

DAT565/DIT407 Assignment 2

Sebastian Miles
miless@chalmers.se

Olle Lapidus
ollelap@chalmers.se

2024-09-12

The goal of this assignment is to practice data cleaning by extracting data from a source that is not in the format we want. Then, we produce plots using the extracted data. The code for generating the .csv file containing the extracted data can be found in the section **Extraction Code**. And the code for generating the plots is located in the section **Plot Code**.

After extraction, we analyze the closing prices in 2022 by calculating the five-number summary. The summary is shown in Table 1. Next, we plot the histogram of the closing prices in 2022. The number of bins was manually chosen as $N = 25$ so that the bins would not contain any gaps and are not overly grouped. Thus, $N = 25$ was considered an appropriate number of bins.

When analyzing Figure 2, we can clearly see a general trend: the larger the living area, the higher the closing price. Although there are some deviations, the majority of listings follow this trend. This makes sense intuitively, as a larger area typically results in a higher cost. In Figure 3, we observe that as the living area increases, the number of rooms also increases.

Minimum	1,650,000
First Quartile (25%)	4,012,500
Median (50%)	5,000,000
Third Quartile (75%)	5,795,000
Maximum	10,500,000

Table 1: Five-Number Summary of Closing Prices in 2022 (in SEK)

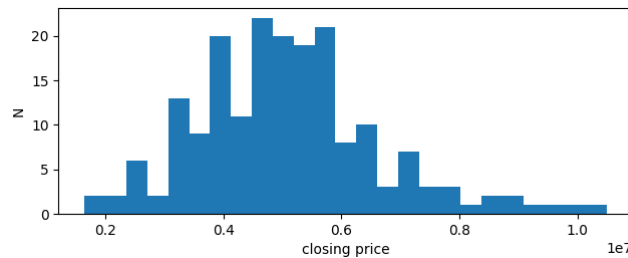


Figure 1: Histogram showing the closing prices of houses in 2022.

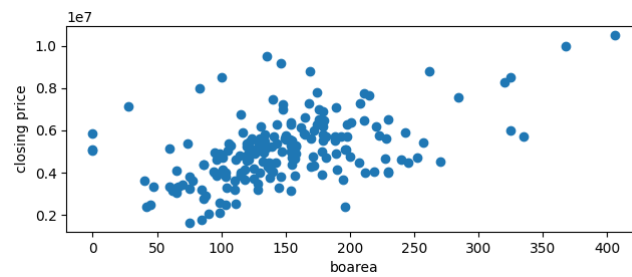


Figure 2: Scatter plot over the closing prices of 2022 in Kungälv.

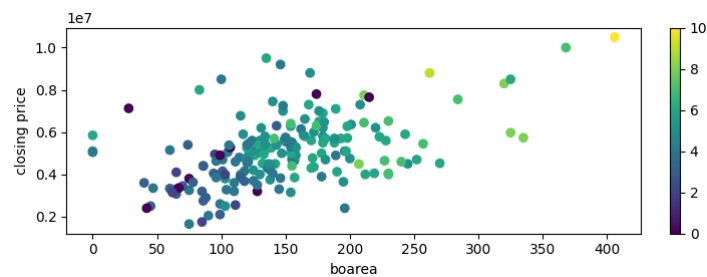


Figure 3: Scatter plot over closing prices including the number of rooms during 2022 in Kungälv.

Extraction code

```

1 from bs4 import BeautifulSoup
2 import re
3
4 import datetime
5 import pandas as pd
6 import os
7
8 def ParseDate(date):
9     monthDict = {
10         'januari' : 1,
11         'februari' : 2,
12         'mars' : 3,
13         'april' : 4,
14         'maj' : 5,
15         'juni' : 6,
16         'juli' : 7,
17         'augusti' : 8,
18         'september' : 9,
19         'oktober' : 10,
20         'november' : 11,

```

```

21         'december': 12
22     }
23     parts = date.split("_")
24
25     return datetime.datetime(day=int(parts[1]), month=
        ↳ monthDict[parts[2]], year=int(parts[3]))
26
27 def Process(html_doc):
28     soup = BeautifulSoup(html_doc, 'html.parser')
29
30     cols = 7
31     data = [[]*cols for i in range(cols)]
32     for listing in soup.find_all('li', class_='sold-
        ↳ results__normal-hit'):
33         data[0].append(ParseDate(listing.find('span',
            ↳ class_='hcl-label--sold-at').text.strip()
            ↳ ()))
34         data[1].append(listing.find('h2', class_='sold
            ↳ -property-listing__heading').text.strip()
            ↳ ())
35         # Get raw location
36         location = listing.select('div.sold-property-
            ↳ listing__location>div')[0].text.strip()
37         location = re.sub(r'VillaVilla\s+|\n|_', '',
            ↳ location) # Remove garbage
38         data[2].append(location)
39
40         area = listing.select('div.sold-property-
            ↳ listing__area')[0].text.strip()
41         text = re.sub(r'\s+', '', area).split('^2')
42         rooms = re.sub(r'rum', '', text[len(text)-1])
43         area = 0
44         for i in range(len(text)-1):
45             for temp in text[i].split('+'):
46                 area += float(re.sub(',+', '.', re.sub
                    ↳ ('m+|m^2+', '', temp)))
47
48         data[3].append(area)
49         data[4].append(rooms)
50
51         if(listing.find('div', class_='sold-property-
            ↳ listing__land-area')):
52             landarea = listing.select('div.sold-
                ↳ property-listing__land-area')[0].
                ↳ text.strip()
53             landarea = re.sub(r'm^2_tomt|\s+', '',
                ↳ landarea) # Remove garbage
54             data[5].append(landarea)
55         else:

```

```

56         data[5].append(0) # Land area doesnt exist
57         ↪ , maybe an apartment?
58     price = re.sub(r'\s+|Slutpris|kr', '', listing
59         ↪ .find('span', 'hcl-text').text.strip())
60     data[6].append(price)
61
62     df = pd.DataFrame({
63         'date': data[0],
64         'address': data[1],
65         'location': data[2],
66         'boarea': data[3],
67         'rooms': data[4],
68         'landarea': data[5],
69         'closing_price': data[6],
70     })
71     return df
72
73 folder_path = 'kungalv_slutpriser'
74 files = os.listdir(folder_path)
75
76 table = list()
77
78 # For each filename
79 for file_name in files:
80     file_path = os.path.join(folder_path, file_name)
81     f = open(file_path, 'r', encoding='utf-8')
82     table.append(Process(f.read()))
83
84 df = pd.concat(table)
85
86 df.to_csv('listings.csv', index = None)

```

Plot code

```

1  import pandas as pd
2  import matplotlib.pyplot as plt
3  import numpy as np
4
5  # Read from file
6  df = pd.read_csv('listings.csv')
7
8  df2 = (df[df['date'].str.contains('2022')])
9
10 print(df2['closing_price'].describe().loc[['min', '25%',
11     ↪ , '50%', '75%', 'max']])

```

```

12 fig, axs = plt.subplots(3, constrained_layout = True ,
    ↪     figsize = (7 ,8))
13
14 rooms = np.nan_to_num(np.array(pd.to_numeric(df2['
    ↪     rooms'])))
15 axs [0].hist (df2['closing_price'], bins = 25)
16 sp = axs[2].scatter(df2['boarea'], df2['closing_price'
    ↪     ], c = rooms , cmap = "viridis")
17
18 axs[0].set_xlabel("closing_price")
19 axs[0].set_ylabel("N")
20 axs[1].scatter(df2['boarea'], df2['closing_price'])
21 axs[1].set_xlabel("boarea")
22 axs[1].set_ylabel("Closing_price")
23 axs[2].set_xlabel("boarea")
24 axs[2].set_ylabel("Closing_price")
25
26 fig.colorbar(sp)
27 plt.show()

```