

Chapter 2: Data handling and Linear Regression models

Werner Bauer



- 1 Loading Data with Pandas
- 2 Visualize and discover the data
- 3 Prepare the data for machine learning
- 4 Linear regression models
- 5 Metrics for Regression models

Pandas provides useful functions to import various file types:

- .csv (comma-separated values) file

```
1 data = pd.read_csv('titanic.csv')
2 # Output:
3      PassengerId  Survived  Pclass  ...    Fare Cabin  Embarked
4  0              1         0       3  ...    7.2500   NaN        S
5  1              2         1       1  ...   71.2833   C85        C
6  2              3         1       3  ...    7.9250   NaN        S
7  3              4         1       1  ...   53.1000  C123        S
8  4              5         0       3  ...    8.0500   NaN        S
```

- .txt files (without column names and separated by spaces):

```
1 data = pd.read_csv('file.txt', sep = "\s+", names = ['column1_name', 'column2_name']) # use '\s+'
   if spacing is not known, assign names to columns, ....)
```

- excel (.xlsx) files

```
1 banking = pd.read_excel('banking.xlsx', sheet_name='bank-additional-full')
2 # Output:
3      age      job  marital  ...  euribor3m  nr_employed  y
4  0     56  housemaid  married  ...     4.857      5191.0  no
5  1     57  services  married  ...     4.857      5191.0  no
6  2     37  services  married  ...     4.857      5191.0  no
```

- Loading files from internet sources via url's

Example: loading the housing.txt file (used later)

We want to read in the following txt file:

```
1 0.00632 18.00 2.310 0 0.5380 6.5750 65.20 4.0900 1 296.0 15.30 396.90 4.98 24.00
2 0.02731 0.00 7.070 0 0.4690 6.4210 78.90 4.9671 2 242.0 17.80 396.90 9.14 21.60
3 0.02729 0.00 7.070 0 0.4690 7.1850 61.10 4.9671 2 242.0 17.80 392.83 4.03 34.70
4 0.03237 0.00 2.180 0 0.4580 6.9980 45.80 6.0622 3 222.0 18.70 394.63 2.94 33.40
5 . . . . . . . . . . . . . .
6 . . . . . . . . . . . . . .
```

We can use `.read_csv()` with options as shown next:

```
1 # as file has no header and values separated by whitespaces
2 df = pd.read_csv('housing.txt', delimiter= '\s+', header = None) # OR
3 df = pd.read_csv('housing.txt', delim_whitespace=True, header = None) # OR
4 df = pd.read_fwf('housing.txt', header = None) # fwf for 'fixed-width formatted'
5 # Output:
6      0      1      2      3      4      ...      9      10      11      12      13
7 0    0.00632 18.0    2.31    0  0.538  ...  296.0  15.3  396.90  4.98  24.0
8 1    0.02731  0.0    7.07    0  0.469  ...  242.0  17.8  396.90  9.14  21.6
9
10 # To set names rather than indices for the COLUMNS:
11 df.columns = ['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT', 'MEDV'] # OR
12 df = pd.read_csv('housing.txt', delim_whitespace=True, header = None, names = ['CRIM', ..., 'MEDV'])
13 # Output:
14      CRIM      ZN  INDUS  CHAS      NOX      ...      TAX  PTRATIO      B  LSTAT  MEDV
15 0    0.00632 18.0    2.31    0  0.538  ...  296.0    15.3  396.90  4.98  24.0
16 1    0.02731  0.0    7.07    0  0.469  ...  242.0    17.8  396.90  9.14  21.6
```

Example (Hands-On ML book 46-47): loading (another) housing data

- We write functions to download and decompress the file:

```
1 import os
2 import tarfile
3 import urllib
4 import pandas as pd
5
6 DOWNLOAD_ROOT = "https://raw.githubusercontent.com/ageron/handson-ml2/master/"
7 HOUSING_PATH = os.path.join("datasets", "housing")
8 HOUSING_URL = DOWNLOAD_ROOT + "datasets/housing/housing.tgz"
9
10 def fetch_housing_data(housing_url=HOUSING_URL, housing_path=HOUSING_PATH):
11     os.makedirs(housing_path, exist_ok=True)
12     tgz_path = os.path.join(housing_path, "housing.tgz")
13     urllib.request.urlretrieve(housing_url, tgz_path)
14     housing_tgz = tarfile.open(tgz_path)
15     housing_tgz.extractall(path=housing_path)
16     housing_tgz.close()
17
18 def load_housing_data(housing_path=HOUSING_PATH):
19     csv_path = os.path.join(housing_path, "housing.csv")
20     return pd.read_csv(csv_path)
21
22 fetch_housing_data() # fetch the data
23 housing = load_housing_data()
```

```
1 housing.head() # Outputs top 5 lines:
2   longitude  latitude  ...  median_house_value  ocean_proximity
3   0    -122.23    37.88  ...           452600.0             NEAR BAY
4   1    -122.22    37.86  ...           358500.0             NEAR BAY
5   2    -122.24    37.85  ...           352100.0             NEAR BAY
6   3    -122.25    37.85  ...           341300.0             NEAR BAY
7   4    -122.25    37.85  ...           342200.0             NEAR BAY
```

Setting and modifying the column types (e.g. boolean, category)

Types are not always correctly recognised:

```
1 print(titanic) # Output:
2      PassengerId  Survived  Pclass  ...    Fare Cabin  Embarked
3      0           1         0       3  ...    7.2500   NaN        S
4      1           2         1       1  ...   71.2833   C85        C
5      2           3         1       3  ...    7.9250   NaN        S
6 print(titanic.dtypes) # Output:
7 PassengerId      int64 # OK
8 Survived        int64 # should be categorical
9 ...
10
11 titanic = titanic.astype({'Survived' : 'category'}, copy=True) # DEFAULT: copy=True
12 print(titanic.dtypes) # Output:
13 PassengerId      int64 # OK
14 Survived        category # OK
```

```
1 print(housing)
2      longitude  latitude  ...  median_house_value  ocean_proximity
3      0      -122.23    37.88  ...           452600.0      NEAR BAY
4      1      -122.22    37.86  ...           358500.0      NEAR BAY
5      2      -122.24    37.85  ...           352100.0      NEAR BAY
6 print(housing.dtypes)
7 longitude      float64 # OK
8 ...
9 ocean_proximity    object # Could be any Python object, but repetitive entries suggest it is
   categorical!
10 housing['ocean_proximity'].value_counts()
11 <1H OCEAN          9136
12 INLAND            6551
13 NEAR OCEAN        2658
14 NEAR BAY          2290
15 ISLAND             5
16 # THEN: modify type as suggested above! Discussed later again!
```

- Use `.head()` to look at top 5 lines
- `.describe()` gives overview of each numerical attribute

```
1      longitude      latitude  ...  median_income  median_house_value
2  count    20640.000000    20640.000000  ...    20640.000000         20640.000000
3  mean     -119.569704     35.631861  ...         3.870671         206855.816909
4  std        2.003532        2.135952  ...         1.899822         115395.615874
5  min     -124.350000     32.540000  ...         0.499900         14999.000000
6  25%     -121.800000     33.930000  ...         2.563400         119600.000000
7  50%     -118.490000     34.260000  ...         3.534800         179700.000000
8  75%     -118.010000     37.710000  ...         4.743250         264725.000000
9  max     -114.310000     41.950000  ...        15.000100         500001.000000
```

- count, mean, min, max are self-explanatory
- std stands for (standard deviation)
- 25%, 50%, 75% for percentile (e.g. 25% of the districts have a median_income lower than 2.563400)

- The `.info()` method gives quick description of data

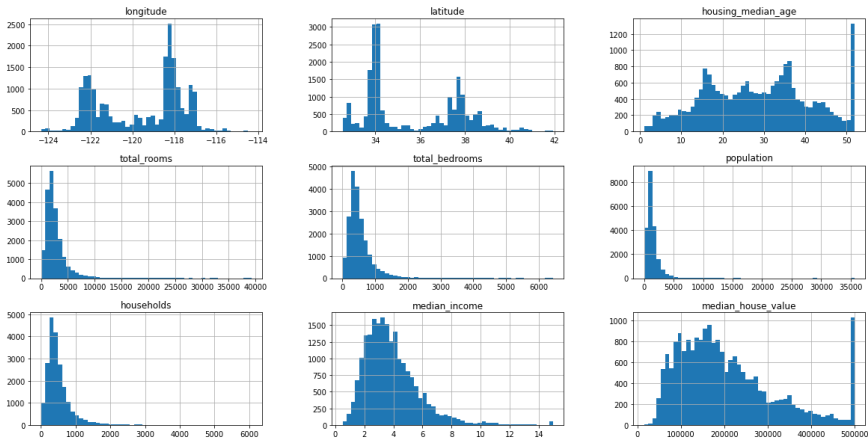
```
1 housing = load_housing_data()
2 print(housing.info()) # Output:
3 <class 'pandas.core.frame.DataFrame'>
4 RangeIndex: 20640 entries, 0 to 20639
5 Data columns (total 10 columns):
6  #   Column                Non-Null Count  Dtype
7
8  0   longitude             20640 non-null  float64
9  1   latitude              20640 non-null  float64
10  2   housing_median_age    20640 non-null  float64
11  3   total_rooms           20640 non-null  float64
12  4   total_bedrooms        20433 non-null  float64
13  5   population            20640 non-null  float64
14  6   households            20640 non-null  float64
15  7   median_income         20640 non-null  float64
16  8   median_house_value    20640 non-null  float64
17  9   ocean_proximity       20640 non-null  object
18 dtypes: float64(9), object(1)
19 memory usage: 1.6+ MB
20 None
```

Gives total number of rows and columns, the attribute types and number of non-null values.

- Note: `total_bedrooms` is missing some values (we will need to address this)

- Using a histogram plot to view numerical attributes:

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



At this stage it is time to set aside part of the data. Why?

- Having test data is important to avoid overfitting of the ML model
- If your brain detects pattern when looking at the test data, you might realize pattern that influences your choice of ML model
- When you estimate the generalization error using the test set, your estimate will be too optimistic and it does not generalize this well to new data (data snooping bias)

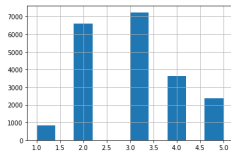
Creating a test set (cf. Hands-On ML book on pages 51-55)

We can use a scikit-learn function to split datasets.

```
1 from sklearn.model_selection import train_test_split
2
3 # purely random sampling method:
4 train_set, test_set = train_test_split(housing, test_size=0.2, random_state=42)
```

- Potentially with smaller datasets, random sampling introduces bias! Hence:
we must make sure that the test set is representative for the various categories. E.g. consider the median_income attribute:

```
1 housing["income_cat"] = pd.cut(housing["median_income"],
2                               bins=[0., 1.5, 3.0, 4.5, 6., np.inf],
3                               labels=[1, 2, 3, 4, 5])
4
5 from sklearn.model_selection import StratifiedShuffleSplit
6
7 split = StratifiedShuffleSplit(n_splits=1, test_size=0.2,
8                               random_state=42)
9 for train_index, test_index in split.split(housing, housing["
10 income_cat"]):
11     strat_train_set = housing.loc[train_index]
12     strat_test_set = housing.loc[test_index]
13
14 # remove income_cat to bring data back to original state:
15 for set_ in (strat_train_set, strat_test_set):
16     set_.drop("income_cat", axis=1, inplace=True)
```



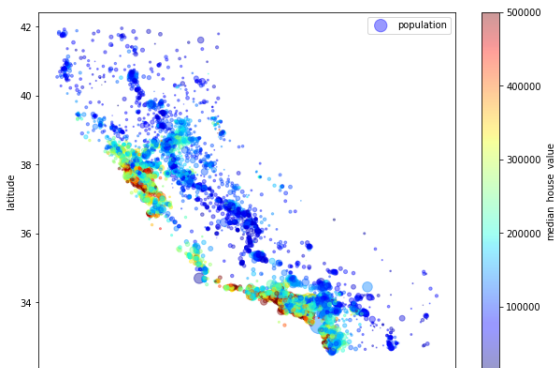
```
1 ??_train_set['income_cat'].
2 value_counts()/len(??_train_set)
```

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

- 1 Loading Data with Pandas
- 2 Visualize and discover the data
- 3 Prepare the data for machine learning
- 4 Linear regression models
- 5 Metrics for Regression models

Example: scatter plot of housing dataset

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



California housing
prices: red is
expensive, blue is
cheap;
large circles for
areas with large
population

- Housing prices are related very much to location!

Looking for correlations: *standard correlation coefficient*

- If dataset is not too large: calculate standard correlation coefficient (Pearson's r):

```
1 corr_matrix = housing.corr()
2 print(corr_matrix) # Output:
3      longitude  latitude  ...  median_income  median_house_value
4 longitude      1.000000 -0.924478  ...      -0.019615      -0.047466
5 latitude      -0.924478  1.000000  ...      -0.075146      -0.142673
6 housing_median_age -0.105823  0.005737  ...      -0.111315       0.114146
7 ...
8 households        0.063146 -0.077765  ...        0.010869       0.064590
9 median_income     -0.019615 -0.075146  ...        1.000000       0.687151
10 median_house_value -0.047466 -0.142673  ...        0.687151       1.000000
```

- select an attribute of interest and sort values:

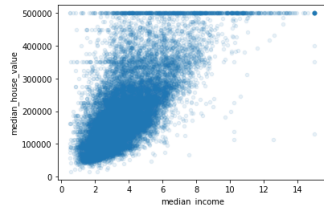
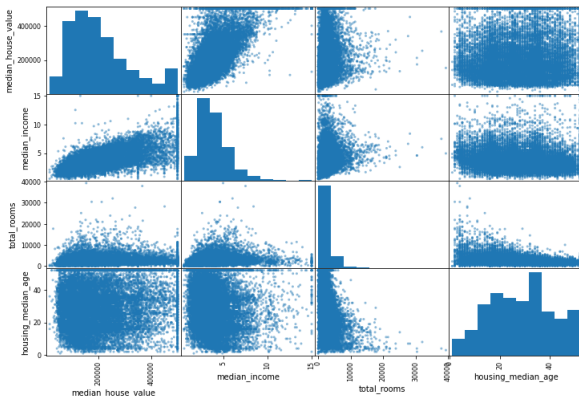
```
1 corr_matrix['median_house_value'].sort_values(ascending=False)
2 # Output:
3 median_house_value    1.000000
4 median_income         0.687151
5 total_rooms           0.135140
6 housing_median_age    0.114146
7 households            0.064590
8 total_bedrooms        0.047781
9 population            -0.026882
10 longitude            -0.047466
11 latitude             -0.142673
```

- Range $[-1,1]$: -1 negative, +1 positive, 0 no correlation
- works only for lin. correlation: see examples in Hands-On ML bookpage 59

Looking for correlations: pandas *scatter matrix*

- Plot attributes most correlated to median_house_value

```
1 from pandas.plotting import scatter_matrix
2
3 attributes = ["median_house_value", "median_income", "total_rooms", "housing_median_age"]
4 scatter_matrix(housing[attributes], figsize=(12, 8)) # LEFT
5
6 housing.plot(kind="scatter", x="median_income", y="median_house_value", alpha=0.1) # RIGHT
```



- Note: diagonal shows histograms rather than correlation

Testing out attribute combinations

- Try out some attribute combinations to increase correlation with the target attribute, here median_house_value.
- Example: total number of rooms in a district not very useful without knowing how many households there are!
- Following attribute combinations seem promising:

```
1 housing["rooms_per_household"] = housing["total_rooms"]/housing["households"]
2 housing["bedrooms_per_room"] = housing["total_bedrooms"]/housing["total_rooms"]
3 housing["population_per_household"]=housing["population"]/housing["households"]
4
5 corr_matrix = housing.corr()
6 corr_matrix["median_house_value"].sort_values(ascending=False)
7 # Output:
8 median_house_value      1.000000
9 median_income           0.687151
10 rooms_per_household     0.146255
11 total_rooms             0.135140
12 housing_median_age      0.114146
13 households              0.064590
14 total_bedrooms          0.047781
15 population_per_household -0.021991
16 population              -0.026882
17 longitude               -0.047466
18 latitude                -0.142673
19 bedrooms_per_room       -0.259952
```

We found via bedrooms_per_room: houses with lower bedroom to room ratio are more expensive!

- 1 Loading Data with Pandas
- 2 Visualize and discover the data
- 3 Prepare the data for machine learning**
- 4 Linear regression models
- 5 Metrics for Regression models

- 1 Loading Data with Pandas
- 2 Visualize and discover the data
- 3 Prepare the data for machine learning
- 4 Linear regression models**
- 5 Metrics for Regression models

- 1 Loading Data with Pandas
- 2 Visualize and discover the data
- 3 Prepare the data for machine learning
- 4 Linear regression models
- 5 Metrics for Regression models**