances, attributes), reverse=True)[:k]
```

```
[(0.3790092248170967, 'median_income'),
 (0.16570630316895876, 'INLAND'),
 (0.10703132208204354, 'pop_per_hhold'),
 (0.06965425227942929, 'longitude'),
 (0.0604213840080722, 'latitude')]
```

Looking good... Now let's create a new pipeline that runs the previously defined preparation pipeline, and adds top k feature selection:

```
1 preparation_and_feature_selection_pipeline = Pipeline([
2     ('preparation', full_pipeline),
3     ('feature_selection', TopFeatureSelector(feature_importances, k))
4 ])
```

```
1 housing_prepared_top_k_features = preparation_and_feature_selection_pipeline.fit_transform(housing)
```

Let's look at the features of the first 3 instances:

```
1 housing_prepared_top_k_features[0:3]
```

```
array([[-0.94135046,  1.34743822, -0.8936472 ,  0.00622264,  1.        ],
       [ 1.17178212, -1.19243966,  1.292168  , -0.04081077,  0.        ],
       [ 0.26758118, -0.1259716 , -0.52543365, -0.07537122,  1.        ]])
```

Now let's double check that these are indeed the top k features:

```
1 housing_prepared[0:3, top_k_feature_indices]
```

```
array([[-0.94135046,  1.34743822, -0.8936472 ,  0.00622264,  1.        ],
       [ 1.17178212, -1.19243966,  1.292168  , -0.04081077,  0.        ],
       [ 0.26758118, -0.1259716 , -0.52543365, -0.07537122,  1.        ]])
```

Works great! :)

## 4.

Question: Try creating a single pipeline that does the full data preparation plus the final prediction.

```
1 prepare_select_and_predict_pipeline = Pipeline([
2     ('preparation', full_pipeline),
3     ('feature_selection', TopFeatureSelector(feature_importances, k)),
4     ('svm_reg', SVR(**rnd_search.best_params_))
5 ])
```

```
1 prepare_select_and_predict_pipeline.fit(housing, housing_labels)
```

```
Pipeline(steps=[('preparation',
                 ColumnTransformer(transformers=[('num',
                                                  Pipeline(steps=[('imputer',
                                                                   SimpleImputer(strategy='median')),
                                                                  ('attribs_adder',
                                                                   FunctionTransformer(func=<function add_extra_features at
       0x7f86bc86aca0>)),
                                                                  ('std_scaler',
                                                                   StandardScaler())]),
                                                  ['longitude', 'latitude',
                                                   'housing_median_age',
                                                   'total_rooms',
                                                   'total_bedrooms',
```

```
                                                'population', 'househ...
                  TopFeatureSelector(feature_importances=array([6.96542523e-02, 6.04213840e-02, 4.21882202e-02, 1.52450557e-
      02,
           1.55545295e-02, 1.58491147e-02, 1.49346552e-02, 3.79009225e-01,
           5.47789150e-02, 1.07031322e-01, 4.82031213e-02, 6.79266007e-03,
           1.65706303e-01, 7.83480660e-05, 1.52473276e-03, 3.02816106e-03]),
                                      k=5)),
                  ('svm_reg',
                   SVR(C=157055.10989448498, gamma=0.26497040005002437))])
```

Let's try the full pipeline on a few instances:

```
1 some_data = housing.iloc[:4]
2 some_labels = housing_labels.iloc[:4]
3
4 print("Predictions:\t", prepare_select_and_predict_pipeline.predict(some_data))
5 print("Labels:\t\t", list(some_labels))
```

```
Predictions:    [ 83384.49158095 299407.90439234  92272.03345144 150173.16199041]
Labels:         [72100.0, 279600.0, 82700.0, 112500.0]
```

Well, the full pipeline seems to work fine. Of course, the predictions are not fantastic: they would be better if we used the best
`RandomForestRegressor` that we found earlier, rather than the best `SVR`.

## ▾ 5.

Question: Automatically explore some preparation options using `GridSearchCV`.

```
1   param_grid = [{
2       'preparation__num__imputer__strategy': ['mean', 'median', 'most_frequent'],
3       'feature_selection__k': list(range(1, len(feature_importances) + 1))
4   }]
5
6   grid_search_prep = GridSearchCV(prepare_select_and_predict_pipeline, param_grid, cv=5,
7                                   scoring='neg_mean_squared_error', verbose=2, n_jobs=4)
8   grid_search_prep.fit(housing, housing_labels)
```

```
Fitting 5 folds for each of 48 candidates, totalling 240 fits
/usr/local/lib/python3.8/dist-packages/sklearn/model_selection/_validation.py:372: FitFailedWarning:
9 fits failed out of a total of 240.
The score on these train-test partitions for these parameters will be set to nan.
If these failures are not expected, you can try to debug them by setting error_score='raise'.

Below are more details about the failures:
--------------------------------------------------------------------------
9 fits failed with the following error:
Traceback (most recent call last):
  File "/usr/local/lib/python3.8/dist-packages/sklearn/model_selection/_validation.py", line 680, in _fit_and_score
    estimator.fit(X_train, y_train, **fit_params)
  File "/usr/local/lib/python3.8/dist-packages/sklearn/pipeline.py", line 390, in fit
    Xt = self._fit(X, y, **fit_params_steps)
  File "/usr/local/lib/python3.8/dist-packages/sklearn/pipeline.py", line 348, in _fit
    X, fitted_transformer = fit_transform_one_cached(
  File "/usr/local/lib/python3.8/dist-packages/joblib/memory.py", line 349, in __call__
    return self.func(*args, **kwargs)
  File "/usr/local/lib/python3.8/dist-packages/sklearn/pipeline.py", line 893, in _fit_transform_one
    res = transformer.fit_transform(X, y, **fit_params)
  File "/usr/local/lib/python3.8/dist-packages/sklearn/base.py", line 855, in fit_transform
    return self.fit(X, y, **fit_params).transform(X)
  File "<ipython-input-126-6a801ecaa128>", line 14, in transform
IndexError: index 15 is out of bounds for axis 1 with size 15

  warnings.warn(some_fits_failed_message, FitFailedWarning)
/usr/local/lib/python3.8/dist-packages/sklearn/model_selection/_search.py:969: UserWarning: One or more of the test scores a
 nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan
 nan nan nan nan nan nan nan nan nan nan nan nan]
  warnings.warn(
GridSearchCV(cv=5,
             estimator=Pipeline(steps=[('preparation',
                                        ColumnTransformer(transformers=[('num',
                                                                         Pipeline(steps=[('imputer',
                                                                                          SimpleImputer(strategy='median')),
                                                                                         ('attribs_adder',
                                                                                          FunctionTransformer(func=<function
      add_extra_features at 0x7f86bc86aca0>)),
```

```
                                                                    ('std_scaler',
                                                                     StandardScaler())]),
                                              ['longitude',
                                               'latitude',
                                               'housing_median_age',
                                               'total_rooms',
                                               'total_be...
                   5.47789150e-02, 1.07031322e-01, 4.82031213e-02, 6.79266007e-03,
                   1.65706303e-01, 7.83480660e-05, 1.52473276e-03, 3.02816106e-03]),
                                                     k=5)),
                              ('svm_reg',
                               SVR(C=157055.10989448498,
                                   gamma=0.26497040005002437))]),
           n_jobs=4,
           param_grid=[{'feature_selection__k': [1, 2, 3, 4, 5, 6, 7, 8, 9,
                                                 10, 11, 12, 13, 14, 15, 16],
                        'preparation__num__imputer__strategy': ['mean',
                                                                'median',
                                                                'most_frequent']}],
           scoring='neg mean squared error', verbose=2)
```

```
1   grid_search_prep.best_params_
```

```
{'feature_selection__k': 1, 'preparation__num__imputer__strategy': 'mean'}
```

The best imputer strategy is `mean`

```
1
```

✓ 0s    completed at 3:39 PM                                                                    ● ✕