

** End-to-end Machine Learning project**

- ✓ **housing_Price_Prediction**

- ✓ Download the Data

```
1 from pathlib import Path
2 import pandas as pd
3 import tarfile
4 import urllib.request
5
6 def load_housing_data():
7     tarball_path = Path("datasets/housing.tgz")
8     if not tarball_path.is_file():
9         Path("datasets").mkdir(parents=True, exist_ok=True)
10        url = "https://github.com/ageron/data/raw/main/housing.tgz"
11        urllib.request.urlretrieve(url, tarball_path)
12        with tarfile.open(tarball_path) as housing_tarball:
13            housing_tarball.extractall(path="datasets")
14        return pd.read_csv(Path("datasets/housing/housing.csv"))
15
16 housing = load_housing_data()
```

- ✓ Take a Quick Look at the Data Structure

```
1 housing.head()
```



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



Next steps:

[Generate code with housing](#)

[View recommended plots](#)

[New interactive sheet](#)

```
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()
```



count

ocean_proximity

<1H OCEAN	9136
INLAND	6551
NEAR OCEAN	2658
NEAR BAY	2290
ISLAND	5

dtype: int64

1 housing.describe()



	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_hou
count	20640.000000	20640.000000	20640.000000	20640.000000	20433.000000	20640.000000	20640.000000	20640.000000	2064
mean	-119.569704	35.631861	28.639486	2635.763081	537.870553	1425.476744	499.539680	3.870671	20685
std	2.003532	2.135952	12.585558	2181.615252	421.385070	1132.462122	382.329753	1.899822	11539
min	-124.350000	32.540000	1.000000	2.000000	1.000000	3.000000	1.000000	0.499900	1495
25%	-121.800000	33.930000	18.000000	1447.750000	296.000000	787.000000	280.000000	2.563400	1196C
50%	-118.490000	34.260000	29.000000	2127.000000	435.000000	1166.000000	409.000000	3.534800	1797C
75%	-118.010000	37.710000	37.000000	3148.000000	647.000000	1725.000000	605.000000	4.743250	26472
max	-114.310000	41.950000	52.000000	39320.000000	6445.000000	35682.000000	6082.000000	15.000100	5000C



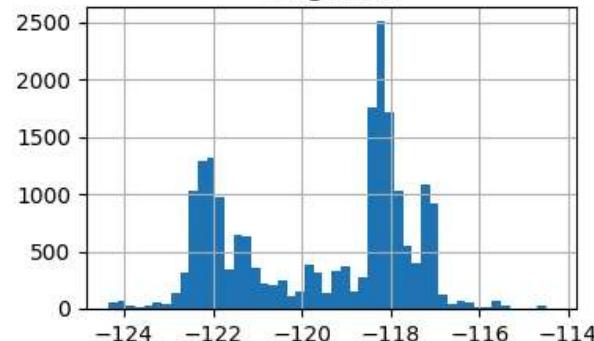
The following cell is not shown either in the book. It creates the `images/end_to_end_project` folder (if it doesn't already exist), and it defines the `save_fig()` function which is used through this notebook to save the figures in high-res for the book.

```
1 # extra code - code to save the figures as high-res PNGs for the book
2
3 IMAGES_PATH = Path() / "images" / "end_to_end_project"
4 IMAGES_PATH.mkdir(parents=True, exist_ok=True)
5
6 def save_fig(fig_id, tight_layout=True, fig_extension="png", resolution=300):
7     path = IMAGES_PATH / f"{fig_id}.{fig_extension}"
8     if tight_layout:
9         plt.tight_layout()
10    plt.savefig(path, format=fig_extension, dpi=resolution)

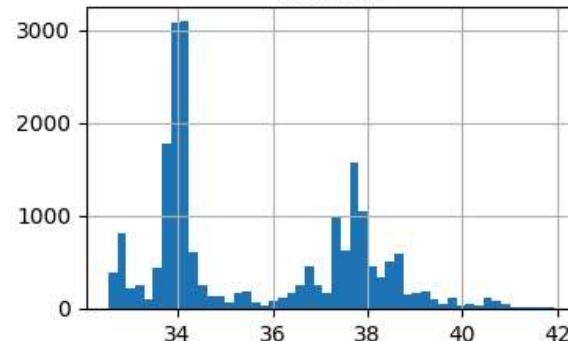
1 import matplotlib.pyplot as plt
2
3 # extra code - the next 5 lines define the default font sizes
4 plt.rc('font', size=14)
5 plt.rc('axes', labelsize=14, titlesize=14)
6 plt.rc('legend', fontsize=14)
7 plt.rc('xtick', labelsize=10)
8 plt.rc('ytick', labelsize=10)
9
10 housing.hist(bins=50, figsize=(12, 8))
11 save_fig("attribute_histogram_plots") # extra code
12 plt.show()
```



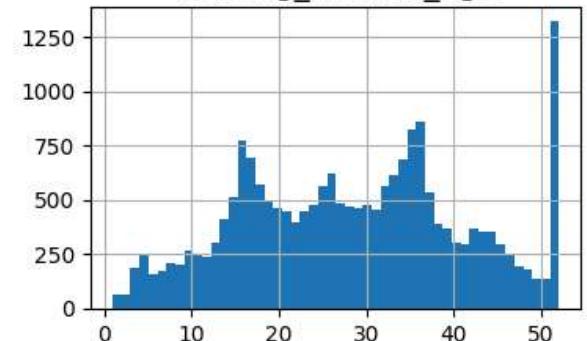
longitude



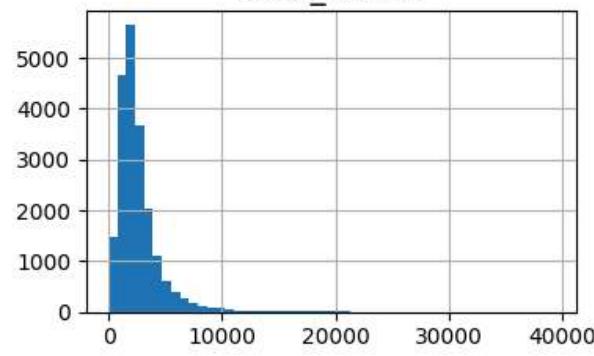
latitude



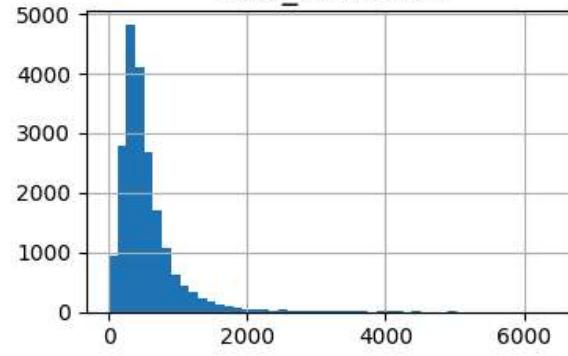
housing_median_age



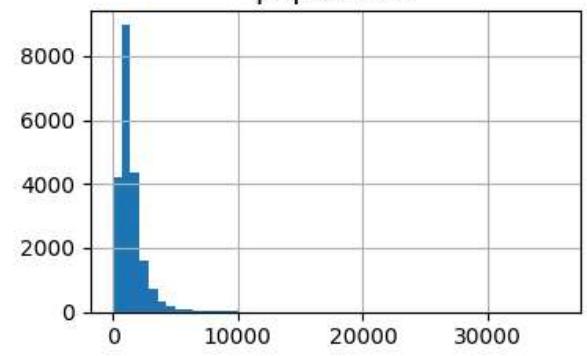
total_rooms



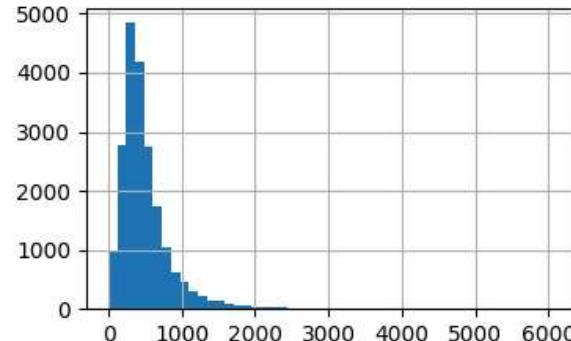
total_bedrooms



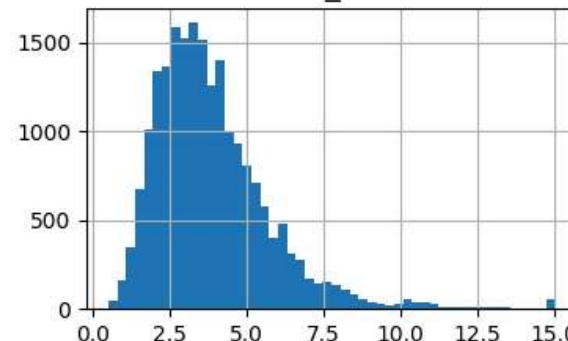
population



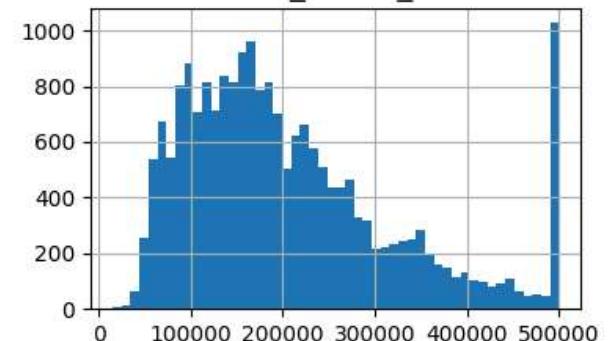
households



median_income



median_house_value



▼ Create a Test Set

```
1 import numpy as np
2
3 def shuffle_and_split_data(data, test_ratio):
4     shuffled_indices = np.random.permutation(len(data))
5     test_set_size = int(len(data) * test_ratio)
6     test_indices = shuffled_indices[:test_set_size]
7     train_indices = shuffled_indices[test_set_size:]
8     return data.iloc[train_indices], data.iloc[test_indices]
```

```
1 train_set, test_set = shuffle_and_split_data(housing, 0.2)
2 len(train_set)
```

→ 16512

```
1 len(test_set)
```

→ 4128

To ensure that this notebook's outputs remain the same every time we run it, we need to set the random seed:

```
1 np.random.seed(42)

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

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

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

```
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["total_bedrooms"].isnull().sum()
```

44

To find the probability that a random sample of 1,000 people contains less than 48.5% female or more than 53.5% female when the population's female ratio is 51.1%, we use the [binomial distribution](#). The `cdf()` method of the binomial distribution gives us the probability that the number of females will be equal or less than the given value.

```
1 # extra code - shows how to compute the 10.7% proba of getting a bad sample
2
3 from scipy.stats import binom
4
5 sample_size = 1000
6 ratio_female = 0.511
7 proba_too_small = binom(sample_size, ratio_female).cdf(485 - 1)
8 proba_too_large = 1 - binom(sample_size, ratio_female).cdf(535)
9 print(proba_too_small + proba_too_large)
```

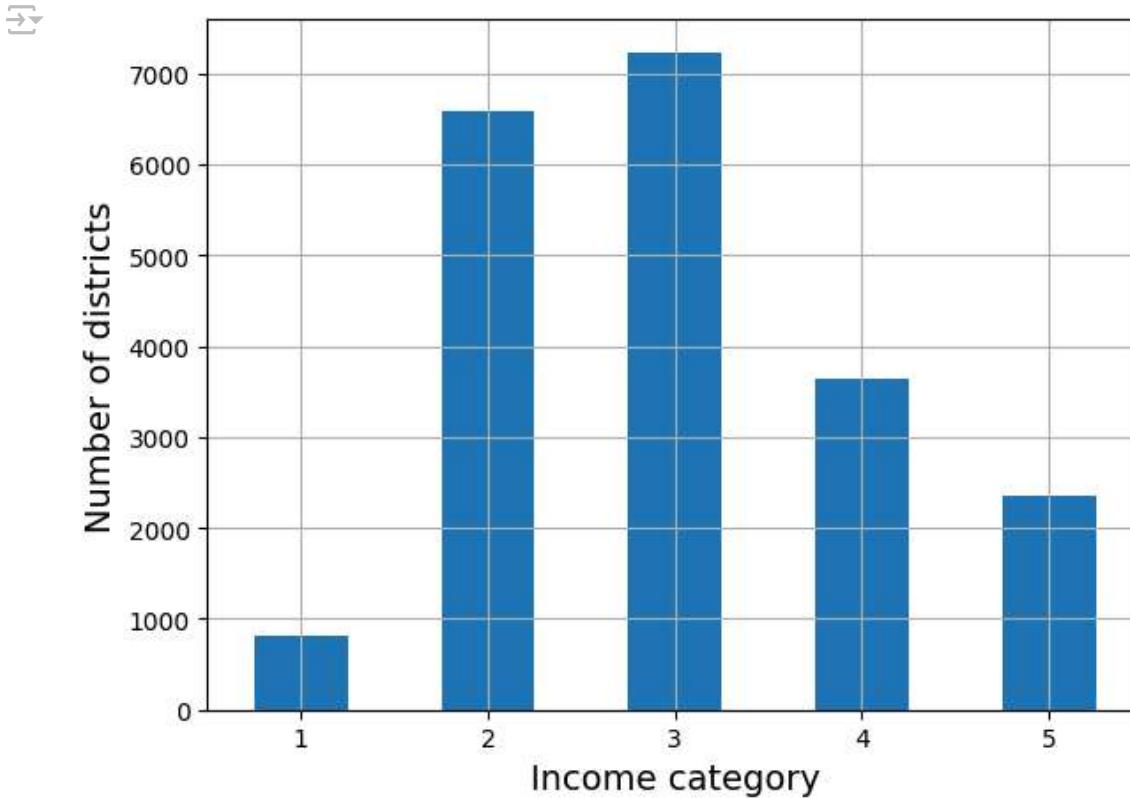
→ 0.10736798530929942

If you prefer simulations over maths, here's how you could get roughly the same result:

```
1 # extra code - shows another way to estimate the probability of bad sample
2
3 np.random.seed(42)
4
5 samples = (np.random.rand(100_000, sample_size) < ratio_female).sum(axis=1)
6 ((samples < 485) | (samples > 535)).mean()
```

→ 0,1071

```
1 housing["income_cat"].value_counts().sort_index().plot.bar(rot=0, grid=True)
2 plt.xlabel("Income category")
3 plt.ylabel("Number of districts")
4 save_fig("housing_income_cat_bar_plot") # extra code
5 plt.show()
```



```
1 from sklearn.model_selection import StratifiedShuffleSplit
2
3 splitter = StratifiedShuffleSplit(n_splits=10, test_size=0.2, random_state=42)
4 strat_splits = []
5 for train_index, test_index in splitter.split(housing, housing["income_cat"]):
6     strat_train_set_n = housing.iloc[train_index]
7     strat_test_set_n = housing.iloc[test_index]
8     strat_splits.append([strat_train_set_n, strat_test_set_n])
9
10 strat_train_set, strat_test_set = strat_splits[0]
```

It's much shorter to get a single stratified split:

```
1 strat_train_set, strat_test_set = train_test_split(  
2     housing, test_size=0.2, stratify=housing["income_cat"], random_state=42)
```

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

income_cat	count
3	0.350533
2	0.318798
4	0.176357
5	0.114341
1	0.039971

dtype: float64

```
1 # extra code - computes the data for Figure 2-10  
2  
3 def income_cat_proportions(data):  
4     return data["income_cat"].value_counts() / len(data)  
5  
6 train_set, test_set = train_test_split(housing, test_size=0.2, random_state=42)  
7  
8 compare_props = pd.DataFrame({  
9     "Overall %": income_cat_proportions(housing),  
10    "Stratified %": income_cat_proportions(strat_test_set),  
11    "Random %": income_cat_proportions(test_set),  
12 }).sort_index()  
13 compare_props.index.name = "Income Category"  
14 compare_props["Strat. Error %"] = (compare_props["Stratified %"] /  
15                                     compare_props["Overall %"] - 1)  
16 compare_props["Rand. Error %"] = (compare_props["Random %"] /  
17                                     compare_props["Overall %"] - 1)  
18 (compare_props * 100).round(2)
```



Overall % Stratified % Random % Strat. Error % Rand. Error %



Income Category

Income Category	Overall %	Stratified %	Random %	Strat. Error %	Rand. Error %
1	3.98	4.00	4.24	0.36	6.45
2	31.88	31.88	30.74	-0.02	-3.59
3	35.06	35.05	34.52	-0.01	-1.53
4	17.63	17.64	18.41	0.03	4.42
5	11.44	11.43	12.09	-0.08	5.63

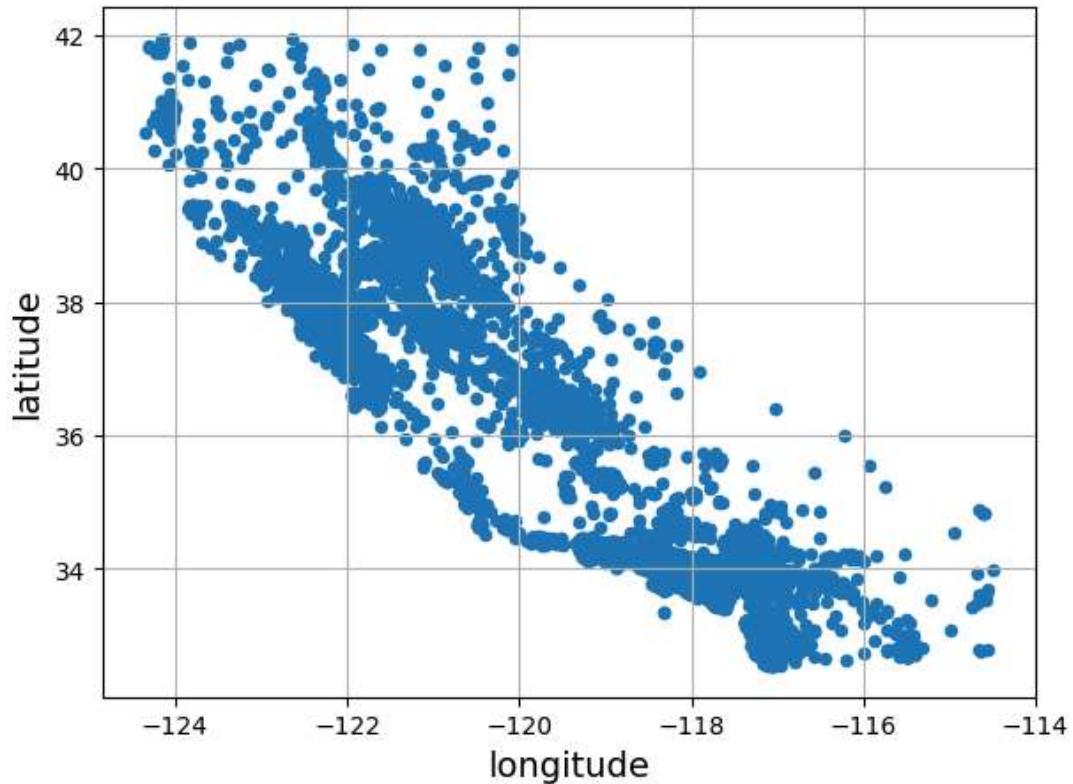
```
1 for set_ in (strat_train_set, strat_test_set):
2     set_.drop("income_cat", axis=1, inplace=True)
```

✓ Discover and Visualize the Data to Gain Insights

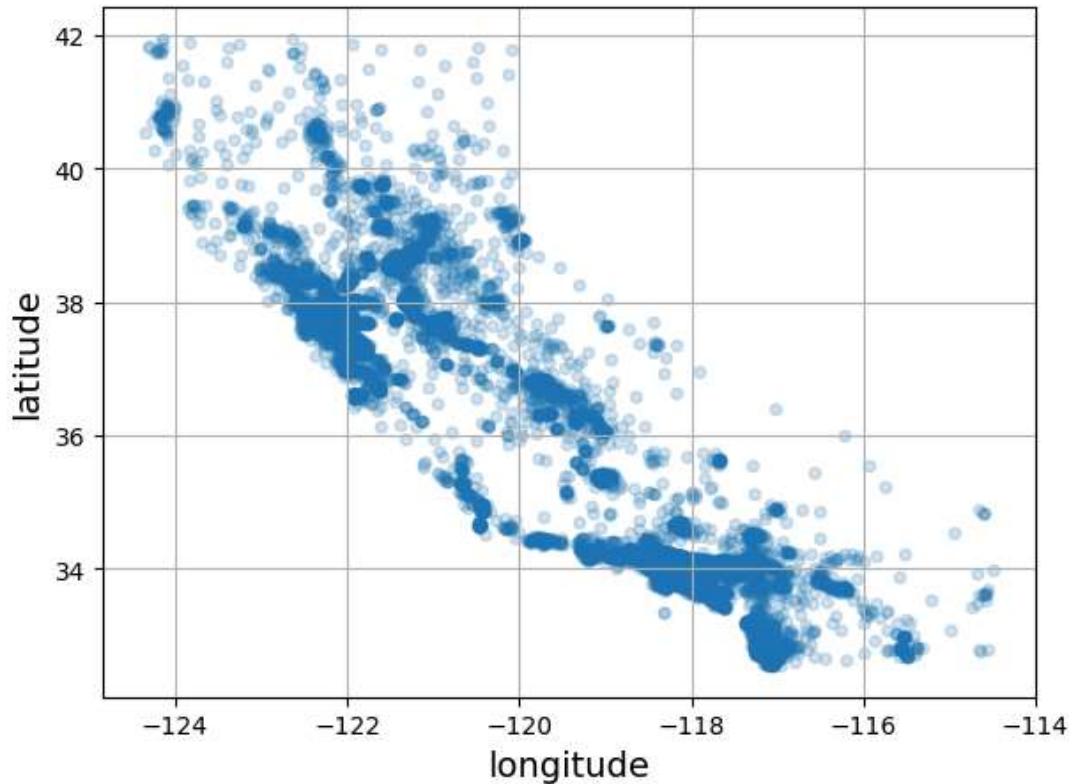
```
1 housing = strat_train_set.copy()
```

✓ Visualizing Geographical Data

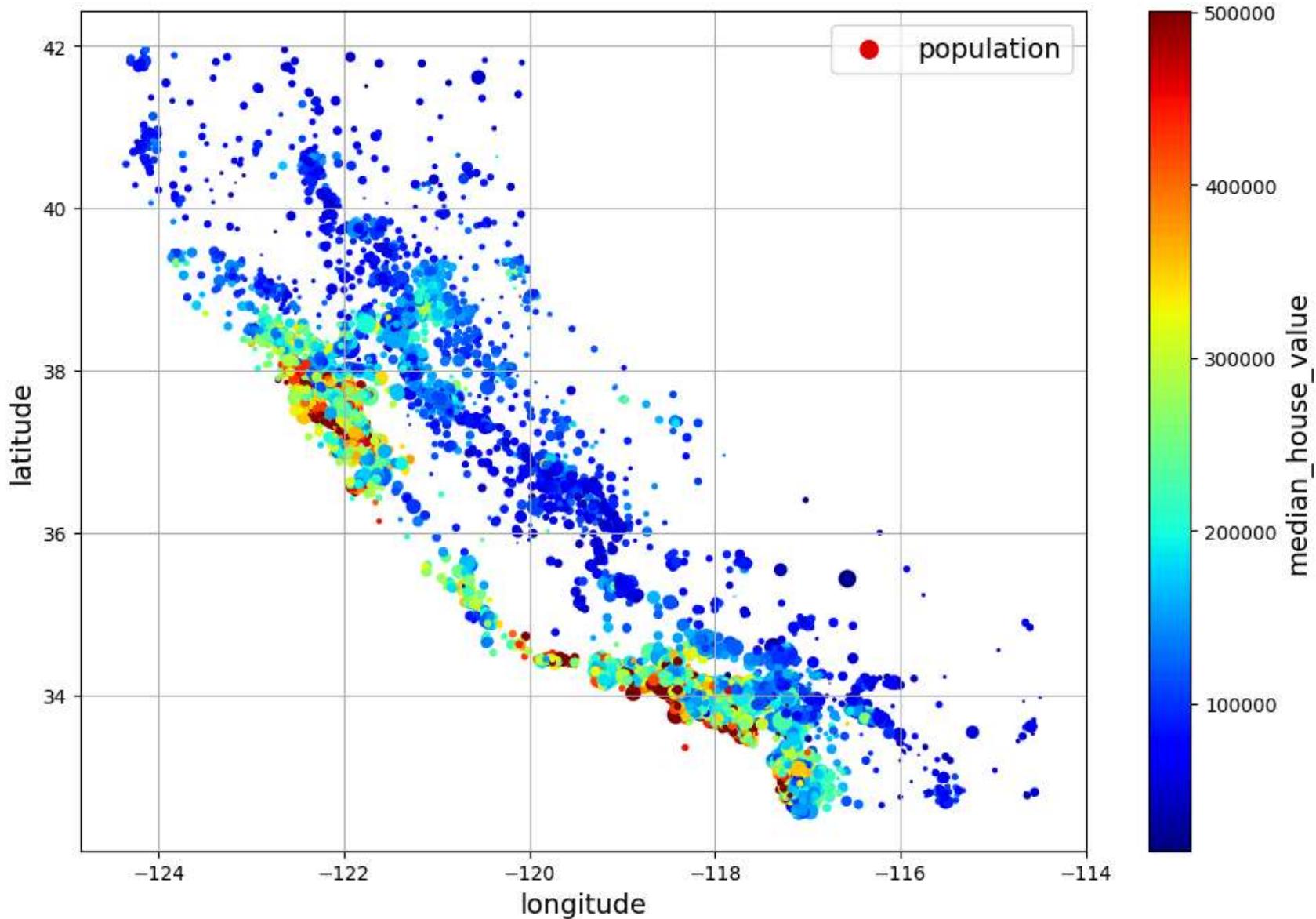
```
1 housing.plot(kind="scatter", x="longitude", y="latitude", grid=True)
2 save_fig("bad_visualization_plot") # extra code
3 plt.show()
```



```
1 housing.plot(kind="scatter", x="longitude", y="latitude", grid=True, alpha=0.2)
2 save_fig("better_visualization_plot") # extra code
3 plt.show()
```

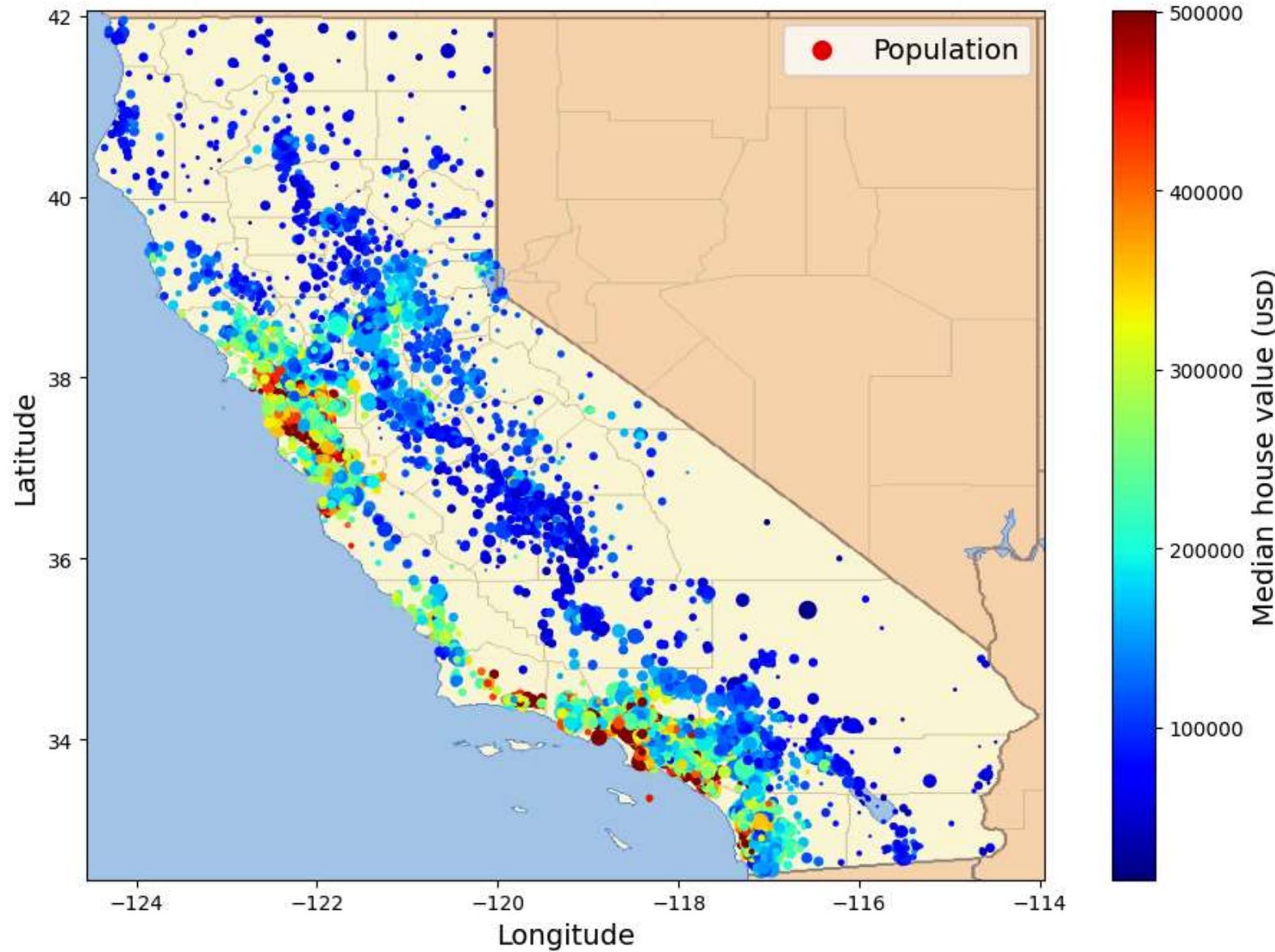


```
1 housing.plot(kind="scatter", x="longitude", y="latitude", grid=True,
2                 s=housing["population"] / 100, label="population",
3                 c="median_house_value", cmap="jet", colorbar=True,
4                 legend=True, sharex=False, figsize=(10, 7))
5 save_fig("housing_prices_scatterplot")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



```
1 # extra code - this cell generates the first figure in the chapter
2
3 # Download the California image
4 filename = "california.png"
5 if not (IMAGES_PATH / filename).is_file():
6     homl3_root = "https://github.com/ageron/handson-ml3/raw/main/"
7     url = homl3_root + "images/end_to_end_project/" + filename
8     print("Downloading", filename)
9     urllib.request.urlretrieve(url, IMAGES_PATH / filename)
10
11 housing_renamed = housing.rename(columns={
12     "latitude": "Latitude", "longitude": "Longitude",
13     "population": "Population",
14     "median_house_value": "Median house value (usd)"})
15 housing_renamed.plot(
16     kind="scatter", x="Longitude", y="Latitude",
17     s=housing_renamed["Population"] / 100, label="Population",
18     c="Median house value (usd)", cmap="jet", colorbar=True,
19     legend=True, sharex=False, figsize=(10, 7))
20
21 california_img = plt.imread(IMAGES_PATH / filename)
22 axis = -124.55, -113.95, 32.45, 42.05
23 plt.axis(axis)
24 plt.imshow(california_img, extent=axis)
25
26 save_fig("california_housing_prices_plot")
27 plt.show()
```

⤵ Downloading california.png



⌄ Looking for Correlations

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

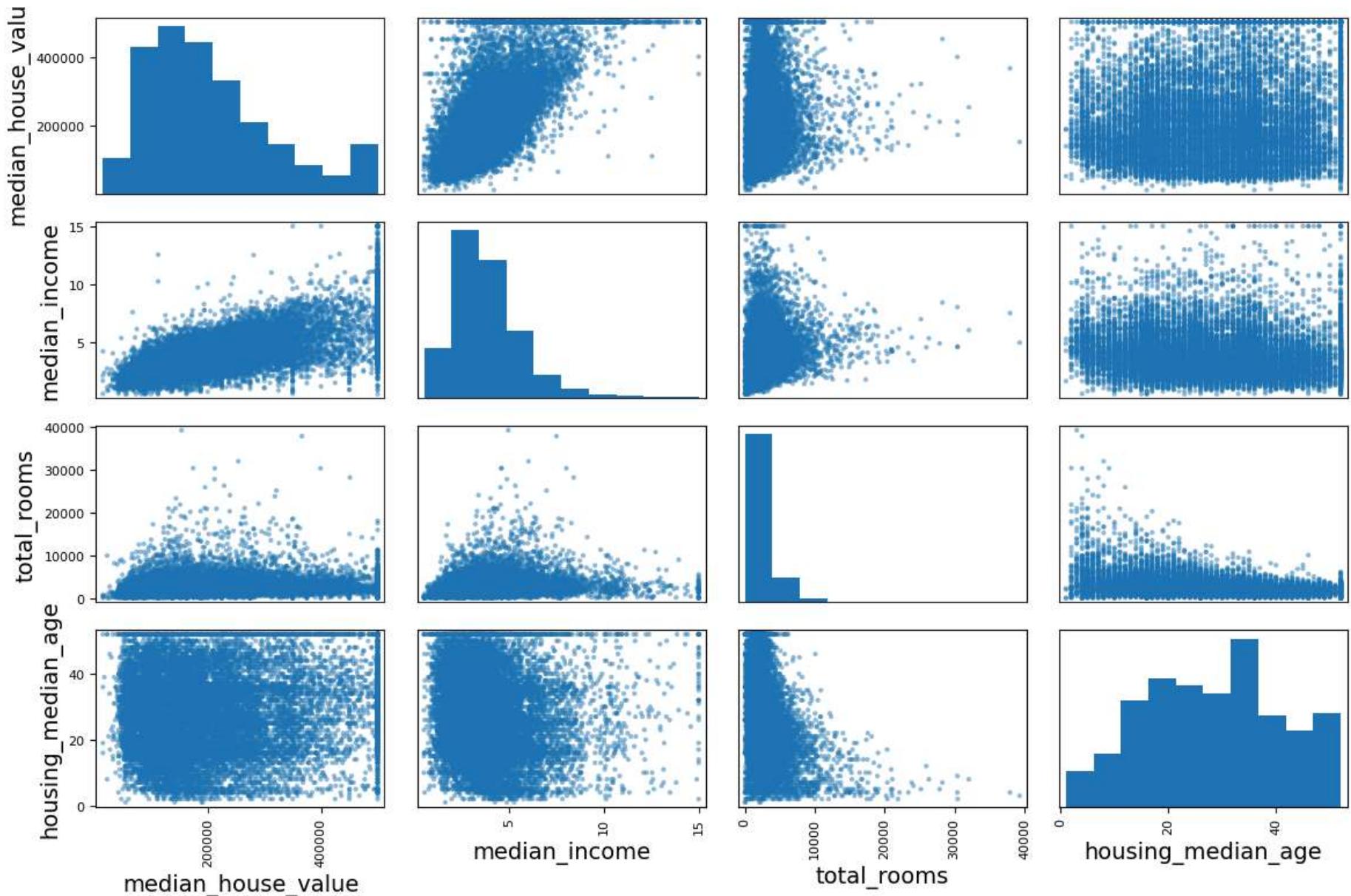
```
1 corr_matrix["median_house_value"].sort_values(ascending=False)
```

	median_house_value
median_house_value	1.000000
median_income	0.688380
total_rooms	0.137455
housing_median_age	0.102175
households	0.071426
total_bedrooms	0.054635
population	-0.020153
longitude	-0.050859
latitude	-0.139584

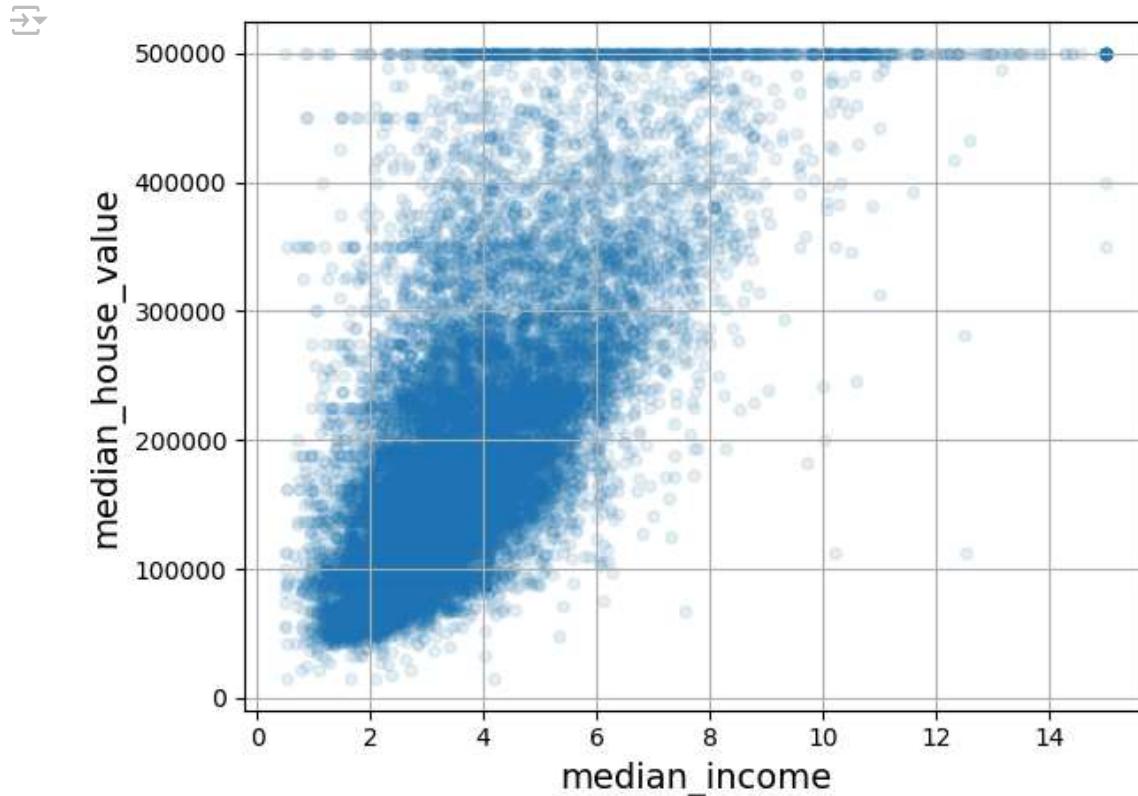
dtype: float64

```
1 from pandas.plotting import scatter_matrix
2
3 attributes = ["median_house_value", "median_income", "total_rooms",
4                 "housing_median_age"]
5 scatter_matrix(housing[attributes], figsize=(12, 8))
6 save_fig("scatter_matrix_plot") # extra code
7 plt.show()
```

[↓]



```
1 housing.plot(kind="scatter", x="median_income", y="median_house_value",
2                 alpha=0.1, grid=True)
3 save_fig("income_vs_house_value_scatterplot") # extra code
4 plt.show()
```



▼ Experimenting with Attribute Combinations

```
1 housing["rooms_per_house"] = housing["total_rooms"] / housing["households"]
2 housing["bedrooms_ratio"] = housing["total_bedrooms"] / housing["total_rooms"]
3 housing["people_per_house"] = housing["population"] / housing["households"]

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



	median_house_value
median_house_value	1.000000
median_income	0.688380
rooms_per_house	0.143663
total_rooms	0.137455
housing_median_age	0.102175
households	0.071426
total_bedrooms	0.054635
population	-0.020153
people_per_house	-0.038224
longitude	-0.050859
latitude	-0.139584
bedrooms_ratio	-0.256397

dtype: float64

▼ Prepare the Data for Machine Learning Algorithms

Let's revert to the original training set and separate the target (note that `strat_train_set.drop()` creates a copy of `strat_train_set` without the column, it doesn't actually modify `strat_train_set` itself, unless you pass `inplace=True`):

```
1 housing = strat_train_set.drop("median_house_value", axis=1)
2 housing_labels = strat_train_set["median_house_value"].copy()
```

▼ Data Cleaning

In the book 3 options are listed to handle the NaN values:

```

housing.dropna(subset=["total_bedrooms"], inplace=True)      # option 1

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

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

```

For each option, we'll create a copy of `housing` and work on that copy to avoid breaking `housing`. We'll also show the output of each option, but filtering on the rows that originally contained a `Nan` value.

```

1 null_rows_idx = housing.isnull().any(axis=1)
2 housing.loc>null_rows_idx].head()

```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	ocean_proximity
14452	-120.67	40.50	15.0	5343.0	NaN	2503.0	902.0	3.5962	INLAND
18217	-117.96	34.03	35.0	2093.0	NaN	1755.0	403.0	3.4115	<1H OCEAN
11889	-118.05	34.04	33.0	1348.0	NaN	1098.0	257.0	4.2917	<1H OCEAN
20325	-118.88	34.17	15.0	4260.0	NaN	1701.0	669.0	5.1033	<1H OCEAN
14360	-117.87	33.62	8.0	1266.0	NaN	375.0	183.0	9.8020	<1H OCEAN

```

1 housing_option1 = housing.copy()
2
3 housing_option1.dropna(subset=["total_bedrooms"], inplace=True)    # option 1
4
5 housing_option1.loc>null_rows_idx].head()

```

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

```

1 housing_option2 = housing.copy()
2
3 housing_option2.drop("total_bedrooms", axis=1, inplace=True)    # option 2
4
5 housing_option2.loc>null_rows_idx].head()

```

→

	longitude	latitude	housing_median_age	total_rooms	population	households	median_income	ocean_proximity
14452	-120.67	40.50	15.0	5343.0	2503.0	902.0	3.5962	INLAND
18217	-117.96	34.03	35.0	2093.0	1755.0	403.0	3.4115	<1H OCEAN
11889	-118.05	34.04	33.0	1348.0	1098.0	257.0	4.2917	<1H OCEAN
20325	-118.88	34.17	15.0	4260.0	1701.0	669.0	5.1033	<1H OCEAN
14360	-117.87	33.62	8.0	1266.0	375.0	183.0	9.8020	<1H OCEAN

←

```

1 housing_option3 = housing.copy()
2
3 median = housing["total_bedrooms"].median()
4 housing_option3["total_bedrooms"].fillna(median, inplace=True) # option 3
5
6 housing_option3.loc>null_rows_idx].head()

```

→

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	ocean_proximity
14452	-120.67	40.50	15.0	5343.0	434.0	2503.0	902.0	3.5962	INLAND
18217	-117.96	34.03	35.0	2093.0	434.0	1755.0	403.0	3.4115	<1H OCEAN
11889	-118.05	34.04	33.0	1348.0	434.0	1098.0	257.0	4.2917	<1H OCEAN
20325	-118.88	34.17	15.0	4260.0	434.0	1701.0	669.0	5.1033	<1H OCEAN
14360	-117.87	33.62	8.0	1266.0	434.0	375.0	183.0	9.8020	<1H OCEAN

← →

```

1 from sklearn.impute import SimpleImputer
2
3 imputer = SimpleImputer(strategy="median")

```

Separating out the numerical attributes to use the "median" strategy (as it cannot be calculated on text attributes like ocean_proximity):

```

1 housing_num = housing.select_dtypes(include=[np.number])
2
1 imputer.fit(housing_num)

```

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

```
1 imputer.statistics_
array([-118.51 ,  34.26 ,  29.    , 2125.    ,  434.    , 1167.    ,
       408.    ,  3.5385])
```

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.    , 2125.    ,  434.    , 1167.    ,
       408.    ,  3.5385])
```

Transform the training set:

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

1 imputer.feature_names_in_
array(['longitude', 'latitude', 'housing_median_age', 'total_rooms',
       'total_bedrooms', 'population', 'households', 'median_income'],
      dtype=object)

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

1 housing_tr.loc>null_rows_idx].head()
```

→

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income
14452	-120.67	40.50	15.0	5343.0	434.0	2503.0	902.0	3.5962
18217	-117.96	34.03	35.0	2093.0	434.0	1755.0	403.0	3.4115
11889	-118.05	34.04	33.0	1348.0	434.0	1098.0	257.0	4.2917
20325	-118.88	34.17	15.0	4260.0	434.0	1701.0	669.0	5.1033
14360	-117.87	33.62	8.0	1266.0	434.0	375.0	183.0	9.8020

```
1 imputer.strategy
```

→ 'median'

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

```
1 housing_tr.loc>null_rows_idx].head() # not shown in the book
```

→

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income
14452	-120.67	40.50	15.0	5343.0	434.0	2503.0	902.0	3.5962
18217	-117.96	34.03	35.0	2093.0	434.0	1755.0	403.0	3.4115
11889	-118.05	34.04	33.0	1348.0	434.0	1098.0	257.0	4.2917
20325	-118.88	34.17	15.0	4260.0	434.0	1701.0	669.0	5.1033
14360	-117.87	33.62	8.0	1266.0	434.0	375.0	183.0	9.8020



Now let's drop some outliers:

```
1 from sklearn.ensemble import IsolationForest
2
3 isolation_forest = IsolationForest(random_state=42)
4 outlier_pred = isolation_forest.fit_predict(X)
```

```
1 outlier_pred
```

```
→ array([-1,  1,  1, ...,  1,  1,  1])
```

If you wanted to drop outliers, you would run the following code:

```
1 #housing = housing.iloc[outlier_pred == 1]
2 #housing_labels = housing_labels.iloc[outlier_pred == 1]
```

✓ Handling Text and Categorical Attributes

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

```
1 housing_cat = housing[["ocean_proximity"]]
2 housing_cat.head(8)
```

→

	ocean_proximity	grid
13096	NEAR BAY	grid
14973	<1H OCEAN	
3785	INLAND	
14689	INLAND	
20507	NEAR OCEAN	
1286	INLAND	
18078	<1H OCEAN	
4396	NEAR BAY	

Next steps:

[Generate code with `housing_cat`](#)

[!\[\]\(d415b5172fecdbaea44b7ff6524f4d79_img.jpg\) View recommended plots](#)

[New interactive sheet](#)

```
1 from sklearn.preprocessing import OrdinalEncoder
2
3 ordinal_encoder = OrdinalEncoder()
4 housing_cat_encoded = ordinal_encoder.fit_transform(housing_cat)
```

```
1 housing_cat_encoded[:8]
2
3 array([[3.],
4        [0.],
5        [1.],
6        [1.],
7        [4.],
8        [1.],
9        [0.],
10       [3.]])  
  
1 ordinal_encoder.categories_
2
3 [array(['<1H OCEAN', 'INLAND', 'ISLAND', 'NEAR BAY', 'NEAR OCEAN'],
4       dtype=object)]  
  
1 from sklearn.preprocessing import OneHotEncoder
2
3 cat_encoder = OneHotEncoder()
4 housing_cat_1hot = cat_encoder.fit_transform(housing_cat)  
  
1 housing_cat_1hot
2
3 <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()
2
3 array([[0., 0., 0., 1., 0.],
4        [1., 0., 0., 0., 0.],
5        [0., 1., 0., 0., 0.],
6        ...,
7        [0., 0., 0., 0., 1.],
8        [1., 0., 0., 0., 0.],
9        [0., 0., 0., 0., 1.]])  
  
1 cat_encoder = OneHotEncoder(sparse_output=False)
2 housing_cat_1hot = cat_encoder.fit_transform(housing_cat)
3 housing_cat_1hot
```

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

```
1 cat_encoder.categories_
```

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

```
1 df_test = pd.DataFrame({"ocean_proximity": ["INLAND", "NEAR BAY"]})  
2 pd.get_dummies(df_test)
```

	ocean_proximity_INLAND	ocean_proximity_NEAR BAY
0	True	False
1	False	True

```
1 cat_encoder.transform(df_test)
```

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

```
1 df_test_unknown = pd.DataFrame({"ocean_proximity": ["<2H OCEAN", "ISLAND"]})  
2 pd.get_dummies(df_test_unknown)
```

	ocean_proximity_<2H_OCEAN	ocean_proximity_ISLAND
0	True	False
1	False	True

```
1 cat_encoder.handle_unknown = "ignore"  
2 cat_encoder.transform(df_test_unknown)
```

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

```
1 cat_encoder.feature_names_in_
2
3     ↗ array(['ocean_proximity'], dtype=object)

1 cat_encoder.get_feature_names_out()
2
3     ↗ array(['ocean_proximity_<1H OCEAN', 'ocean_proximity_INLAND',
4             'ocean_proximity_ISLAND', 'ocean_proximity_NEAR BAY',
5             'ocean_proximity_NEAR OCEAN'], dtype=object)

1 df_output = pd.DataFrame(cat_encoder.transform(df_test_unknown),
2                           columns=cat_encoder.get_feature_names_out(),
3                           index=df_test_unknown.index)
```

```
1 df_output
```

	ocean_proximity_<1H OCEAN	ocean_proximity_INLAND	ocean_proximity_ISLAND	ocean_proximity_NEAR BAY	ocean_proximity_NEAR OCEAN	
0	0.0	0.0	0.0	0.0	0.0	0.0
1	0.0	0.0	1.0	0.0	0.0	0.0

Next steps: [Generate code with df_output](#) [View recommended plots](#) [New interactive sheet](#)

▼ Feature Scaling

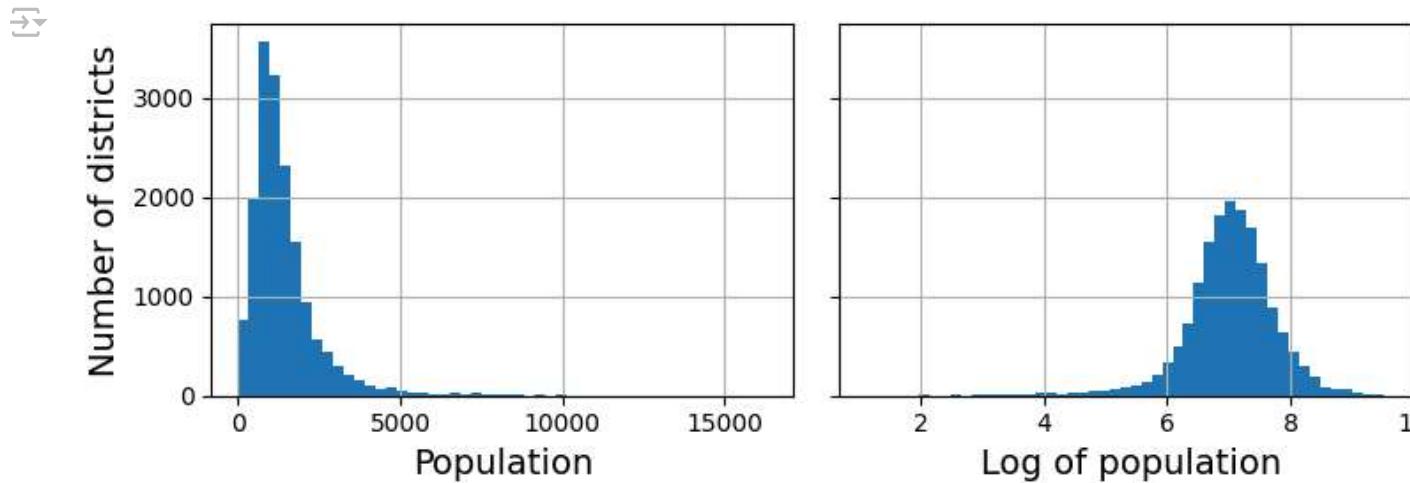
```
1 from sklearn.preprocessing import MinMaxScaler
2
3 min_max_scaler = MinMaxScaler(feature_range=(-1, 1))
4 housing_num_min_max_scaled = min_max_scaler.fit_transform(housing_num)
```

```
1 from sklearn.preprocessing import StandardScaler
2
3 std_scaler = StandardScaler()
4 housing_num_std_scaled = std_scaler.fit_transform(housing_num)
```

```

1 # extra code - this cell generates Figure 2-17
2 fig, axs = plt.subplots(1, 2, figsize=(8, 3), sharey=True)
3 housing["population"].hist(ax=axs[0], bins=50)
4 housing["population"].apply(np.log).hist(ax=axs[1], bins=50)
5 axs[0].set_xlabel("Population")
6 axs[1].set_xlabel("Log of population")
7 axs[0].set_ylabel("Number of districts")
8 save_fig("long_tail_plot")
9 plt.show()

```

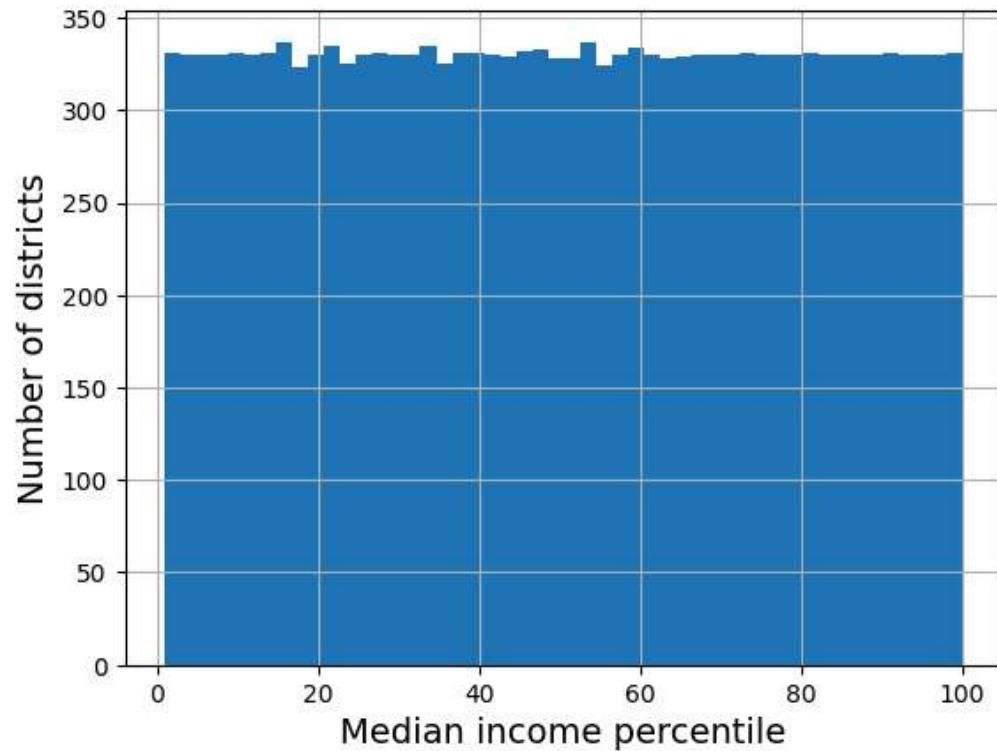


What if we replace each value with its percentile?

```

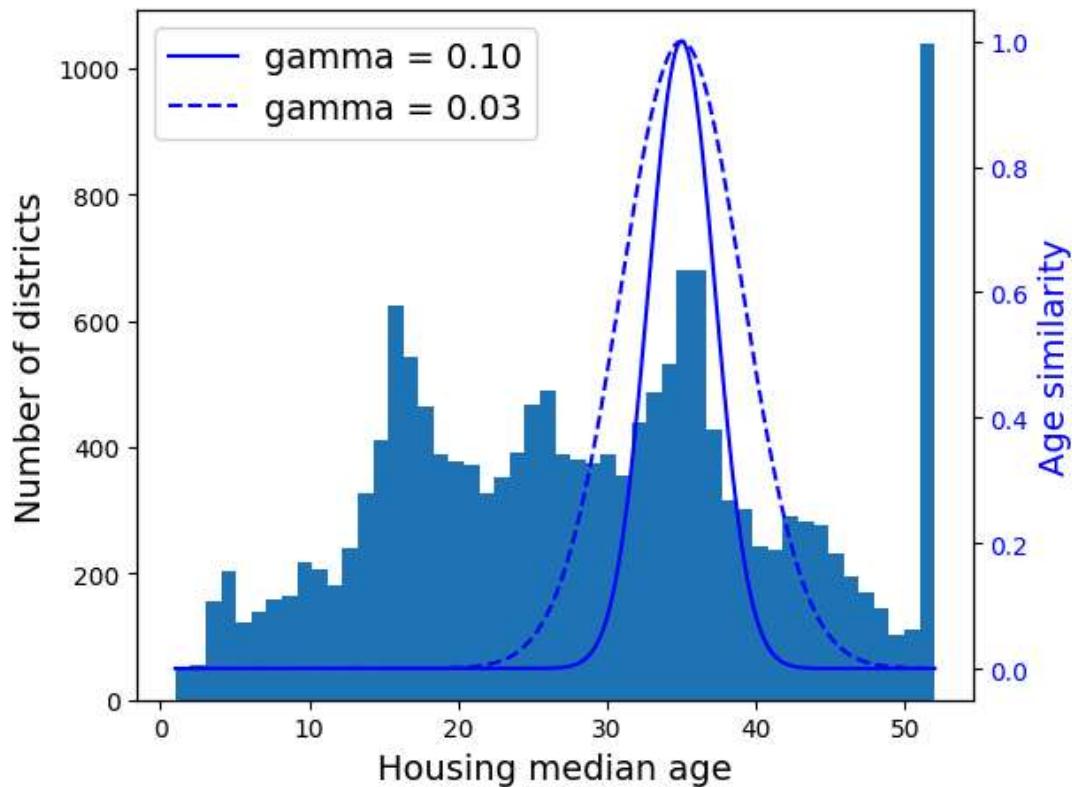
1 # extra code - just shows that we get a uniform distribution
2 percentiles = [np.percentile(housing["median_income"], p)
3                 for p in range(1, 100)]
4 flattened_median_income = pd.cut(housing["median_income"],
5                                   bins=[-np.inf] + percentiles + [np.inf],
6                                   labels=range(1, 100 + 1))
7 flattened_median_income.hist(bins=50)
8 plt.xlabel("Median income percentile")
9 plt.ylabel("Number of districts")
10 plt.show()
11 # Note: incomes below the 1st percentile are labeled 1, and incomes above the
12 # 99th percentile are labeled 100. This is why the distribution below ranges
13 # from 1 to 100 (not 0 to 100).

```



```
1 from sklearn.metrics.pairwise import rbf_kernel  
2  
3 age_simil_35 = rbf_kernel(housing[["housing_median_age"]], [[35]], gamma=0.1)
```

```
1 # extra code - this cell generates Figure 2-18
2
3 ages = np.linspace(housing["housing_median_age"].min(),
4                     housing["housing_median_age"].max(),
5                     500).reshape(-1, 1)
6 gamma1 = 0.1
7 gamma2 = 0.03
8 rbf1 = rbf_kernel(ages, [[35]], gamma=gamma1)
9 rbf2 = rbf_kernel(ages, [[35]], gamma=gamma2)
10
11 fig, ax1 = plt.subplots()
12
13 ax1.set_xlabel("Housing median age")
14 ax1.set_ylabel("Number of districts")
15 ax1.hist(housing["housing_median_age"], bins=50)
16
17 ax2 = ax1.twinx() # create a twin axis that shares the same x-axis
18 color = "blue"
19 ax2.plot(ages, rbf1, color=color, label="gamma = 0.10")
20 ax2.plot(ages, rbf2, color=color, label="gamma = 0.03", linestyle="--")
21 ax2.tick_params(axis='y', labelcolor=color)
22 ax2.set_ylabel("Age similarity", color=color)
23
24 plt.legend(loc="upper left")
25 save_fig("age_similarity_plot")
26 plt.show()
```



```
1 from sklearn.linear_model import LinearRegression
2
3 target_scaler = StandardScaler()
4 scaled_labels = target_scaler.fit_transform(housing_labels.to_frame())
5
6 model = LinearRegression()
7 model.fit(housing[["median_income"]], scaled_labels)
8 some_new_data = housing[["median_income"]].iloc[:5] # pretend this is new data
9
10 scaled_predictions = model.predict(some_new_data)
11 predictions = target_scaler.inverse_transform(scaled_predictions)

1 predictions
[{"array([[131997.15275877],
       [299359.35844434],
       [146023.37185694],}]]
```

```
[138840.33653057],  
[192016.61557639]])
```

```
1 from sklearn.compose import TransformedTargetRegressor  
2  
3 model = TransformedTargetRegressor(LinearRegression(),  
4                                     transformer=StandardScaler())  
5 model.fit(housing[["median_income"]], housing_labels)  
6 predictions = model.predict(some_new_data)  
  
1 predictions  
→ array([131997.15275877, 299359.35844434, 146023.37185694, 138840.33653057,  
       192016.61557639])
```

▼ Custom Transformers

To create simple transformers:

```
1 from sklearn.preprocessing import FunctionTransformer  
2  
3 log_transformer = FunctionTransformer(np.log, inverse_func=np.exp)  
4 log_pop = log_transformer.transform(housing[["population"]])  
  
1 rbf_transformer = FunctionTransformer(rbf_kernel,  
2                                       kw_args=dict(Y=[[35.]], gamma=0.1))  
3 age_simil_35 = rbf_transformer.transform(housing[["housing_median_age"]])  
  
1 age_simil_35  
→ array([[2.81118530e-13],  
       [8.20849986e-02],  
       [6.70320046e-01],  
       ...,  
       [9.55316054e-22],  
       [6.70320046e-01],  
       [3.03539138e-04]])
```

```
1 sf_coords = 37.7749, -122.41
2 sf_transformer = FunctionTransformer(rbf_kernel,
3                                     kw_args=dict(Y=[sf_coords], gamma=0.1))
4 sf_simil = sf_transformer.transform(housing[["latitude", "longitude"]])

1 sf_simil
→ array([[0.999927  ],
       [0.05258419],
       [0.94864161],
       ...,
       [0.00388525],
       [0.05038518],
       [0.99868067]])

1 ratio_transformer = FunctionTransformer(lambda X: X[:, [0]] / X[:, [1]])
2 ratio_transformer.transform(np.array([[1., 2.], [3., 4.]]))

→ array([[0.5 ],
       [0.75]])
```

```
1 from sklearn.base import BaseEstimator, TransformerMixin
2 from sklearn.utils.validation import check_array, check_is_fitted
3
4 class StandardScalerClone(BaseEstimator, TransformerMixin):
5     def __init__(self, with_mean=True): # no *args or **kwargs!
6         self.with_mean = with_mean
7
8     def fit(self, X, y=None): # y is required even though we don't use it
9         X = check_array(X) # checks that X is an array with finite float values
10        self.mean_ = X.mean(axis=0)
11        self.scale_ = X.std(axis=0)
12        self.n_features_in_ = X.shape[1] # every estimator stores this in fit()
13        return self # always return self!
14
15    def transform(self, X):
16        check_is_fitted(self) # looks for learned attributes (with trailing _)
17        X = check_array(X)
18        assert self.n_features_in_ == X.shape[1]
19        if self.with_mean:
20            X = X - self.mean_
21        return X / self.scale_
```

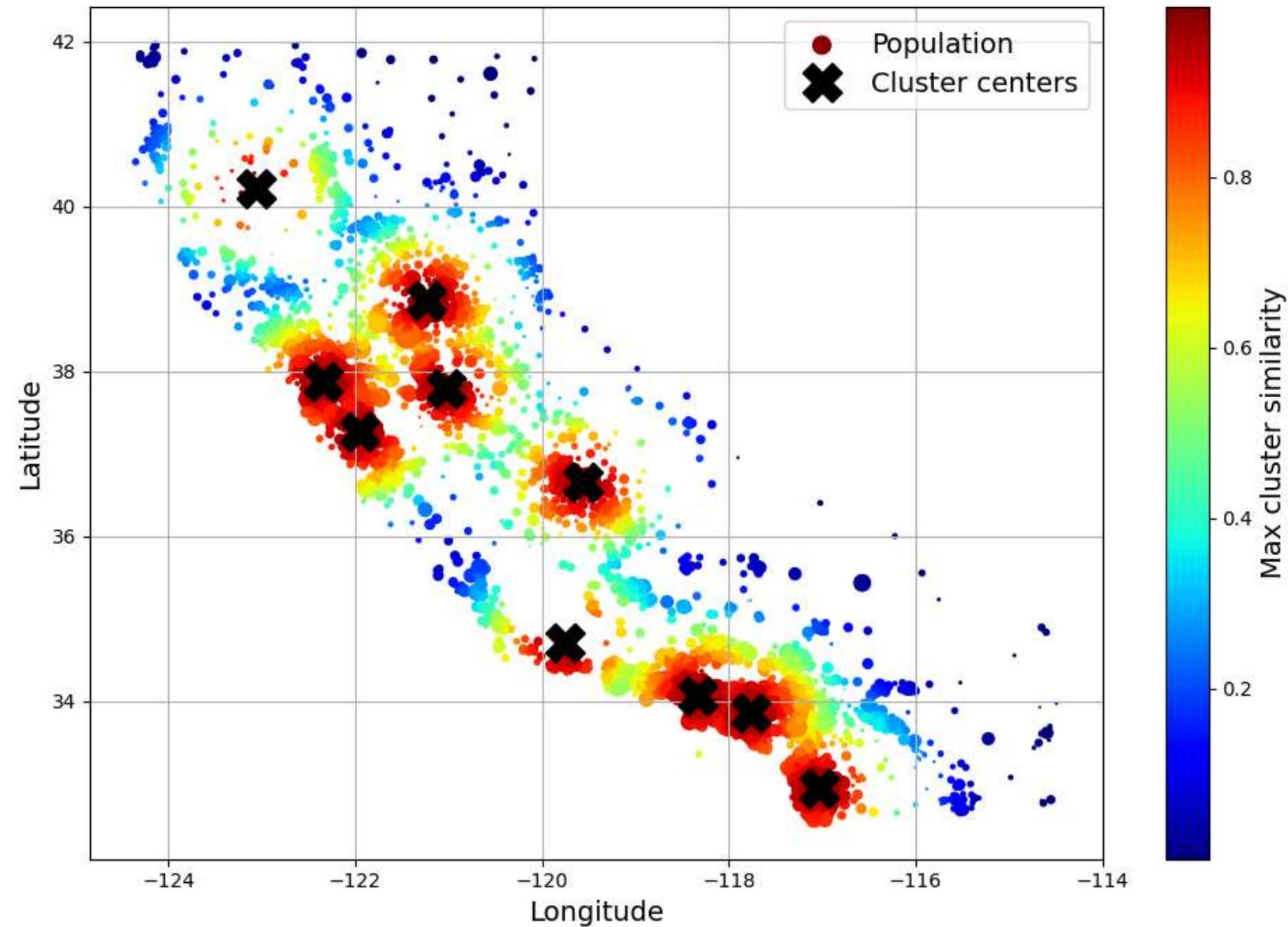
```
1 from sklearn.cluster import KMeans
2
3 class ClusterSimilarity(BaseEstimator, TransformerMixin):
4     def __init__(self, n_clusters=10, gamma=1.0, random_state=None):
5         self.n_clusters = n_clusters
6         self.gamma = gamma
7         self.random_state = random_state
8
9     def fit(self, X, y=None, sample_weight=None):
10        self.kmeans_ = KMeans(self.n_clusters, n_init=10,
11                             random_state=self.random_state)
12        self.kmeans_.fit(X, sample_weight=sample_weight)
13        return self # always return self!
14
15    def transform(self, X):
16        return rbf_kernel(X, self.kmeans_.cluster_centers_, gamma=self.gamma)
17
18    def get_feature_names_out(self, names=None):
19        return [f"Cluster {i} similarity" for i in range(self.n_clusters)]
```

```
1 cluster_simil = ClusterSimilarity(n_clusters=10, gamma=1., random_state=42)
2 similarities = cluster_simil.fit_transform(housing[["latitude", "longitude"]],
3                                         sample_weight=housing_labels)
```

```
1 similarities[:3].round(2)
```

```
→ array([[0.08, 0. , 0.6 , 0. , 0. , 0.99, 0. , 0. , 0. , 0.14],
       [0. , 0.99, 0. , 0.04, 0. , 0. , 0.11, 0. , 0.63, 0. ],
       [0.44, 0. , 0.3 , 0. , 0. , 0.7 , 0. , 0.01, 0. , 0.29]])
```

```
1 # extra code - this cell generates Figure 2-19
2
3 housing_renamed = housing.rename(columns={
4     "latitude": "Latitude", "longitude": "Longitude",
5     "population": "Population",
6     "median_house_value": "Median house value (usd)"})
7 housing_renamed["Max cluster similarity"] = similarities.max(axis=1)
8
9 housing_renamed.plot(kind="scatter", x="Longitude", y="Latitude", grid=True,
10                      s=housing_renamed["Population"] / 100, label="Population",
11                      c="Max cluster similarity",
12                      cmap="jet", colorbar=True,
13                      legend=True, sharex=False, figsize=(10, 7))
14 plt.plot(cluster_simil.kmeans_.cluster_centers_[:, 1],
15           cluster_simil.kmeans_.cluster_centers_[:, 0],
16           linestyle="", color="black", marker="X", markersize=20,
17           label="Cluster centers")
18 plt.legend(loc="upper right")
19 save_fig("district_cluster_plot")
20 plt.show()
```



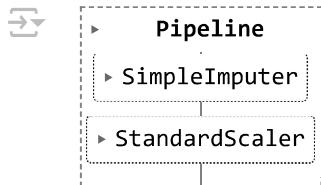
▼ Transformation Pipelines

Now let's build a pipeline to preprocess the numerical attributes:

```
1 from sklearn.pipeline import Pipeline
2
3 num_pipeline = Pipeline([
4     ("impute", SimpleImputer(strategy="median")),
5     ("standardize", StandardScaler()),
6 ])
```

```
1 from sklearn.pipeline import make_pipeline
2
3 num_pipeline = make_pipeline(SimpleImputer(strategy="median"), StandardScaler())
```

```
1 from sklearn import set_config
2
3 set_config(display='diagram')
4
5 num_pipeline
```



```
1 housing_num_prepared = num_pipeline.fit_transform(housing_num)
2 housing_num_prepared[:2].round(2)
```

```
array([[-1.42,  1.01,  1.86,  0.31,  1.37,  0.14,  1.39, -0.94],
       [ 0.6 , -0.7 ,  0.91, -0.31, -0.44, -0.69, -0.37,  1.17]])
```

```
1 def monkey_patch_get_signature_names_out():
2     """Monkey patch some classes which did not handle get_feature_names_out()
3         correctly in Scikit-Learn 1.0.*."""
4     from inspect import Signature, signature, Parameter
5     import pandas as pd
6     from sklearn.impute import SimpleImputer
7     from sklearn.pipeline import make_pipeline, Pipeline
8     from sklearn.preprocessing import FunctionTransformer, StandardScaler
9
10    default_get_feature_names_out = StandardScaler.get_feature_names_out
11
12    if not hasattr(SimpleImputer, "get_feature_names_out"):
13        print("Monkey-patching SimpleImputer.get_feature_names_out()")
14        SimpleImputer.get_feature_names_out = default_get_feature_names_out
15
16    if not hasattr(FunctionTransformer, "get_feature_names_out"):
17        print("Monkey-patching FunctionTransformer.get_feature_names_out()")
18        orig_init = FunctionTransformer.__init__
19        orig_sig = signature(orig_init)
20
21        def __init__(*args, feature_names_out=None, **kwargs):
22            orig_sig.bind(*args, **kwargs)
23            orig_init(*args, **kwargs)
24            args[0].feature_names_out = feature_names_out
25
26            __init__.signature = Signature(
27                list(signature(orig_init).parameters.values()) +
28                [Parameter("feature_names_out", Parameter.KEYWORD_ONLY)])
29
30        def get_feature_names_out(self, names=None):
31            if callable(self.feature_names_out):
32                return self.feature_names_out(self, names)
33            assert self.feature_names_out == "one-to-one"
34            return default_get_feature_names_out(self, names)
35
36        FunctionTransformer.__init__ = __init__
37        FunctionTransformer.get_feature_names_out = get_feature_names_out
38
39 monkey_patch_get_signature_names_out()
```

```
1 df_housing_num_prepared = pd.DataFrame(  
2     housing_num_prepared, columns=num_pipeline.get_feature_names_out(),  
3     index=housing_num.index)
```

```
1 df_housing_num_prepared.head(2) # extra code
```

Next steps: [Generate code with df_housing_num_prepared](#) [View recommended plots](#) [New interactive sheet](#)

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income
13096	-1.423037	1.013606	1.861119	0.311912	1.368167	0.137460	1.394812	-0.936491
14973	0.596394	-0.702103	0.907630	-0.308620	-0.435925	-0.693771	-0.373485	1.171942

```
1 num_pipeline.steps
```

```
[(('simpleimputer', SimpleImputer(strategy='median')),  
 ('standardscaler', StandardScaler()))]
```

```
1 num_pipeline[1]
```

```
StandardScaler()
```

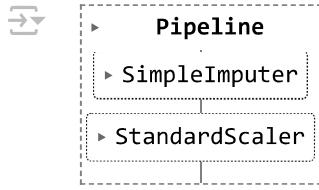
```
1 num_pipeline[:-1]
```

```
Pipeline  
SimpleImputer
```

```
1 num_pipeline.named_steps["simpleimputer"]
```

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

```
1 num_pipeline.set_params(simpleimputer__strategy="median")
```



```
1 from sklearn.compose import ColumnTransformer
2
3 num_attribs = ["longitude", "latitude", "housing_median_age", "total_rooms",
4                 "total_bedrooms", "population", "households", "median_income"]
5 cat_attribs = ["ocean_proximity"]
6
7 cat_pipeline = make_pipeline(
8     SimpleImputer(strategy="most_frequent"),
9     OneHotEncoder(handle_unknown="ignore"))
10
11 preprocessing = ColumnTransformer([
12     ("num", num_pipeline, num_attribs),
13     ("cat", cat_pipeline, cat_attribs),
14 ])
15
16
17 from sklearn.compose import make_column_selector, make_column_transformer
18
19 preprocessing = make_column_transformer(
20     (num_pipeline, make_column_selector(dtype_include=np.number)),
21     (cat_pipeline, make_column_selector(dtype_include=object)),
22 )
23
24
25 housing_prepared = preprocessing.fit_transform(housing)
26
27
28 # extra code - shows that we can get a DataFrame out if we want
29 housing_prepared_fr = pd.DataFrame(
30     housing_prepared,
31     columns=preprocessing.get_feature_names_out(),
32     index=housing.index)
33 housing_prepared_fr.head(2)
```



	pipeline- 1_longitude	pipeline- 1_latitude	pipeline- 1_housing_median_age	pipeline- 1_total_rooms	pipeline- 1_total_bedrooms	pipeline- 1_population	pipeline- 1_households	pipeline- 1_median_inco
--	--------------------------	-------------------------	-----------------------------------	----------------------------	-------------------------------	---------------------------	---------------------------	----------------------------

13096	-1.423037	1.013606	1.861119	0.311912	1.368167	0.137460	1.394812	-0.9364
14973	0.596394	-0.702103	0.907630	-0.308620	-0.435925	-0.693771	-0.373485	1.1719



Next steps: [Generate code with housing_prepared_fr](#)

[View recommended plots](#)

[New interactive sheet](#)

```
1 def column_ratio(X):
2     return X[:, [0]] / X[:, [1]]
3
4 def ratio_name(function_transformer, feature_names_in):
5     return ["ratio"] # feature names out
6
7 def ratio_pipeline():
8     return make_pipeline(
9         SimpleImputer(strategy="median"),
10        FunctionTransformer(column_ratio, feature_names_out=ratio_name),
11        StandardScaler())
12
13 log_pipeline = make_pipeline(
14    SimpleImputer(strategy="median"),
15    FunctionTransformer(np.log, feature_names_out="one-to-one"),
16    StandardScaler())
17 cluster_simil = ClusterSimilarity(n_clusters=10, gamma=1., random_state=42)
18 default_num_pipeline = make_pipeline(SimpleImputer(strategy="median"),
19                                      StandardScaler())
20 preprocessing = ColumnTransformer([
21    ("bedrooms", ratio_pipeline(), ["total_bedrooms", "total_rooms"]),
22    ("rooms_per_house", ratio_pipeline(), ["total_rooms", "households"]),
23    ("people_per_house", ratio_pipeline(), ["population", "households"]),
24    ("log", log_pipeline, ["total_bedrooms", "total_rooms", "population",
25                           "households", "median_income"]),
26    ("geo", cluster_simil, ["latitude", "longitude"]),
27    ("cat", cat_pipeline, make_column_selector(dtype_include=object)),
28 ],
29 remainder=default_num_pipeline) # one column remaining: housing_median_age
```

```
1 housing_prepared = preprocessing.fit_transform(housing)
2 housing_prepared.shape
```

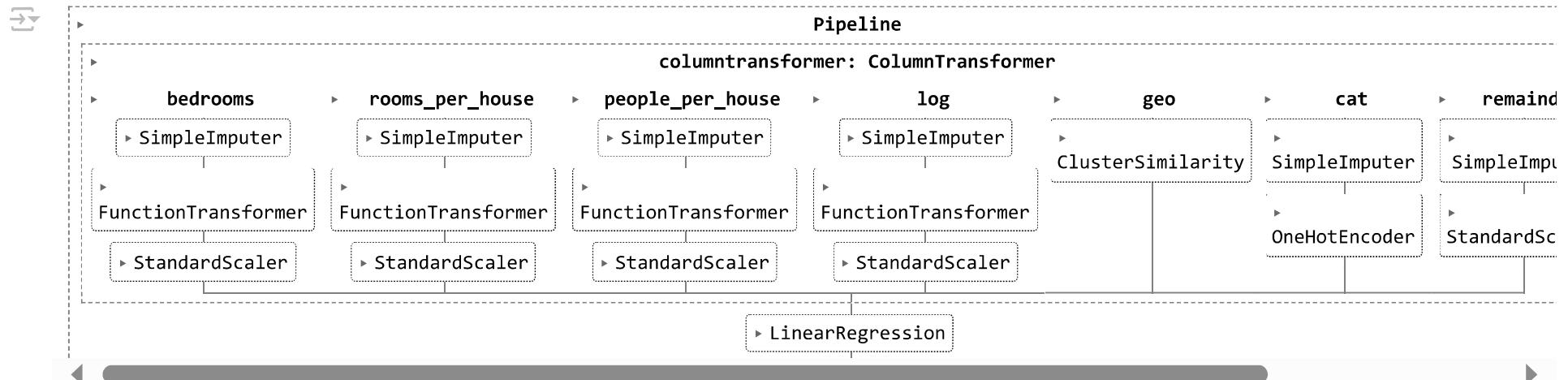
→ (16512, 24)

```
1 preprocessing.get_feature_names_out()
```

→ array(['bedrooms_ratio', 'rooms_per_house_ratio',
'people_per_house_ratio', 'log_total_bedrooms',
'log_total_rooms', 'log_population', 'log_households',
'log_median_income', 'geo_Cluster 0 similarity',
'geo_Cluster 1 similarity', 'geo_Cluster 2 similarity',
'geo_Cluster 3 similarity', 'geo_Cluster 4 similarity',
'geo_Cluster 5 similarity', 'geo_Cluster 6 similarity',
'geo_Cluster 7 similarity', 'geo_Cluster 8 similarity',
'geo_Cluster 9 similarity', 'cat_ocean_proximity_<1H OCEAN',
'cat_ocean_proximity_INLAND', 'cat_ocean_proximity_ISLAND',
'cat_ocean_proximity_NEAR BAY', 'cat_ocean_proximity_NEAR OCEAN',
'remainder_housing_median_age'], dtype=object)

▼ Select and Train a Model

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



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

```
1 housing_predictions = lin_reg.predict(housing)
2 housing_predictions[:5].round(-2) # -2 = rounded to the nearest hundred
→ array([242800., 375900., 127500., 99400., 324600.])
```

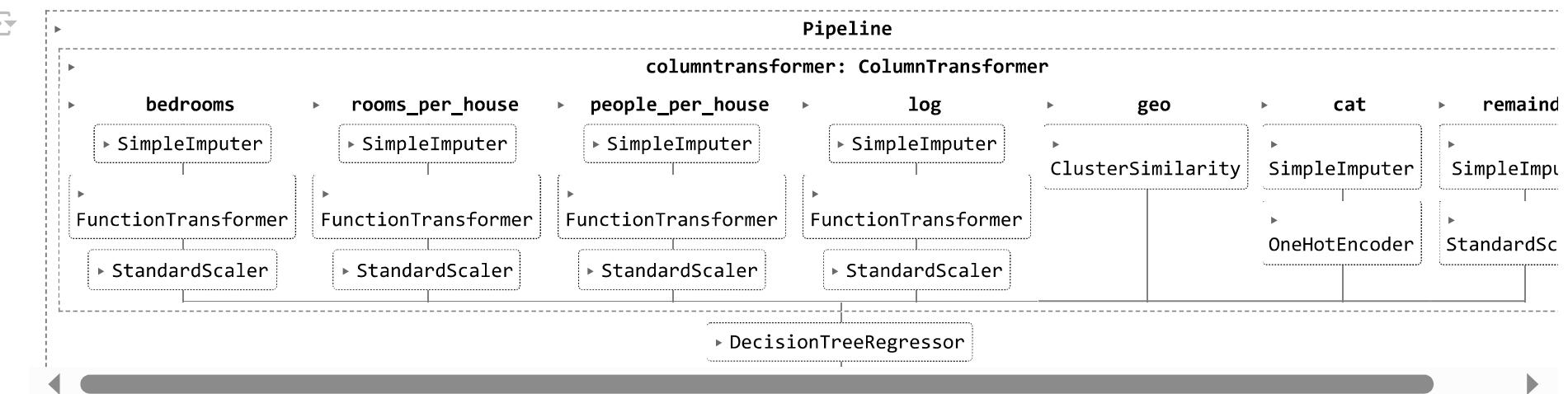
Compare against the actual values:

```
1 housing_labels.iloc[:5].values
→ array([458300., 483800., 101700., 96100., 361800.])

1 # extra code - computes the error ratios discussed in the book
2 error_ratios = housing_predictions[:5].round(-2) / housing_labels.iloc[:5].values - 1
3 print(", ".join(["f'{100 * ratio:.1f}%" for ratio in error_ratios]))
→ -47.0%, -22.3%, 25.4%, 3.4%, -10.3%
```

```
1 from sklearn.metrics import mean_squared_error
2
3 lin_rmse = mean_squared_error(housing_labels, housing_predictions,
4                                squared=False)
5 lin_rmse
→ 68647.95686706658
```

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



```

1 housing_predictions = tree_reg.predict(housing)
2 tree_rmse = mean_squared_error(housing_labels, housing_predictions,
3                                squared=False)
4 tree_rmse

```

0.0

▼ Better Evaluation Using Cross-Validation

```

1 from sklearn.model_selection import cross_val_score
2
3 tree_rmses = -cross_val_score(tree_reg, housing, housing_labels,
4                               scoring="neg_root_mean_squared_error", cv=10)
5
6 pd.Series(tree_rmses).describe()

```



0

```
count    10.000000
mean   67153.318273
std     1963.580924
min    63925.253106
25%    66083.277180
50%    66795.829871
75%    68074.018403
max    70664.635833
```

dtype: float64

```
1 # extra code - computes the error stats for the linear model
2 lin_rmses = -cross_val_score(lin_reg, housing, housing_labels,
3                               scoring="neg_root_mean_squared_error", cv=10)
4 pd.Series(lin_rmses).describe()
```



0

```
count    10.000000
mean   69847.923224
std     4078.407329
min    65659.761079
25%    68088.799156
50%    68697.591463
75%    69800.966364
max    80685.254832
```

dtype: float64

```
1 from sklearn.ensemble import RandomForestRegressor  
2  
3 forest_reg = make_pipeline(preprocessing,  
4                             RandomForestRegressor(random_state=42))  
5 forest_rmses = -cross_val_score(forest_reg, housing, housing_labels,  
6                                 scoring="neg_root_mean_squared_error", cv=10)  
  
1 pd.Series(forest_rmses).describe()
```

```
0  
count    10.000000  
mean   47002.931706  
std     1048.451340  
min    45667.064036  
25%    46494.358345  
50%    47093.173938  
75%    47274.873814  
max    49354.705514  
  
dtype: float64
```

Let's compare this RMSE measured using cross-validation (the "validation error") with the RMSE measured on the training set (the "training error"):

```
1 forest_reg.fit(housing, housing_labels)  
2 housing_predictions = forest_reg.predict(housing)  
3 forest_rmse = mean_squared_error(housing_labels, housing_predictions,  
4                                   squared=False)  
5 forest_rmse
```

```
17547.52124624957
```

The training error is much lower than the validation error, which usually means that the model has overfit the training set. Another possible explanation may be that there's a mismatch between the training data and the validation data, but it's not the case here, since both came from

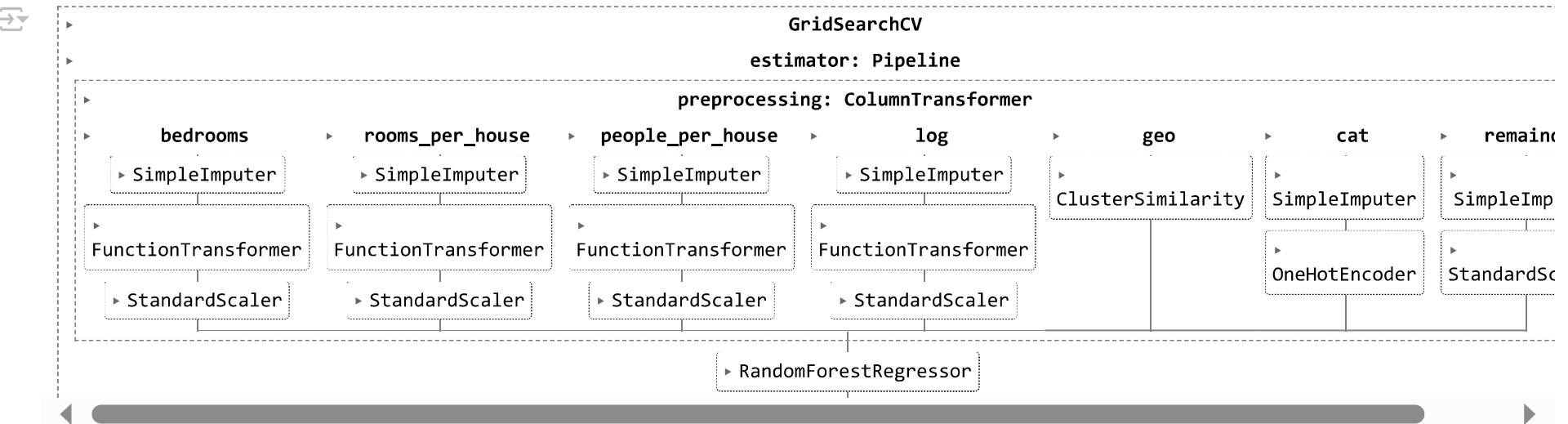
the same dataset that we shuffled and split in two parts.

- ✓ Fine-Tune Your Model

- ✓ Grid Search

Warning: the following cell may take a few minutes to run:

```
1 from sklearn.model_selection import GridSearchCV
2
3 full_pipeline = Pipeline([
4     ("preprocessing", preprocessing),
5     ("random_forest", RandomForestRegressor(random_state=42)),
6 ])
7 param_grid = [
8     {'preprocessing_geo_n_clusters': [5, 8, 10],
9      'random_forest_max_features': [4, 6, 8]},
10    {'preprocessing_geo_n_clusters': [10, 15],
11      'random_forest_max_features': [6, 8, 10]},
12 ]
13 grid_search = GridSearchCV(full_pipeline, param_grid, cv=3,
14                            scoring='neg_root_mean_squared_error')
15 grid_search.fit(housing, housing_labels)
```



You can get the full list of hyperparameters available for tuning by looking at `full_pipeline.get_params().keys()`:

```

1 # extra code - shows part of the output of get_params().keys()
2 print(str(full_pipeline.get_params().keys()[:1000] + "..."))

```

The best hyperparameter combination found:

```

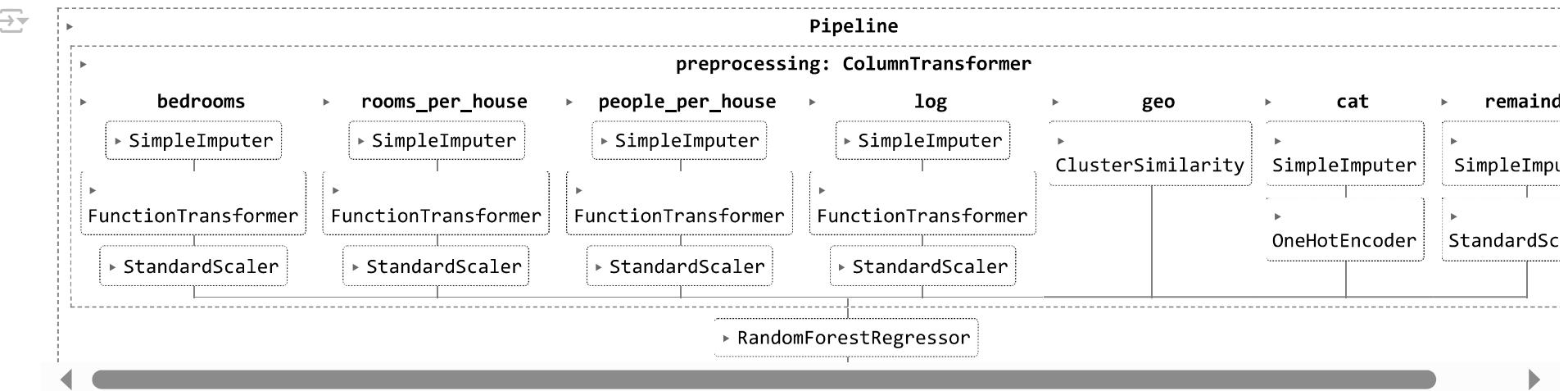
1 grid_search.best_params_

```

```

1 grid_search.best_estimator_

```



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

```

1 cv_res = pd.DataFrame(grid_search.cv_results_)
2 cv_res.sort_values(by="mean_test_score", ascending=False, inplace=True)
3
4 # extra code - these few lines of code just make the DataFrame look nicer
5 cv_res = cv_res[["param_preprocessing_geo_n_clusters",
6                 "param_random_forest_max_features", "split0_test_score",
7                 "split1_test_score", "split2_test_score", "mean_test_score"]]
8 score_cols = ["split0", "split1", "split2", "mean_test_rmse"]
9 cv_res.columns = ["n_clusters", "max_features"] + score_cols
10 cv_res[score_cols] = -cv_res[score_cols].round().astype(np.int64)
11
12 cv_res.head()

```

	n_clusters	max_features	split0	split1	split2	mean_test_rmse	grid
12	15		6	43536	43753	44569	43953
13	15		8	44084	44205	44863	44384
14	15		10	44368	44496	45200	44688
7	10		6	44251	44628	45857	44912
9	10		6	44251	44628	45857	44912

Next steps:

[Generate code with cv_res](#)

[View recommended plots](#)

[New interactive sheet](#)

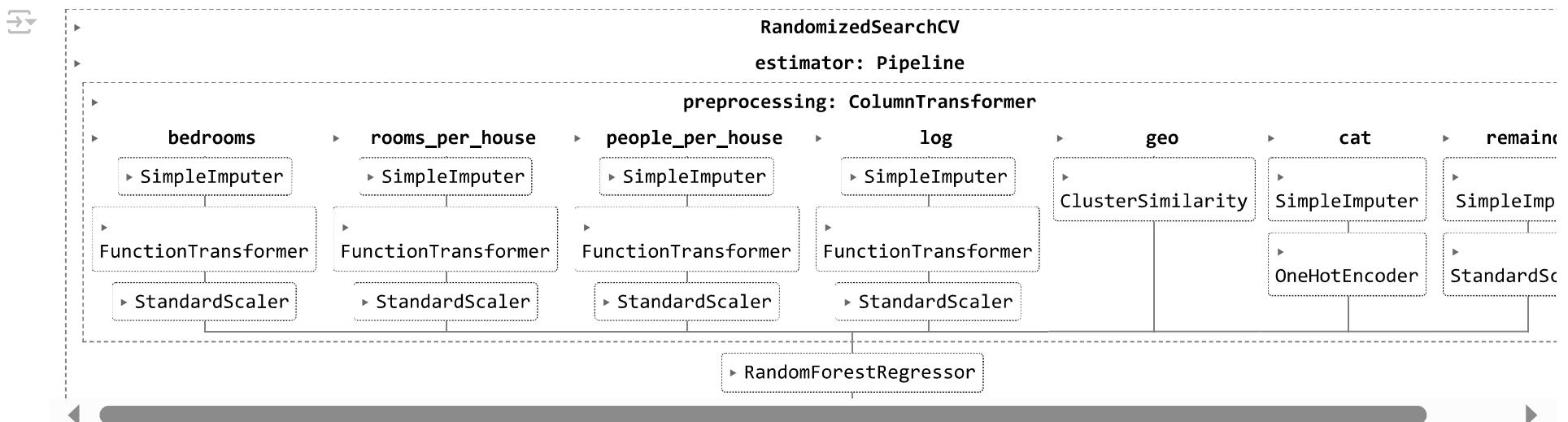
Randomized Search

```
1 from sklearn.experimental import enable_halving_search_cv  
2 from sklearn.model_selection import HalvingRandomSearchCV
```

Try 30 ($n_{\text{iter}} \times cv$) random combinations of hyperparameters:

Warning: the following cell may take a few minutes to run:

```
1 from sklearn.model_selection import RandomizedSearchCV  
2 from scipy.stats import randint  
3  
4 param_distributions = {'preprocessing__geo__n_clusters': randint(low=3, high=50),  
5                         'random_forest__max_features': randint(low=2, high=20)}  
6  
7 rnd_search = RandomizedSearchCV(  
8     full_pipeline, param_distributions=param_distributions, n_iter=10, cv=3,  
9     scoring='neg_root_mean_squared_error', random_state=42)  
10  
11 rnd_search.fit(housing, housing_labels)
```



```

1 # extra code - displays the random search results
2 cv_res = pd.DataFrame(rnd_search.cv_results_)
3 cv_res.sort_values(by="mean_test_score", ascending=False, inplace=True)
4 cv_res = cv_res[["param_preprocessing_geo_n_clusters",
5                 "param_random_forest_max_features", "split0_test_score",
6                 "split1_test_score", "split2_test_score", "mean_test_score"]]
7 cv_res.columns = ["n_clusters", "max_features"] + score_cols
8 cv_res[score_cols] = -cv_res[score_cols].round().astype(np.int64)
9 cv_res.head()

```

	n_clusters	max_features	split0	split1	split2	mean_test_rmse	
1	45	9	41115	42151	42695	41987	
8	32	7	41604	42200	43219	42341	
0	41	16	42106	42743	43443	42764	
5	42	4	41812	42925	43557	42765	
2	23	8	42421	43094	43856	43124	

Next steps: [Generate code with cv_res](#) [View recommended plots](#) [New interactive sheet](#)

Bonus section: how to choose the sampling distribution for a hyperparameter

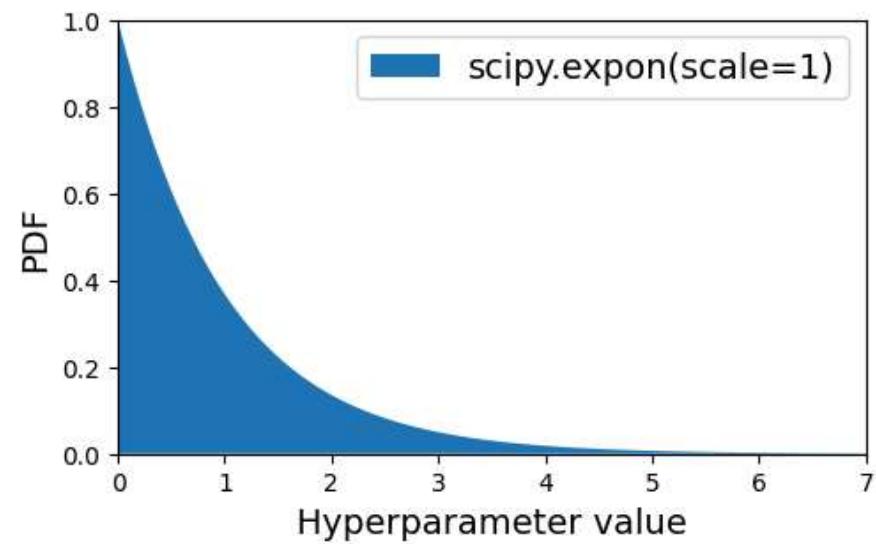
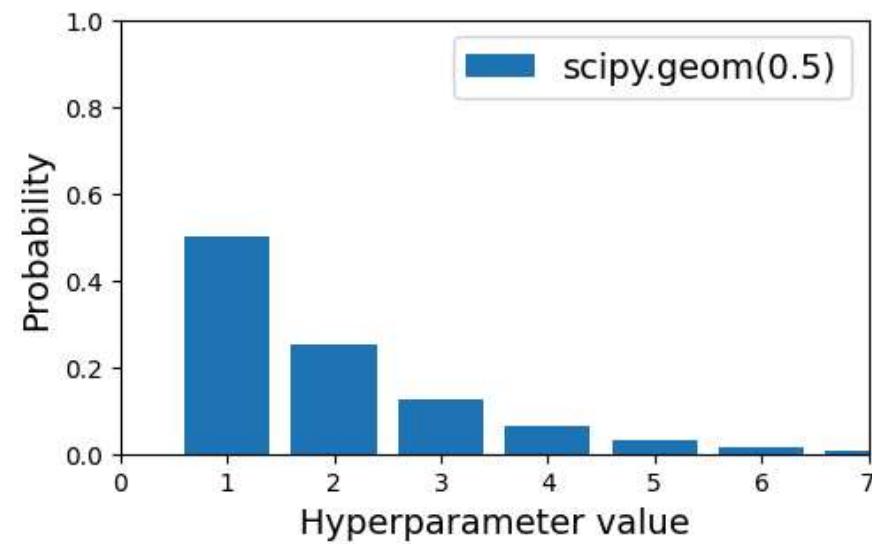
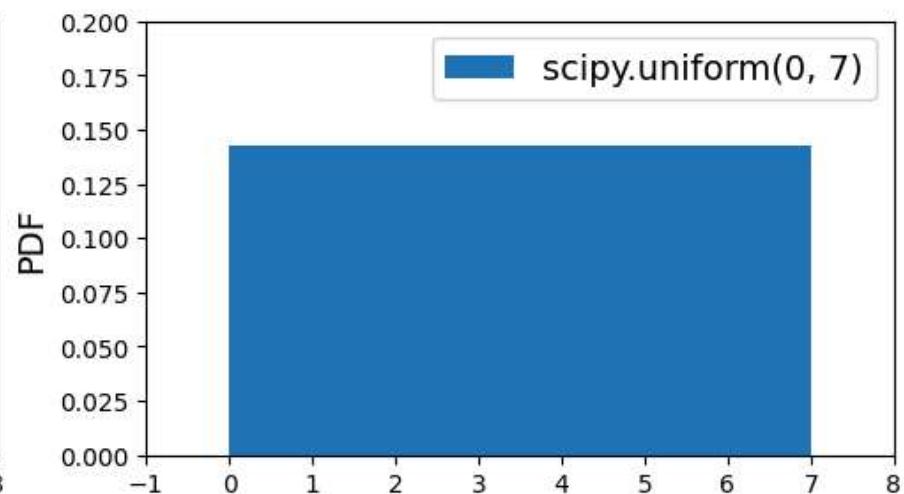
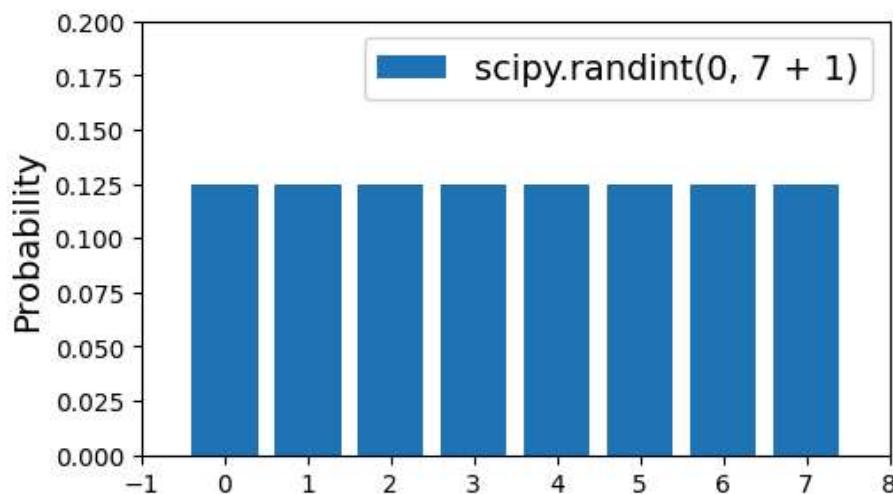
- `scipy.stats.randint(a, b+1)`: for hyperparameters with *discrete* values that range from a to b, and all values in that range seem equally likely.
- `scipy.stats.uniform(a, b)`: this is very similar, but for *continuous* hyperparameters.
- `scipy.stats.geom(1 / scale)`: for discrete values, when you want to sample roughly in a given scale. E.g., with `scale=1000` most samples will be in this ballpark, but ~10% of all samples will be <100 and ~10% will be >2300.
- `scipy.stats.expon(scale)`: this is the continuous equivalent of `geom`. Just set `scale` to the most likely value.
- `scipy.stats.loguniform(a, b)`: when you have almost no idea what the optimal hyperparameter value's scale is. If you set `a=0.01` and `b=100`, then you're just as likely to sample a value between 0.01 and 0.1 as a value between 10 and 100.

Here are plots of the probability mass functions (for discrete variables), and probability density functions (for continuous variables) for `randint()`, `uniform()`, `geom()` and `expon()`:

```
1 # extra code - plots a few distributions you can use in randomized search
2
3 from scipy.stats import randint, uniform, geom, expon
4
5 xs1 = np.arange(0, 7 + 1)
6 randint_distrib = randint(0, 7 + 1).pmf(xs1)
7
8 xs2 = np.linspace(0, 7, 500)
9 uniform_distrib = uniform(0, 7).pdf(xs2)
10
11 xs3 = np.arange(0, 7 + 1)
12 geom_distrib = geom(0.5).pmf(xs3)
13
14 xs4 = np.linspace(0, 7, 500)
15 expon_distrib = expon(scale=1).pdf(xs4)
16
17 plt.figure(figsize=(12, 7))
18
19 plt.subplot(2, 2, 1)
20 plt.bar(xs1, randint_distrib, label="scipy.randint(0, 7 + 1)")
21 plt.ylabel("Probability")
22 plt.legend()
23 plt.axis([-1, 8, 0, 0.2])
24
25 plt.subplot(2, 2, 2)
26 plt.fill_between(xs2, uniform_distrib, label="scipy.uniform(0, 7)")
27 plt.ylabel("PDF")
28 plt.legend()
29 plt.axis([-1, 8, 0, 0.2])
30
31 plt.subplot(2, 2, 3)
32 plt.bar(xs3, geom_distrib, label="scipy.geom(0.5)")
33 plt.xlabel("Hyperparameter value")
34 plt.ylabel("Probability")
35 plt.legend()
36 plt.axis([0, 7, 0, 1])
37
38 plt.subplot(2, 2, 4)
39 plt.fill_between(xs4, expon_distrib, label="scipy.expon(scale=1)")
40 plt.xlabel("Hyperparameter value")
41 plt.ylabel("PDF")
42 plt.legend()
43 plt.axis([0, 7, 0, 1])
```

44

45 plt.show()

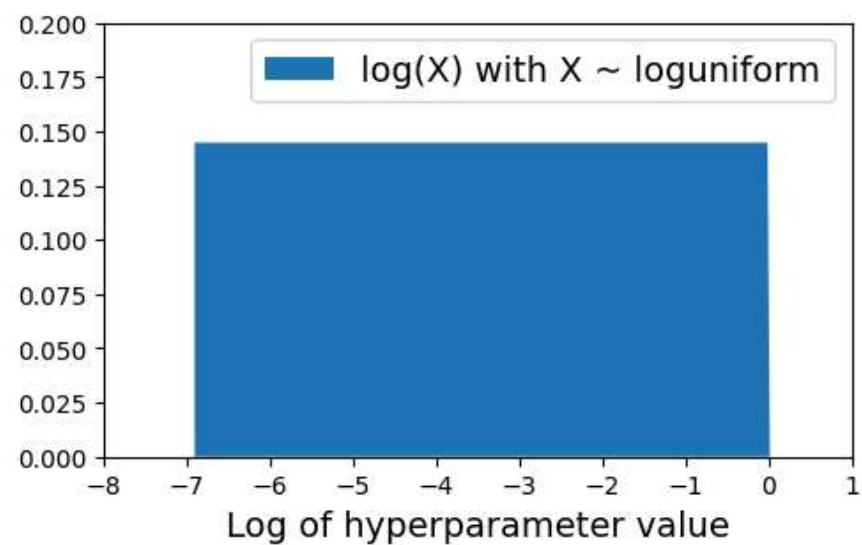
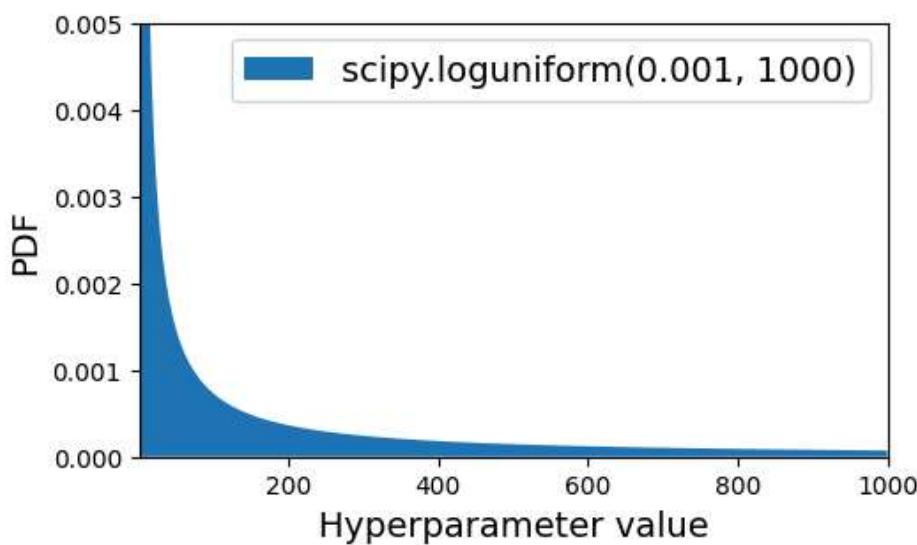
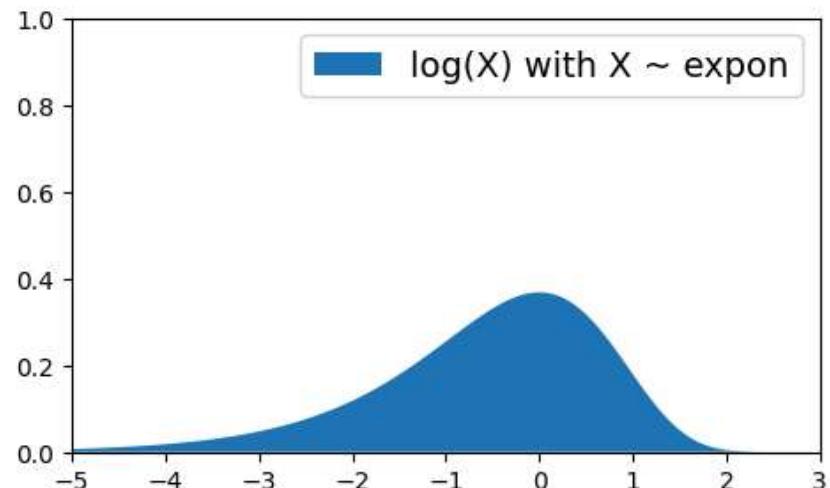
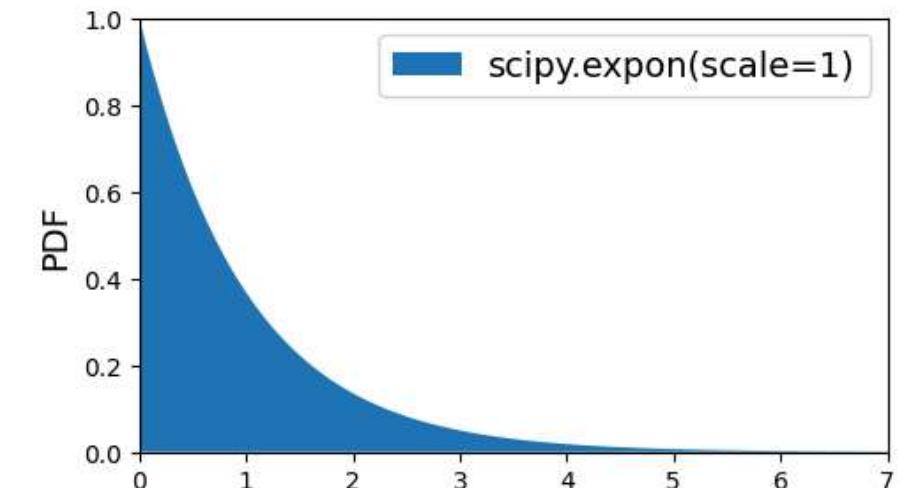


Here are the PDF for `expon()` and `loguniform()` (left column), as well as the PDF of $\log(X)$ (right column). The right column shows the distribution of hyperparameter *scales*. You can see that `expon()` favors hyperparameters with roughly the desired scale, with a longer tail towards the smaller scales. But `loguniform()` does not favor any scale, they are all equally likely:

```
1 # extra code – shows the difference between expon and loguniform
```

```
2
3 from scipy.stats import loguniform
4
5 xs1 = np.linspace(0, 7, 500)
6 expon_distrib = expon(scale=1).pdf(xs1)
7
8 log_xs2 = np.linspace(-5, 3, 500)
9 log_expon_distrib = np.exp(log_xs2 - np.exp(log_xs2))
10
11 xs3 = np.linspace(0.001, 1000, 500)
12 loguniform_distrib = loguniform(0.001, 1000).pdf(xs3)
13
14 log_xs4 = np.linspace(np.log(0.001), np.log(1000), 500)
15 log_loguniform_distrib = uniform(np.log(0.001), np.log(1000)).pdf(log_xs4)
16
17 plt.figure(figsize=(12, 7))
18
19 plt.subplot(2, 2, 1)
20 plt.fill_between(xs1, expon_distrib,
21                  label="scipy.expon(scale=1)")
22 plt.ylabel("PDF")
23 plt.legend()
24 plt.axis([0, 7, 0, 1])
25
26 plt.subplot(2, 2, 2)
27 plt.fill_between(log_xs2, log_expon_distrib,
28                  label="log(X) with X ~ expon")
29 plt.legend()
30 plt.axis([-5, 3, 0, 1])
31
32 plt.subplot(2, 2, 3)
33 plt.fill_between(xs3, loguniform_distrib,
34                  label="scipy.loguniform(0.001, 1000)")
35 plt.xlabel("Hyperparameter value")
36 plt.ylabel("PDF")
37 plt.legend()
38 plt.axis([0.001, 1000, 0, 0.005])
39
40 plt.subplot(2, 2, 4)
41 plt.fill_between(log_xs4, log_loguniform_distrib,
42                  label="log(X) with X ~ loguniform")
43 plt.xlabel("Log of hyperparameter value")
44 plt.legend()
45 plt.axis([-8, 1, 0, 0.2])
--
```

```
46  
47 plt.show()
```



▼ Analyze the Best Models and Their Errors

```
1 final_model = rnd_search.best_estimator_ # includes preprocessing  
2 feature_importances = final_model["random_forest"].feature_importances_  
3 feature_importances.round(2)
```

```
→ array([0.06, 0.06, 0.05, 0.01, 0.01, 0.01, 0.01, 0.19, 0.01, 0.01, 0.02,
       0.04, 0.01, 0. , 0.02, 0.01, 0.01, 0.01, 0.01, 0.01, 0. , 0. ,
       0.01, 0. , 0.01, 0.02, 0.02, 0.01, 0.01, 0.01, 0.03, 0.01, 0.01,
       0.01, 0.01, 0.01, 0. , 0.01, 0.01, 0.02, 0.01, 0.01, 0.01, 0.01,
       0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.08,
       0. , 0. , 0. , 0.01])
```

```
1 sorted(zip(feature_importances,
2             final_model["preprocessing"].get_feature_names_out()),
3             reverse=True)
```

```
→ (0.07709175866873944, 'cat_ocean_proximity_INLAND'),
(0.06455488601956336, 'bedrooms_ratio'),
(0.056936146643377976, 'rooms_per_house_ratio'),
(0.0490294770805355, 'people_per_house_ratio'),
(0.03807069074492323, 'geo_Cluster 3 similarity'),
(0.025643913400094476, 'geo_Cluster 22 similarity'),
(0.02179127543243723, 'geo_Cluster 17 similarity'),
(0.021575251507503695, 'geo_Cluster 6 similarity'),
(0.017868654556924362, 'geo_Cluster 2 similarity'),
(0.017431400050755975, 'geo_Cluster 32 similarity'),
(0.015981159400591683, 'geo_Cluster 18 similarity'),
(0.01488846425739688, 'geo_Cluster 40 similarity'),
(0.014488389218107143, 'geo_Cluster 43 similarity'),
(0.014252940099964142, 'geo_Cluster 7 similarity'),
(0.014038173319370725, 'geo_Cluster 21 similarity'),
(0.013846025114732157, 'geo_Cluster 38 similarity'),
(0.01362570996472274, 'geo_Cluster 34 similarity'),
(0.013547297167034428, 'geo_Cluster 41 similarity'),
(0.012900089026066918, 'geo_Cluster 24 similarity'),
(0.012620908145579916, 'geo_Cluster 10 similarity'),
(0.011621275372313349, 'remainder_housing_median_age'),
(0.01159076532124518, 'geo_Cluster 42 similarity'),
(0.011475471902949968, 'geo_Cluster 30 similarity'),
(0.011145862167873194, 'geo_Cluster 31 similarity'),
(0.011019562821385253, 'geo_Cluster 16 similarity'),
(0.010872453237288457, 'geo_Cluster 1 similarity'),
(0.010683824851366689, 'geo_Cluster 25 similarity'),
(0.01060046411739516, 'geo_Cluster 26 similarity'),
(0.010233275929521545, 'geo_Cluster 20 similarity'),
(0.01000904987677275, 'geo_Cluster 35 similarity'),
(0.009948648710662759, 'geo_Cluster 14 similarity'),
(0.009237024861905358, 'geo_Cluster 39 similarity'),
(0.009095509036677479, 'geo_Cluster 37 similarity'),
(0.009039692793107918, 'geo_Cluster 0 similarity'),
```

```
(0.0088028519810115, 'geo_Cluster 9 similarity'),
(0.00865735444375151, 'geo_Cluster 36 similarity'),
(0.008468774649648216, 'geo_Cluster 28 similarity'),
(0.008132158692323862, 'geo_Cluster 44 similarity'),
(0.00806872019824202, 'geo_Cluster 4 similarity'),
(0.00790076467114424, 'geo_Cluster 11 similarity'),
(0.007654000251882609, 'log_total_rooms'),
(0.007117524274343567, 'log_population'),
(0.006828512079872657, 'log_total_bedrooms'),
(0.006436845508698041, 'log_households'),
(0.00606171524521791, 'geo_Cluster 23 similarity'),
(0.005534308195396936, 'geo_Cluster 19 similarity'),
(0.00549408919549185, 'geo_Cluster 27 similarity'),
(0.0052089915629445595, 'geo_Cluster 33 similarity'),
(0.004814425966149223, 'geo_Cluster 15 similarity'),
(0.004210239717026306, 'geo_Cluster 12 similarity'),
(0.004038656871498407, 'geo_Cluster 13 similarity'),
(0.00373031587904682, 'geo_Cluster 29 similarity'),
(0.003249851526492918, 'cat_ocean_proximity_<1H OCEAN'),
(0.001973755321832911, 'geo_Cluster 5 similarity'),
(0.0018976773551793616, 'cat_ocean_proximity_NEAR OCEAN'),
(0.0002358514810131377, 'cat_ocean_proximity_NEAR BAY'),
(6.124671466559031e-05, 'cat_ocean_proximity ISLAND')]
```

✓ Evaluate Your System on the Test Set

```
1 X_test = strat_test_set.drop("median_house_value", axis=1)
2 y_test = strat_test_set["median_house_value"].copy()
3
4 final_predictions = final_model.predict(X_test)
5
6 final_rmse = mean_squared_error(y_test, final_predictions, squared=False)
7 print(final_rmse)
```

→ 41549.20158097943

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

```
1 from scipy import stats
2
3 confidence = 0.95
4 squared_errors = (final_predictions - y_test) ** 2
5 np.sqrt(stats.t.interval(confidence, len(squared_errors) - 1,
6                         loc=squared_errors.mean(),
7                         scale=stats.sem(squared_errors)))
→ array([39395.35475927, 43596.76969025])
```

We could compute the interval manually like this:

```
1 # extra code – shows how to compute a confidence interval for the RMSE
2 m = len(squared_errors)
3 mean = squared_errors.mean()
4 tscore = stats.t.ppf((1 + confidence) / 2, df=m - 1)
5 tmargin = tscore * squared_errors.std(ddof=1) / np.sqrt(m)
6 np.sqrt(mean - tmargin), np.sqrt(mean + tmargin)
→ (39395.35475926931, 43596.76969025394)
```

Alternatively, we could use a z-score rather than a t-score. Since the test set is not too small, it won't make a big difference:

```
1 # extra code – computes a confidence interval again using a z-score
2 zscore = stats.norm.ppf((1 + confidence) / 2)
3 zmargin = zscore * squared_errors.std(ddof=1) / np.sqrt(m)
4 np.sqrt(mean - zmargin), np.sqrt(mean + zmargin)
→ (39396.00369767951, 43596.18328117898)
```

✓ Model persistence using joblib

Save the final model:

```
1 import joblib
2
3 joblib.dump(final_model, "my_california_housing_model.pkl")
4
5
```

Double-click (or enter) to edit

Now you can deploy this model to production. For example, the following code could be a script that would run in production:

```
1 import joblib
2
3 # extra code - excluded for conciseness
4 from sklearn.cluster import KMeans
5 from sklearn.base import BaseEstimator, TransformerMixin
6 from sklearn.metrics.pairwise import rbf_kernel
7
8 def column_ratio(X):
9     return X[:, [0]] / X[:, [1]]
10
11 #class ClusterSimilarity(BaseEstimator, TransformerMixin):
12 #    [...]
13
14 final_model_reloaded = joblib.load("my_california_housing_model.pkl")
15
16 new_data = housing.iloc[:5] # pretend these are new districts
17 predictions = final_model_reloaded.predict(new_data)

1 predictions
→ array([439808.14, 455211.06, 109492. , 98208. , 340021.04])
```

You could use pickle instead, but joblib is more efficient.

✓ Exercise solutions

▼ 1.

Exercise: 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). Note that SVMs don't scale well to large datasets, so you should probably train your model on just the first 5,000 instances of the training set and use only 3-fold cross-validation, or else it will take hours. Don't worry about what the hyperparameters mean for now (see the SVM notebook if you're interested). How does the best SVR predictor perform?

```
1 from sklearn.model_selection import GridSearchCV
2 from sklearn.svm import SVR
3
4 param_grid = [
5     {'svr_kernel': ['linear'], 'svr_C': [10., 30., 100., 300., 1000.,
6                                         3000., 10000., 30000.0]},
7     {'svr_kernel': ['rbf'], 'svr_C': [1.0, 3.0, 10., 30., 100., 300.,
8                                     1000.0],
9      'svr_gamma': [0.01, 0.03, 0.1, 0.3, 1.0, 3.0]},
10    ]
11
12 svr_pipeline = Pipeline([('preprocessing', preprocessing), ("svr", SVR())])
13 grid_search = GridSearchCV(svr_pipeline, param_grid, cv=3,
14                            scoring='neg_root_mean_squared_error')
15 grid_search.fit(housing.iloc[:5000], housing_labels.iloc[:5000])
```

The best model achieves the following score (evaluated using 3-fold cross validation):

```
1 svr_grid_search_rmse = -grid_search.best_score_
2 svr_grid_search_rmse
→ 69814.13889867254
```

That's much worse than the `RandomForestRegressor` (but to be fair, we trained the model on much less data). Let's check the best hyperparameters found:

```
1 grid_search.best_params_
```

```
['svr_C': 10000.0, 'svr_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.

Exercise: Try replacing the `GridSearchCV` with a `RandomizedSearchCV`.

Warning: the following cell will take several minutes to run. You can specify `verbose=2` when creating the `RandomizedSearchCV` if you want to see the training details.

```
1 from sklearn.model_selection import RandomizedSearchCV
2 from scipy.stats import expon, loguniform
3
4 # see https://docs.scipy.org/doc/scipy/reference/stats.html
5 # for `expon()` and `loguniform()` documentation and more probability distribution functions.
6
7 # Note: gamma is ignored when kernel is "linear"
8 param_distribs = {
9     'svr_kernel': ['linear', 'rbf'],
10    'svr_C': loguniform(20, 200_000),
11    'svr_gamma': expon(scale=1.0),
12}
13
14 rnd_search = RandomizedSearchCV(svr_pipeline,
15                                  param_distributions=param_distribs,
16                                  n_iter=50, cv=3,
17                                  scoring='neg_root_mean_squared_error',
18                                  random_state=42)
19 rnd_search.fit(housing.iloc[:5000], housing_labels.iloc[:5000])
```

```
('standardscaler',
 StandardScaler())]),  
transformers=[('bedrooms',
 Pipeline(steps=[('simpleimputer',
 SimpleImputer(strategy='median')),  
('functiontransformer',
 FunctionTransformer(feature_names_...  
  
<sklearn.compose._column_transformer.make_column_selector object at 0x1a57e3a00>]])),  
('svr', SVR())),  
n_iter=50,  
param_distributions={'svr__C': <scipy.stats._distn_infrastructure.rv_continuous_frozen object at 0x1a5d4c3a0>,  
'svr__gamma': <scipy.stats._distn_infrastructure.rv_continuous_frozen object at 0x1a5d9ca00>,  
'svr__kernel': ['linear', 'rbf']},  
random_state=42, scoring='neg_root_mean_squared_error')
```

The best model achieves the following score (evaluated using 3-fold cross validation):

```
1 svr_rnd_search_rmse = -rnd_search.best_score_  
2 svr_rnd_search_rmse
```

```
→ 55853.88100300133
```

Now that's really much better, but still far from the `RandomForestRegressor`'s performance. Let's check the best hyperparameters found:

```
1 rnd_search.best_params_  
  
→ {'svr__C': 157055.10989448498,  
'svr__gamma': 0.26497040005002437,  
'svr__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.

Note that we used the `expon()` distribution for `gamma`, with a scale of 1, so `RandomSearch` mostly searched for values roughly of that scale: about 80% of the samples were between 0.1 and 2.3 (roughly 10% were smaller and 10% were larger):

```
1 np.random.seed(42)
2
3 s = expon(scale=1).rvs(100_000) # get 100,000 samples
4 ((s > 0.105) & (s < 2.29)).sum() / 100_000
→ 0.80066
```

We used the `loguniform()` distribution for `c`, meaning we did not have a clue what the optimal scale of `c` was before running the random search. It explored the range from 20 to 200 just as much as the range from 2,000 to 20,000 or from 20,000 to 200,000.

▼ 3.

Exercise: *Try adding a `SelectFromModel` transformer in the preparation pipeline to select only the most important attributes.*

Let's create a new pipeline that runs the previously defined preparation pipeline, and adds a `SelectFromModel` transformer based on a `RandomForestRegressor` before the final regressor:

```
1 from sklearn.feature_selection import SelectFromModel
2
3 selector_pipeline = Pipeline([
4     ('preprocessing', preprocessing),
5     ('selector', SelectFromModel(RandomForestRegressor(random_state=42),
6                                   threshold=0.005)), # min feature importance
7     ('svr', SVR(C=rnd_search.best_params_["svr_C"],
8                  gamma=rnd_search.best_params_["svr_gamma"],
9                  kernel=rnd_search.best_params_["svr_kernel"])),
10    ])
11
12 selector_rmses = -cross_val_score(selector_pipeline,
13                                   housing.iloc[:5000],
14                                   housing_labels.iloc[:5000],
15                                   scoring="neg_root_mean_squared_error",
16                                   cv=3)
17 pd.Series(selector_rmses).describe()
→ count      3.000000
→ mean      56211.362086
```

```
std      1922.002802
min     54150.008629
25%    55339.929909
50%    56529.851189
75%    57242.038815
max     57954.226441
dtype: float64
```

Oh well, feature selection does not seem to help. But maybe that's just because the threshold we used was not optimal. Perhaps try tuning it using random search or grid search?

▼ 4.

Exercise: Try creating a custom transformer that trains a k-Nearest Neighbors regressor (`skLearn.neighbors.KNeighborsRegressor`) in its `fit()` method, and outputs the model's predictions in its `transform()` method. Then add this feature to the preprocessing pipeline, using latitude and longitude as the inputs to this transformer. This will add a feature in the model that corresponds to the housing median price of the nearest districts.

Rather than restrict ourselves to k-Nearest Neighbors regressors, let's create a transformer that accepts any regressor. For this, we can extend the `MetaEstimatorMixin` and have a required `estimator` argument in the constructor. The `fit()` method must work on a clone of this estimator, and it must also save `feature_names_in_`. The `MetaEstimatorMixin` will ensure that `estimator` is listed as a required parameters, and it will update `get_params()` and `set_params()` to make the estimator's hyperparameters available for tuning. Lastly, we create a `get_feature_names_out()` method: the output column name is the ...

```
1 from sklearn.neighbors import KNeighborsRegressor
2 from sklearn.base import MetaEstimatorMixin, clone
3
4 class FeatureFromRegressor(MetaEstimatorMixin, BaseEstimator, TransformerMixin):
5     def __init__(self, estimator):
6         self.estimator = estimator
7
8     def fit(self, X, y=None):
9         estimator_ = clone(self.estimator)
10        estimator_.fit(X, y)
11        self.estimator_ = estimator_
12
13    def transform(self, X):
14        return self.estimator_.predict(X)
```