# CALIFORNIA HOUSING PRICE PREDICTION

Vishal Reddy

# Task in Hand:

In this chapter we chose the California Housing Prices dataset. This dataset was based on data from the 1990 California census. This data has metrics such as the population, median income, median housing price.

# Why Machine Learning?

District housing prices are currently estimated manually by experts: a team gathers up-to-date information about a district, and when they cannot get the median housing price, they estimate it using complex rules.
This is costly and time-consuming, and their estimates are not great; in cases where they manage to find out the actual median housing price, they often realize that their estimates were off by more than 20%. This is why firms  think that it would
be useful to train a model to predict a district's median housing price given other data about that district. The census data is a great dataset to exploit for this objective, since it includes the median housing prices of thousands of districts, as well as other data.

This is clearly a typical supervised learning task since we are given labeled training examples (each instance comes with the expected output, i.e., the district's median housing price). Moreover, it is also a typical regression task, since we are supposed to predict a value. More specifically, this is a multiple regression problem since the system uses multiple features to make a prediction (it uses the district's population, the median income, etc.). It is also a univariate regression problem since we are only trying to predict a single value for each district. If we were trying to predict multiple values per district, it would be a multivariate regression problem.

# Selecting a Performance Metric:

A typical performance measure for regression problems is the Root Mean Square Error (RMSE). It gives an idea of how
much error the system typically makes in its predictions, with a higher weight for large errors.

$$\text{RMSE}(\mathbf{X}, h) = \sqrt{\frac{1}{m}\sum_{i=1}^{m}\left(h\left(\mathbf{x}^{(i)}\right) - y^{(i)}\right)^2}$$

- $m$ is the number of instances in the dataset you are measuring the RMSE on.
  - For example, if you are evaluating the RMSE on a validation set of 2,000 districts, then $m = 2,000$.
- $x^{(i)}$ is a vector of all the feature values (excluding the label) of the $i^{th}$ instance in the dataset, and $y^{(i)}$ is its label (the desired output value for that instance).
  - For example, if the first district in the dataset is located at longitude $-118.29°$, latitude $33.91°$, and it has 1,416 inhabitants with a median income of $38,372, and the median house value is $156,400 (ignoring the other features for now), then:

Even though the RMSE is generally the preferred performance measure for regression tasks, in some contexts we may prefer to use another function. For example, suppose that there are many outlier districts. In that case, you may use the Mean Absolute Error.

$$\text{MAE}(X, h) = \frac{1}{m} \sum_{i=1}^{m} \left| h(x^{(i)}) - y^{(i)} \right|$$

# Obtaining Data and creating Dataframe using Pandas

Each row represents one district. There are 10 attributes (you can see the first 6 in the screenshot): longitude, latitude, housing_median_age, total_rooms, total_bed rooms, population, households, median_income, median_house_value, and Ocean_proximity.

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

Out[5]:

| | longitude | latitude | housing_median_age | total_rooms | total_bedrooms | population |
|---|---|---|---|---|---|---|
| 0 | -122.23 | 37.88 | 41.0 | 880.0 | 129.0 | 322.0 |
| 1 | -122.22 | 37.86 | 21.0 | 7099.0 | 1106.0 | 2401.0 |
| 2 | -122.24 | 37.85 | 52.0 | 1467.0 | 190.0 | 496.0 |
| 3 | -122.25 | 37.85 | 52.0 | 1274.0 | 235.0 | 558.0 |
| 4 | -122.25 | 37.85 | 52.0 | 1627.0 | 280.0 | 565.0 |

The info() method is useful to get a quick description of the data, in particular the total number of rows, and each attribute's type and number of non-null values

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

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

There are 20,640 instances in the dataset. The total_bed
rooms attribute has only 20,433 non-null values, meaning that 207 districts are missing
this feature.

All attributes are numerical, except the ocean_proximity field. Its type is object, so it
could hold any kind of Python object. The values in the ocean_proximity column were
repetitive, which means that it is a categorical attribute. We can find out what categories
exist and how many districts belong to each category by using the
value_counts() method:

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

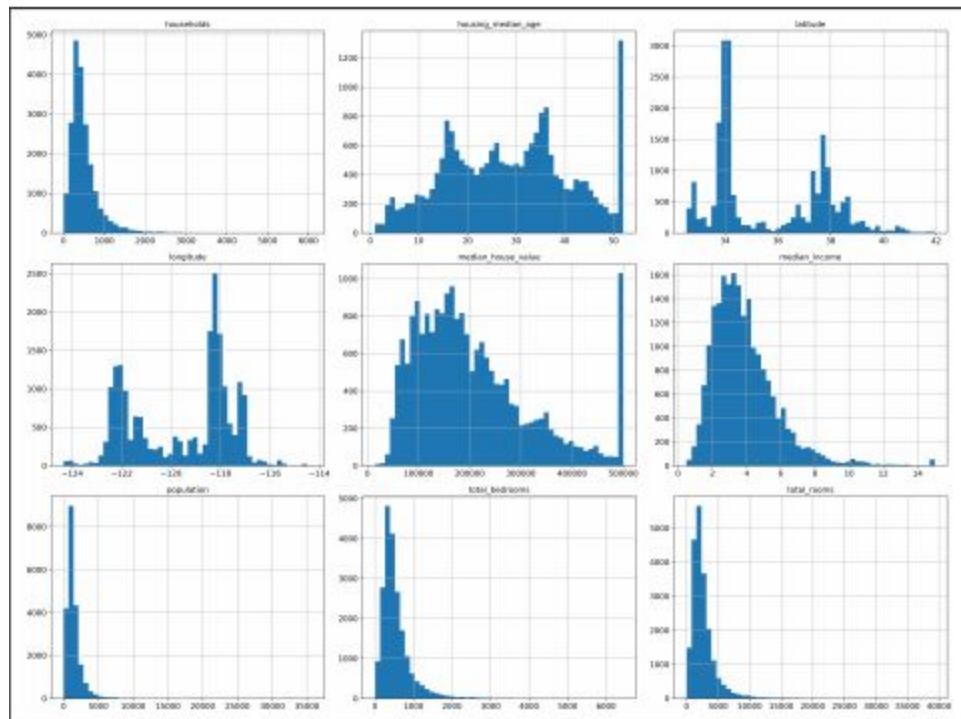The describe() method shows a summary of the numerical attributes:

```
In [8]: housing.describe()
```
Out[8]:

| | longitude | latitude | housing_median_age | total_rooms | total_bedr |
|---|---|---|---|---|---|
| count | 20640.000000 | 20640.000000 | 20640.000000 | 20640.000000 | 20433.0000 |
| mean | -119.569704 | 35.631861 | 28.639486 | 2635.763081 | 537.870553 |
| std | 2.003532 | 2.135952 | 12.585558 | 2181.615252 | 421.385070 |
| min | -124.350000 | 32.540000 | 1.000000 | 2.000000 | 1.000000 |
| 25% | -121.800000 | 33.930000 | 18.000000 | 1447.750000 | 296.000000 |
| 50% | -118.490000 | 34.260000 | 29.000000 | 2127.000000 | 435.000000 |
| 75% | -118.010000 | 37.710000 | 37.000000 | 3148.000000 | 647.000000 |
| max | -114.310000 | 41.950000 | 52.000000 | 39320.000000 | 6445.0000 |

Another quick way to get a feel of the type of data you are dealing with is to plot a histogram for each numerical attribute. The hist() method on the whole dataset will plot a histogram for each numerical attribute. For example, we can see that slightly over 800 districts have a median_house_value equal to about $100,000.

```python
import matplotlib.pyplot as plt
housing.hist(bins=50, figsize=(20,15))
plt.show()
```

A few things to notice in the histogram:
1. First, the median income attribute is not expressed in US dollars
(USD). The data has been scaled and capped at 15 (actually 15.0001) for higher
median incomes, and at 0.5 (actually 0.4999) for lower median incomes. The numbers
represent roughly tens of thousands of dollars (e.g., 3 actually means about
$30,000).
2. The housing median age and the median house value were also capped. The latter
may be a serious problem since it is our target attribute ( labels). The
Machine Learning algorithms may learn that prices never go beyond that limit.
If we need precise predictions even beyond $500,000, then we have mainly two options:
a. Collect proper labels for the districts whose labels were capped.
b. Remove those districts from the training set (and also from the test set, since
the system should not be evaluated poorly if it predicts values beyond
$500,000).
3. These attributes have very different scales.
4.  Many histograms are tail heavy: they extend much farther to the right of
the median than to the left. This may make it a bit harder for some Machine
Learning algorithms to detect patterns. We tried transforming these attributes
 to have more bell-shaped distributions.

# Creating a Training and Test Set

Scikit-Learn provides a few functions to split datasets into multiple subsets in various
ways. The simplest function is train_test_split, which does pretty much the same thing
as the function split_train_test defined earlier, with a couple of additional features. First
there is a random_state parameter that allows you to set the random generator seed as
explained previously, and second you can pass it multiple datasets with an identical
number of rows, and it will split them on the same indices (this is very useful, for
example, if you have a separate DataFrame for labels):

```python
from sklearn.model_selection import train_test_split

train_set, test_set = train_test_split(housing, test_size=0.2, random_state=42)
```
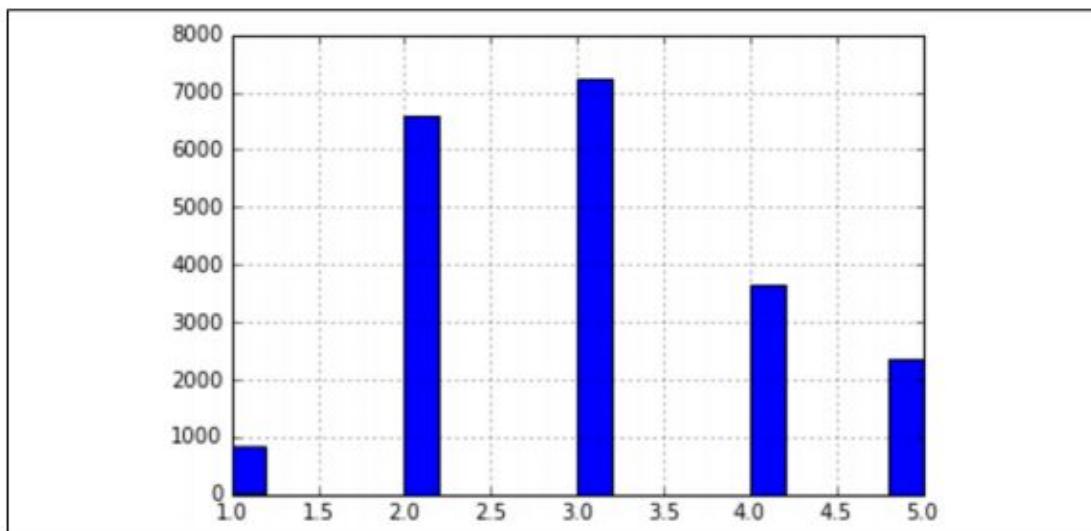
The median income is a very important attribute to predict median housing prices. We
may want to ensure that the test set is representative of the various categories of
incomes in the whole dataset.

Since the median income is a continuous numerical attribute, we first need to create an income category attribute. Looking at the median income histogram more closely: most median income values are clustered around 1.5 to 6 (i.e.,$15,000–$60,000), but some median incomes go far beyond 6. It is important to have a sufficient number of instances in your dataset for each stratum, or else the estimate of the stratum's importance may be biased. This means that we should not have too many strata, and each stratum should be large enough. The following code uses the pd.cut() function to create an income category attribute with 5 categories (labeled from 1 to 5): category 1 ranges from 0 to 1.5 (i.e., less than $15,000), category 2 from 1.5 to 3, and so on:

```
housing["income_cat"] = pd.cut(housing["median_income"],
                               bins=[0., 1.5, 3.0, 4.5, 6., np.inf],
                               labels=[1, 2, 3, 4, 5])
```

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



Stratified sampling is done based on the income category. For this Scikit-Learn's StratifiedShuffleSplit class is used:

```
from sklearn.model_selection import StratifiedShuffleSplit

split = StratifiedShuffleSplit(n_splits=1, test_size=0.2, random_state=42)
for train_index, test_index in split.split(housing, housing["income_cat"]):
    strat_train_set = housing.loc[train_index]
    strat_test_set = housing.loc[test_index]
```

The income category proportions in the test set:

With similar code we can measure the income category proportions in the full data-set. Figure below compares the income category proportions in the overall dataset, in the test set generated with stratified sampling, and in a test set generated using purely random sampling. The test set generated using stratified sampling has income category proportions almost identical to those in the full dataset, whereas the test set generated using purely random sampling is quite skewed.

|       | Overall  | Random   | Stratified | Rand. %error | Strat. %error |
|-------|----------|----------|------------|--------------|---------------|
| 1.0   | 0.039826 | 0.040213 | 0.039738   | 0.973236     | -0.219137     |
| 2.0   | 0.318847 | 0.324370 | 0.318876   | 1.732260     | 0.009032      |
| 3.0   | 0.350581 | 0.358527 | 0.350618   | 2.266446     | 0.010408      |
| 4.0   | 0.176308 | 0.167393 | 0.176399   | -5.056334    | 0.051717      |
| 5.0   | 0.114438 | 0.109496 | 0.114369   | -4.318374    | -0.060464     |

the income_cat attribute should be removed so the data is back to its original State:'
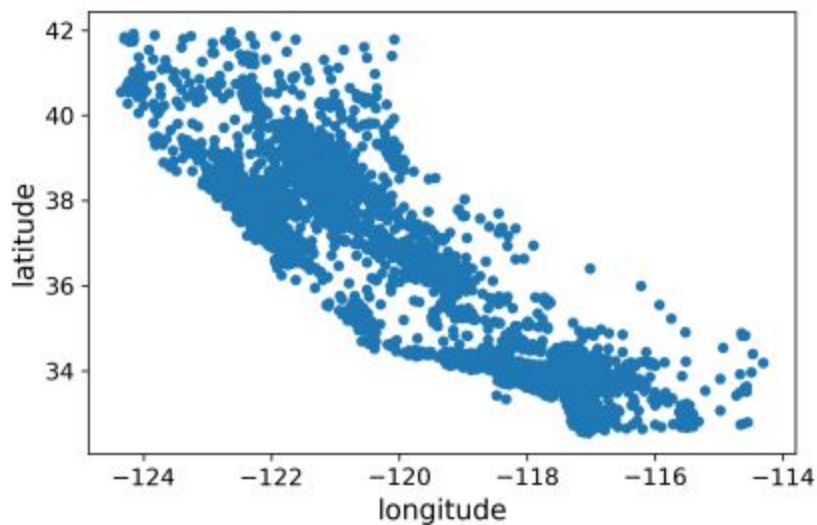
```
for set_ in (strat_train_set, strat_test_set):
    set_.drop("income_cat", axis=1, inplace=True)
```

We spent quite a bit of time on test set generation: this is an often neglected but critical part of a Machine Learning project.
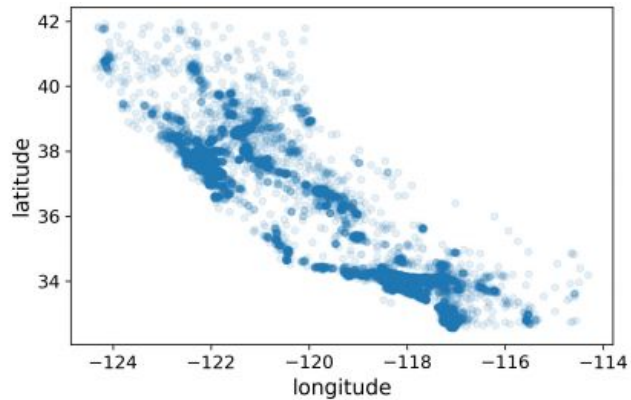
# Data Visualisation

Since there is geographical information (latitude and longitude), a scatterplot of all districts is created to visualize the data.

```
housing.plot(kind="scatter", x="longitude", y="latitude")
```



To make it easier to detect patterns, we highlight the points which are densely populated. For this we set the alpha option to 0.1.
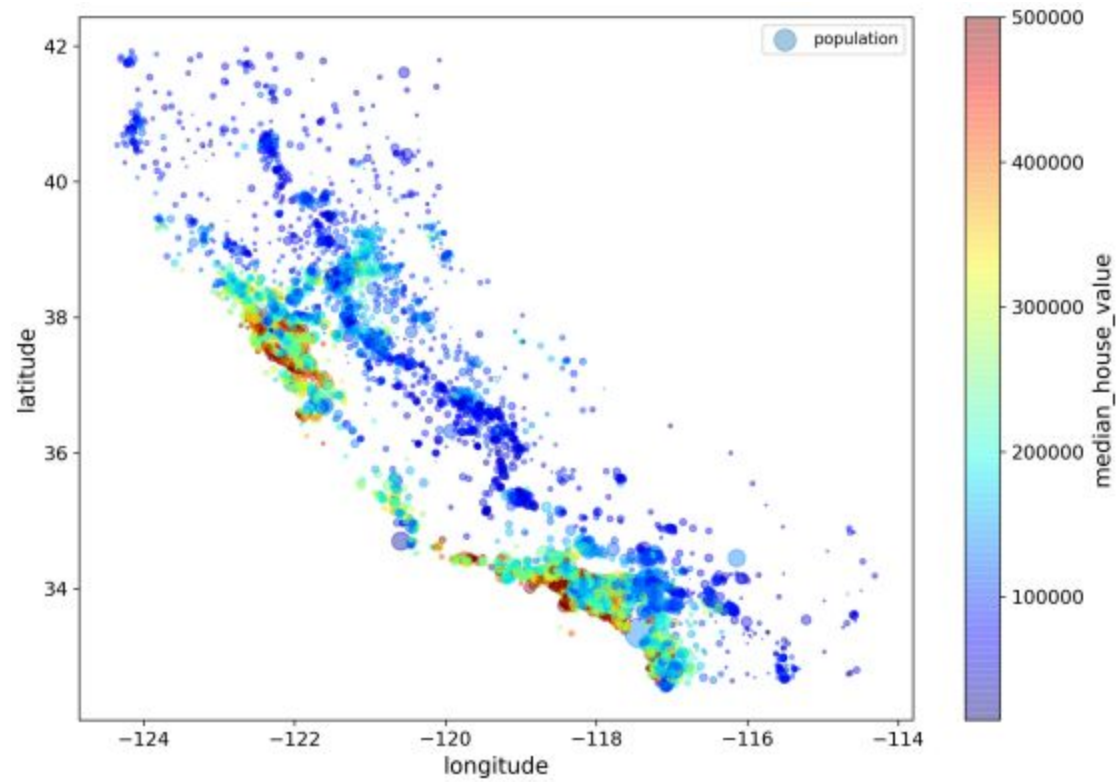
```
housing.plot(kind="scatter", x="longitude", y="latitude", alpha=0.1)
```

 the high-density areas are clearly seen, namely the BayArea and around Los Angeles and San Diego, plus a long line of fairly high density in the Central Valley, in particular around Sacramento and Fresno.

The radius of each circle represents the district's population (option s), and the color represents the price (option c). A predefined color map (option cmap) called jet was used, which ranges from blue low values) to red (high prices):

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

This image tells that the housing prices are very much related to the location (e.g., close to the ocean) and to the population density.

# Correlations in the Data

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

Shown below is how much each attribute correlates with the median house value.
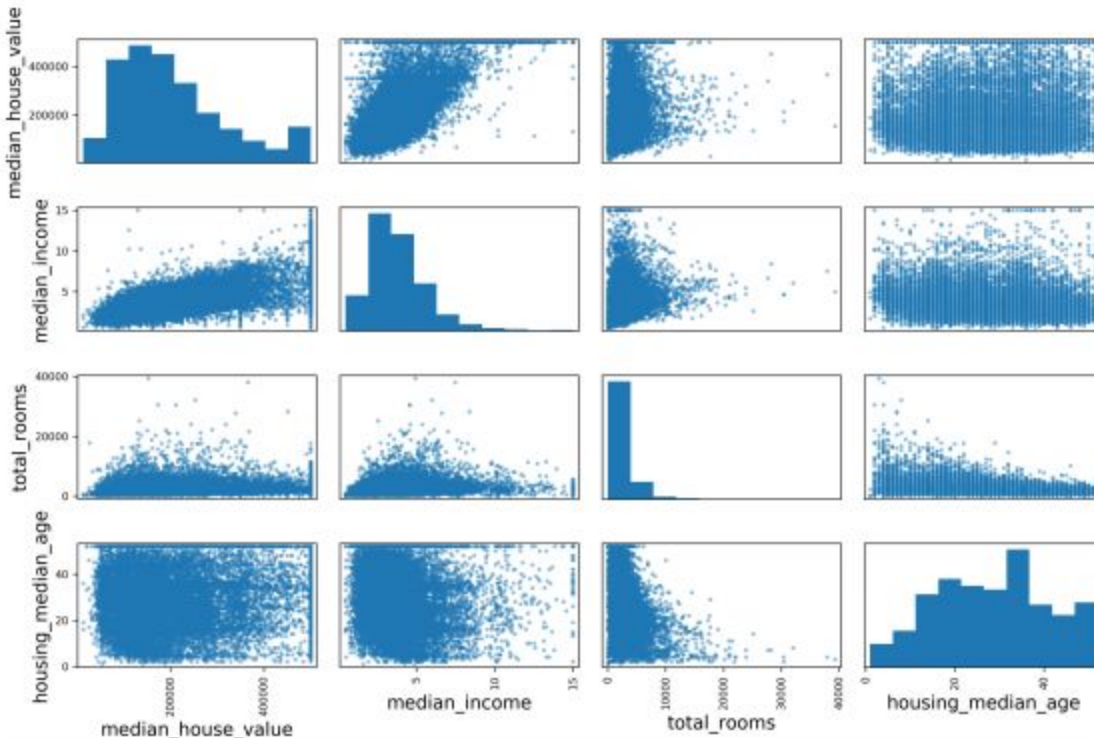
```
>>> corr_matrix["median_house_value"].sort_values(ascending=False)
median_house_value    1.000000
median_income         0.687170
total_rooms           0.135231
housing_median_age    0.114220
households            0.064702
total_bedrooms        0.047865
population            -0.026699
longitude             -0.047279
latitude              -0.142826
Name: median_house_value, dtype: float64
```

The correlation coefficient ranges from –1 to 1. When it is close to 1, it means that there is a strong positive correlation; for example, the median house value tends to go up when the median income goes up. When the coefficient is close to –1, it means that there is a strong negative correlation. A small negative correlation can be seen between the latitude and the median house value (i.e., prices have a slight tendency to go down when you go north). Finally, coefficients close to zero mean that there is no linear correlation.

Another way to check for correlation between attributes is to use Pandas' scatter_matrix function, which plots every numerical attribute against every other numerical attribute. Since there are now 11 numerical attributes, we would get 11^2 =121 plots, which would not fit on a page, so we take a few relevant attributes that seem most correlated with the median housing value.

```
from pandas.plotting import scatter_matrix

attributes = ["median_house_value", "median_income", "total_rooms",
              "housing_median_age"]
scatter_matrix(housing[attributes], figsize=(12, 8))
```
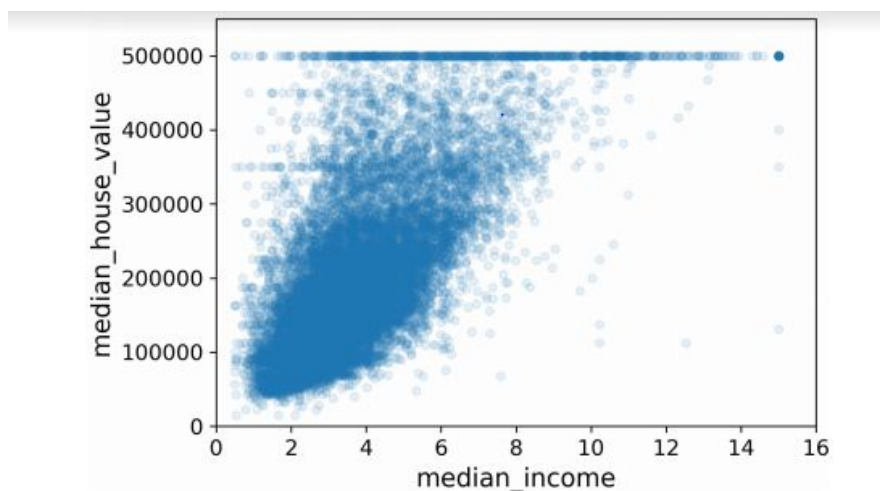
The main diagonal (top left to bottom right) would be full of straight lines if Pandas plotted each variable against itself, which would not be very useful. So instead Pandas displays a histogram of each attribute.

The most promising attribute to predict the median house value is the median Income:

```python
housing.plot(kind="scatter", x="median_income", y="median_house_value",
             alpha=0.1)
```

This plot reveals a few things. First, the correlation is very strong; you can clearly see the upward trend and the points are not too dispersed. Second, the price cap that we noticed earlier is clearly visible as a horizontal line at $500,000. But this plot reveals other less obvious straight lines: a horizontal line around $450,000, another around $350,000, one around $280,000, and a few more below that. We must remove the corresponding districts to prevent the algorithms from learning to reproduce these data anomalies.

There are few attributes which are not very well correlated with median house value. In this case we try to add or remove some attributes which make more sense to the given problem and help us in getting better results. For example, the total number of rooms in a district is not very useful if you don't know how many households there are. What we really want is the number of rooms per household. Similarly, the total number of bedrooms by itself is not very useful. We want to compare it to the number of rooms. And the population per household also is an interesting and useful attribute to include.

```python
housing["rooms_per_household"] = housing["total_rooms"]/housing["households"]
housing["bedrooms_per_room"] = housing["total_bedrooms"]/housing["total_rooms"]
housing["population_per_household"]=housing["population"]/housing["households"]
```

looking at the correlation again:

```
>>> corr_matrix = housing.corr()
>>> corr_matrix["median_house_value"].sort_values(ascending=False)
median_house_value          1.000000
```

```
median_income               0.687160
rooms_per_household         0.146285
total_rooms                 0.135097
housing_median_age          0.114110
households                  0.064506
total_bedrooms              0.047689
population_per_household    -0.021985
population                  -0.026920
longitude                   -0.047432
latitude                    -0.142724
bedrooms_per_room           -0.259984
Name: median_house_value, dtype: float64
```

The new bedrooms_per_room attribute is much more correlated with
the median house value than the total number of rooms or bedrooms. Apparently
houses with a lower bedroom/room ratio tend to be more expensive. The number of
rooms per household is also more informative than the total number of rooms in a
district—obviously the larger the houses, the more expensive they are.

# Preparing Data For Machine Learning

We separate the predictors and the labels since we don't want to apply the same
transformations to the predictors and the target values ( drop() creates a copy of the
data and does not affect strat_train_set):

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

Before training the model with the training dataset, the data needs to be cleaned. A lot
of missing values are present in the dataset, especially in the total_bedrooms attribute.
This can be dealt with by deleting that particular district(sample) or deleting that
particular attribute or by filling missing values with any value. This can be done using
DataFrame's dropna(), drop(), and fillna() methods. Scikit-Learn provides a handy class
to take care of missing values: SimpleImputer.
 First, we need to create a SimpleImputer instance, specifying
that you want to replace each attribute's missing values with the median of that
Attribute.

```
from sklearn.impute import SimpleImputer

imputer = SimpleImputer(strategy="median")
```

Since the median can only be computed on numerical attributes, we need to create a
copy of the data without the text attribute ocean_proximity:

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

Now we can fit the imputer instance to the training data using the fit() method:

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

We also need to deal with text attributes. These obviously cannot be replaced with
median values.

```
>>> housing_cat = housing[["ocean_proximity"]]
>>> housing_cat.head(10)
        ocean_proximity
17606         <1H OCEAN
18632         <1H OCEAN
14650        NEAR OCEAN
3230             INLAND
3555          <1H OCEAN
19480            INLAND
8879          <1H OCEAN
13685            INLAND
4937          <1H OCEAN
4861          <1H OCEAN
```

For this, we can use Scikit-Learn's Ordinal Encoder class:

```
>>> from sklearn.preprocessing import OrdinalEncoder
>>> ordinal_encoder = OrdinalEncoder()
```

```
>>> housing_cat_encoded = ordinal_encoder.fit_transform(housing_cat)
>>> housing_cat_encoded[:10]
array([[0.],
       [0.],
       [4.],
       [1.],
       [0.],
       [1.],
       [0.],
       [1.],
       [0.],
       [0.]])
```

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

One issue with this representation is that ML algorithms will assume that two nearby values are more similar than two distant values. This may be fine in some cases (e.g., for ordered categories such as "bad", "average", "good", "excellent"), but it is obviously not the case for the ocean_proximity column (for example, categories 0 and 4 are clearly more similar than categories 0 and 1). To fix this issue, a common solution is to create one binary attribute per category: one attribute equal to 1 when the category is "<1H OCEAN" (and 0 otherwise), another attribute equal to 1 when the category is

"INLAND" (and 0 otherwise), and so on. This is called one-hot encoding, because only one attribute will be equal to 1 (hot), while the others will be 0 (cold). The new attributes are sometimes called dummy attributes. Scikit-Learn provides a OneHotEncoder class to convert categorical values into one-hot vectors

```
>>> from sklearn.preprocessing import OneHotEncoder
>>> cat_encoder = OneHotEncoder()
>>> housing_cat_1hot = cat_encoder.fit_transform(housing_cat)
>>> housing_cat_1hot
<16512x5 sparse matrix of type '<class 'numpy.float64'>'
  with 16512 stored elements in Compressed Sparse Row format>
```

# Feature Scaling

One of the most important transformations we need to apply to our data is feature scaling. With few exceptions, Machine Learning algorithms don't perform well when the input numerical attributes have very different scales. This is the case for the housing data: the total number of rooms ranges from about 6 to 39,320, while the median incomes only range from 0 to 15.

There are two common ways to get all attributes to have the same scale: min-max scaling and standardization.

Scikit-Learn provides a transformer called StandardScaler for standardization.

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
```

# Transformation Pipelines

As we can see, there are many data transformation steps that need to be executed in the right order. Fortunately, Scikit-Learn provides the Pipeline class to help with such sequences of transformations. Here is a small pipeline for the numerical attributes:

```
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler

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

housing_num_tr = num_pipeline.fit_transform(housing_num)
```

So far, we have handled the categorical columns and the numerical columns separately. It would be more convenient to have a single transformer able to handle all columns, applying the appropriate transformations to each column.

```python
from sklearn.compose import ColumnTransformer

num_attribs = list(housing_num)
cat_attribs = ["ocean_proximity"]

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

housing_prepared = full_pipeline.fit_transform(housing)
```

# Training and Evaluating different Models

**Linear Regression:**

```python
from sklearn.linear_model import LinearRegression

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

Let's evaluate on a few instances from the training set:

```python
>>> some_data = housing.iloc[:5]
>>> some_labels = housing_labels.iloc[:5]
>>> some_data_prepared = full_pipeline.transform(some_data)
>>> print("Predictions:", lin_reg.predict(some_data_prepared))
Predictions: [ 210644.6045  317768.8069  210956.4333  59218.9888  189747.5584]
>>> print("Labels:", list(some_labels))
Labels: [286600.0, 340600.0, 196900.0, 46300.0, 254500.0]
```

Let's measure this regression model's RMSE on the whole train- ing set using Scikit-Learn's mean_squared_error function:

```python
>>> from sklearn.metrics import mean_squared_error
>>> housing_predictions = lin_reg.predict(housing_prepared)
>>> lin_mse = mean_squared_error(housing_labels, housing_predictions)
>>> lin_rmse = np.sqrt(lin_mse)
>>> lin_rmse
68628.19819848922
```

The  model is clearly underfitting the training data as the prediction error is pretty high .
When this happens it can mean that the features do not provide enough information to
make good predictions, or that the model is not powerful enough.

In this scenario we could either try to add more features (e.g., the log of the popula-
tion) or try a more complex model to see how it does.

## Decision Tree Regressor:

```
from sklearn.tree import DecisionTreeRegressor

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

 Evaluating it on the training set:

```
>>> housing_predictions = tree_reg.predict(housing_prepared)
>>> tree_mse = mean_squared_error(housing_labels, housing_predictions)
>>> tree_rmse = np.sqrt(tree_mse)
>>> tree_rmse
0.0
```

**Evaluation using Cross-Validation:**

A great alternative is to use Scikit-Learn's K-fold cross-validation feature. The following
code randomly splits the training set into 10 distinct subsets called folds, then it trains
and evaluates the Decision Tree model 10 times, picking a different fold for evaluation
every time and training on the other 9 folds. The result is an array containing the 10
evaluation scores:

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

Results for Decision Tree Regressor:

```
>>> def display_scores(scores):
...     print("Scores:", scores)
...     print("Mean:", scores.mean())
...     print("Standard deviation:", scores.std())
...
>>> display_scores(tree_rmse_scores)
```

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

The Decision Tree regressor seems to perform worse than the Linear Regression model. Cross-validation not only allows us to get an estimate of the performance of our model, but also a measure of how precise this estimate is (i.e., its standard deviation). The Decision Tree has a score of approximately 71,407, generally ±2,439.

Results for Linear Regression:

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

Scores: [66782.73843989 66960.118071    70347.95244419 74739.57052552
 68031.13388938 71193.84183426 64969.63056405 68281.61137997
 71552.91566558 67665.10082067]
Mean: 69052.46136345083
Standard deviation: 2731.674001798348
```

## Random Forest Regressor:

```
#Random Forests
from sklearn.ensemble import RandomForestRegressor

forest_reg = RandomForestRegressor()
forest_reg.fit(housing_prepared, housing_labels)

housing_predictions = forest_reg.predict(housing_prepared)
forest_mse = mean_squared_error(housing_labels, housing_predictions)
forest_rmse = np.sqrt(forest_mse)
forest_scores = cross_val_score(forest_reg, housing_prepared, housing_labels,
        scoring="neg_mean_squared_error", cv=10)
forest_rmse_scores = np.sqrt(-forest_scores)

display_scores(forest_rmse_scores)
```

```
>>> forest_rmse
18603.515021376355
>>> display_scores(forest_rmse_scores)
Scores: [49519.80364233 47461.9115823  50029.02762854 52325.28068953
 49308.39426421 53446.37892622 48634.8036574  47585.73832311
 53490.10699751 50021.5852922 ]
Mean: 50182.303100336096
Standard deviation: 2097.0810550985693
```

Random Forests work by training many Decision Trees on random subsets of the features, then averaging out their predictions. Building a model on top of many other models is called Ensemble Learning

We should note that the score on the training set is still much lower than on the validation sets, meaning that the model is still overfitting the training set. Possible solutions for overfitting are to simplify the model, constrain it (i.e., regularize it), or get a lot more training data.

# Fine Tuning our Model

### Grid Search

Instead of fiddling with the hyperparameters manually we  used GridSearchCV to search. It evaluates all the possible combinations of hyperparameter values, using cross-validation. For example, the following code searches for the best combination of hyperparameter values for the RandomForestRegressor:

```python
from sklearn.model_selection import GridSearchCV

param_grid = [
    {'n_estimators': [3, 10, 30], 'max_features': [2, 4, 6, 8]},
    {'bootstrap': [False], 'n_estimators': [3, 10], 'max_features': [2, 3, 4]},
  ]

forest_reg = RandomForestRegressor()

grid_search = GridSearchCV(forest_reg, param_grid, cv=5,
                           scoring='neg_mean_squared_error',
                           return_train_score=True)

grid_search.fit(housing_prepared, housing_labels)
```

```
>>> grid_search.best_params_
{'max_features': 8, 'n_estimators': 30}

>>> grid_search.best_estimator_
RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
        max_features=8, max_leaf_nodes=None, min_impurity_decrease=0.0,
        min_impurity_split=None, min_samples_leaf=1,
        min_samples_split=2, min_weight_fraction_leaf=0.0,
        n_estimators=30, n_jobs=None, oob_score=False, random_state=None,
        verbose=0, warm_start=False)
```

In this example, we obtain the best solution by setting the max_features hyperparameter to 8, and the n_estimators hyperparameter to 30. The RMSE score for this combination is 49,682, which is slightly better than the score we got earlier using the default hyperparameter values (which was 50,182).