# DIT407 Introduction to data science and AI - Assignment 2

Reynir Siik
reynir.siik@gmail.com

Franco Zambon
guszamfr@student.gu.se

2024-01-31

## Abstract

In this assignment, we practice data cleaning by extracting the interesting data from a source that contains the data, not quite in the format we want. In the second part, we then produce some plots about the data.

# 1 Problem 1: Scraping house prices

## 1.1 Introduction

*Abbreviated citation from assignment:*

The file `kungalv slutpriser.tar.gz` contains the closing prices of houses, sold in Kungälv municipality between 2013 and 2023. The data has been scraped from [1] in October 2023. Your task will be to convert the data into a form that can be analyzed. The tarball contains a number of HTML files. Use Beautiful Soup2 to parse the HTML and extract the relevant pieces of information for each advertisement of a closing price.

The data of interest is: date of sale, address , location, building living area, building utility area, number of rooms, plot area and closing price.

The data consists of closed listings that is generated as the web page is loaded into the browser. There is no guarantee that every value is present, and some values can have a decimal value.

The `html` data from every file is read into an object `soup` by the module `BeautifulSoup`.

```
36  while i<40:
37      i += 1
38      print("File ", i)
39
40      fileSequence = str(i).rjust(2, '0')
41      with open(
42          "Uppgift02\\code\\kungalv_slutpriser\\kungalv_slutpris_page_"
43          + fileSequence
44          + ".html",
45          encoding="utf8") as fp:
46          soup = BeautifulSoup(fp, 'html.parser')
```

Listing 1: Reading the files

## 1.2 Overview of dataset

Every listing is contained in a separate list item which makes it easy to filter them out into separate items. With `soup.findall` the listings are filtered out into a new object `p` as separate items.

```
49   ## Divide into sold listings ##########################################
50       p = soup.find_all('li', class_='sold-results__normal-hit')
```

Listing 2: Filtering out the listings

Within every item most data is contained under specific "labels". The most challenging values to retrieve is building living area, building utility area and number of rooms, since they are not separated by any labels. The values may, or may not be present, and they can also consist of decimal values.

```
94           tag = t.find(
95               'div',
96               class_="sold-property-listing__subheading␣"
97               "sold-property-listing__area")
```

Listing 3: Getting the area section

## 1.3 Retrieving building living area, building utility area and number of rooms

With regex all digit values are retrieved into a list variable.

```
99           pArea = re.findall(r"\d+[,]?\d?", tag.get_text(strip=True))
```

Listing 4: Getting the area values into a list

By looking at "telltales" `[+]` and `[rum]` the data in the list is identified. Data is assigned to the relevant variable and decimal comma is changed to decimal point before writing it to csv-file.

```
102          bia = re.search(r"[+]", tag.get_text(strip=True))
103          rum = re.search(r"rum", tag.get_text(strip=True))
```

Listing 5: Analyzing presence of data

## 1.4 Transforming the date field

The closing date is written in the format; `Såld 21 juni 2023`. The month is in Swedish and does not work well with `datetime`. By changing month to English and three letter abbreviation the problem is solved.

```
168          ## Transform the date
169          ## This could be a function... I know...
170          monthdict =     {
171              "januari": "Jan",
172              "februari": "Feb",
173              "mars": "Mar",

182              "december": "Dec",
183          }
184
185          monthlist = [
186              "januari",
187              "februari",
```

```
195             "oktober",
196             "november",
197             "december"
198         ]
199         for month in monthlist:
200             if date_of_sale.find(month)>0:
201                 date_of_sale = date_of_sale.replace(
202                     month, monthdict[month])
203         ####################
```

Listing 6: Transforming date

## 1.5 Writing the cvs file

By opening a file for writing and make use of the **cvs** module, writing a cvs file is fairly trivial. The file is opened in the beginning of the code, and a **cvs.writer** object with the file descriptor is defined. The first row with the headings is written to the file.

```
22  csvfile = open(csv_file_path, 'w', newline='')
23  writer = csv.writer(csvfile)
```

Listing 7: Opening CVS file

When all values are retrieved for an individual listing the data is written to the file by composing a list with the values, and writing it to the file with the **writer** object. Finally the file is closed.

```
213         ## Put the data in the csv file
214         myList = (
215             datetime.datetime.strptime(date_of_sale, r"%d %b %Y").date(),
216             PropAddress,
217             Location_of_the_estate,
218             Area_boarea,
219             Area_biarea,
220             The_number_of_rooms,
221             Area_of_the_plot,
222             Closing_price)
223         writer.writerow(myList)
```

Listing 8: Writing the CVS file

# 2 Problem 2: Analyzing 2022 house sales

## 2.1 Introduction

*Abbreviated citation from assignment:*

We will now analyze the data, and produce some plots. Select all houses that were sold in 2022, and answer all points below. Include the plots you create in your report and the answers to any other questions.

## 2.2 The five-numbers summary

The five-number summary is a set of descriptive statistics that provides information about a dataset. It consists of the five most important sample percentiles:

- the sample minimum (smallest observation)

- the lower quartile or first quartile

- the median (the middle value)

- the upper quartile or third quartile

- the sample maximum (largest observation)

In order for these statistics to exist, the observations must be from a univariate variable that can be measured on an ordinal, interval or ratio scale.[2]
You can see the five-numbers summary that was generated selecting all the houses that got sold in 2022 in figure 1.

|  | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|
| date_of_sale | NaN | NaN | NaN | NaN | NaN |
| adress | NaN | NaN | NaN | NaN | NaN |
| location | NaN | NaN | NaN | NaN | NaN |
| boarea | 28.0 | 99.0 | 121.0 | 146.0 | 325.0 |
| biarea | 2.0 | 21.0 | 48.0 | 75.0 | 174.0 |
| rum | 2.0 | 4.0 | 5.0 | 6.0 | 10.0 |
| tomtarea | 127.0 | 600.75 | 1118.5 | 1718.5 | 47500.0 |
| slutpris | 1650000.0 | 4012500.0 | 5000000.0 | 5795000.0 | 10500000.0 |

Figure 1: Five-numbers summary

From this summary we can see how 50% of the houses has between 4 and 6 rooms, a price that ranges between 4 and 5.8 millions SEK and an inhabitable area between 99 and 146 sqm, a non-inhabitable area between 21 and 75 sqm and a plot area between 600 and 1700 sqm. Besides that there could be some house with extreme values, such as one with a plot area of 47500 sqm.
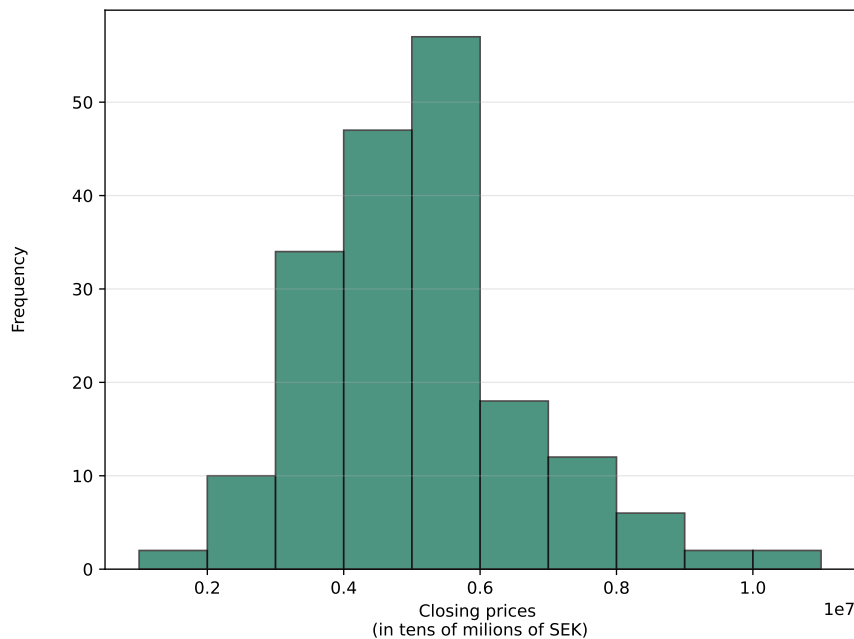
Figure 2: Distribution of closing prices in 2022

## 2.3 The histogram

Plotting the closing prices of all the houses sold in 2022 we obtained the histogram of figure 2.

The number of bins has been chosen to avoid over- or under-fitting the data and each bin was chosen to start and to end with a round number to increase readability. The histogram suggests normally-distributed data with a heavier right tail.

## 2.4 The prices vs area scatter plot

In figure 3 we plotted the closing prices and the boarea (the total inhabitable area) of the houses sold in 2022.

In this graph we could see how, as expected, the boarea is linearly correlated to the closing price. However there's some house which seems to have a price very high compared to its boarea (for example the house that has an area lower than 50 sqm but a price higher than 7 millions SEK) and this is probably caused by other factors such as the position of the house.

## 2.5 The prices vs area scatter plot by number of rooms

Finally in figure 4 we plotted the same graph of above but this time we distinguished the houses by the number of rooms. The color of the dots helps doing that.

As expected we can see in this graph how the houses with higher boareas tend to have more rooms.
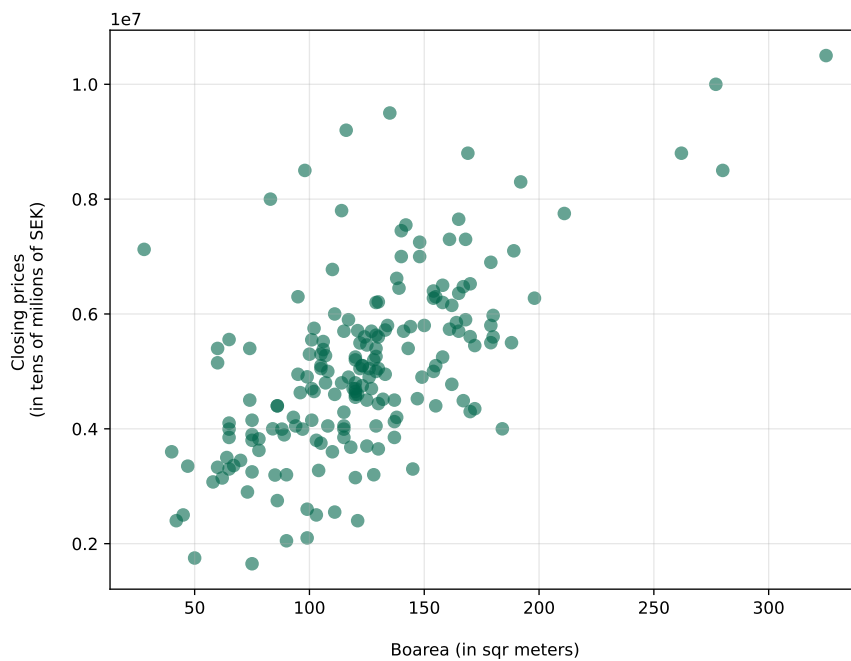
Figure 3: Scatter plot of the boarea and the closing prices



Figure 4: Scatter plot of the boarea and the closing prices differentiated by number of rooms

# References

[1] Hemnet, *Slutpriser*, Made available on Canvas, From assignment handout on Canvas, Oct. 2023. [Online]. Available: `https://www.hemnet.se/salda/bostader?location_ids%5B%5D=17973&item_types%5B%5D=villa`.

[2]  *Five-number summary*, 2023. [Online]. Available: `https : / / en . wikipedia . org / wiki/Five-number_summary`.

# A    Complete source code

This is the complete listing of the source code.

```python
1  # Reynir Siik, Franco Zambon
2  # DIT407 lp3 2024-01-28
3  # Assignment 2
4  # Problem 1 and 2
5  # Scraping Hemnet
6  #
7  ## ####################################################################
8  import datetime
9  from math import nan
10 from os import replace
11 from bs4 import BeautifulSoup
12 import re
13 from matplotlib.dates import date2num
14 import matplotlib.pyplot as plt
15 import numpy  as np
16 import pandas as pd
17 import csv
18 ## ####################################################################
19
20 csv_file_path = "Uppgift02\\code\\huspriser.csv"
21
22 csvfile = open(csv_file_path, 'w', newline='')
23 writer = csv.writer(csvfile)
24 writer.writerow([
25     'date_of_sale',
26     'adress',
27     'location',
28     'boarea',
29     'biarea',
30     'rum',
31     'tomtarea',
32     'slutpris'])
33
34 i = 0
35
36 while i<40:
37     i += 1
38     print("File␣", i)
39
40     fileSequence = str(i).rjust(2, '0')
41     with open(
42         "Uppgift02\\code\\kungalv_slutpriser\\kungalv_slutpris_page_"
43         + fileSequence
44         + ".html",
45         encoding="utf8") as fp:
46         soup = BeautifulSoup(fp, 'html.parser')
47     fp.close()
48
49 ## Divide into sold listings #########################################
50     p = soup.find_all('li', class_='sold-results__normal-hit')
51
52 ## Extract ###########################################################
53     for t in p:
54         ## Date #############
55         tag = t.find(
56             'span',
57             class_="hcl-label␣hcl-label--state␣hcl-label--sold-at")
58         lbl = re.findall(
59             r"\d?\d{1}\s\w+\s\d{4}",
```

```
60              tag.get_text(strip=True))
61
62          date_of_sale = str(lbl).strip("'[]")
63          tag.clear()
64
65          ## Address ##########
66          tag = t.find(
67              'h2',
68              class_="sold-property-listing__heading␣"
69              "qa-selling-price-title␣hcl-card__title")
70          PropAddress = str(
71              re.findall(r".+", tag.get_text(strip=True))).strip("'[]")
72          tag.clear()
73
74          ## Remove tag that will clutter the data
75          tag = t.find('span',
76                      class_="property-icon␣property-icon--result")
77          tag.clear()
78
79          ## Location #########
80          tag = t.find('div', class_="sold-property-listing__location")
81          Location_of_the_estate = str(
82              re.sub(r"\s+", "_",
83                      str(re.findall(r".+", tag.get_text(strip=True)))))
84          Location_of_the_estate = str(
85              re.sub("_", "␣",
86                      re.sub("__",",␣",
87                              re.sub(
88                                  r"[,][']['][,]|\W",
89                                  "",
90                                  Location_of_the_estate))))
91          tag.clear()
92
93          ## Number of rooms, boarea, biarea ##
94          tag = t.find(
95              'div',
96              class_="sold-property-listing__subheading␣"
97              "sold-property-listing__area")
98
99          pArea = re.findall(r"\d+[,]?\d?",  tag.get_text(strip=True))
100         lenP = len(pArea)
101
102         bia = re.search(r"[+]", tag.get_text(strip=True))
103         rum = re.search(r"rum", tag.get_text(strip=True))
104
105         bia_t = not(bia == None)
106         rum_t = not(rum == None)
107         x = (lenP, bia_t, rum_t)
108         match x:
109             case (3, True, True):
110                 Area_boarea = pArea[0]
111                 Area_biarea = pArea[1]
112                 The_number_of_rooms = pArea[2]
113             case (2, False, True):
114                 Area_boarea = pArea[0]
115                 Area_biarea = np.NaN
116                 The_number_of_rooms = pArea[1]
117             case (2, True, True):
118                 Area_boarea = np.NaN
119                 Area_biarea = pArea[0]
120                 The_number_of_rooms = pArea[1]
121             case (2, True, False):
122                 Area_boarea = pArea[0]
```

```
123                    Area_biarea = pArea[1]
124                    The_number_of_rooms = np.NaN
125                case (1, False, False):
126                    Area_boarea = pArea[0]
127                    Area_biarea = np.NaN
128                    The_number_of_rooms = np.NaN
129                case (1, False, True):
130                    Area_boarea = np.NaN
131                    Area_biarea = np.NaN
132                    The_number_of_rooms = pArea[0]
133                case (1, True, False):
134                    Area_boarea = np.NaN
135                    Area_biarea = pArea[0]
136                    The_number_of_rooms = np.NaN
137                case (0, bia_t, rum_t):
138                    Area_boarea = np.NaN
139                    Area_biarea = np.NaN
140                    The_number_of_rooms = np.NaN
141                case _:
142                    Area_boarea = np.NaN
143                    Area_biarea = np.NaN
144                    The_number_of_rooms = np.NaN
145                    print('ERROR ERROR ERROR ERROR ERROR ERROR ERROR ')
146
147            ## Decimal comma to decimal point
148            Area_boarea = re.sub("[,]", "." ,str(Area_boarea))
149            Area_biarea = re.sub("[,]", "." ,str(Area_biarea))
150            The_number_of_rooms = re.sub("[,]", "." ,str(The_number_of_rooms))
151
152            ## Area of the plot #
153            tag = t.find(
154                'div',
155                class_="sold-property-listing__land-area")
156            if len(str(tag))>4:
157                lbl = str(
158                    re.findall(r".+",
159                            tag.get_text(strip=True).strip(r"'[]")))
160                Area_of_the_plot = str(
161                    re.findall(r"\d+",
162                            re.sub(r"\W*|[x][a][0]|[t]{1}[o,m,t]+", "",
163                                str(lbl)))).strip(r"'[]")
164                tag.clear()
165            else:
166                Area_of_the_plot = np.NaN
167
168            ## Transform the date
169            ## This could be a function... I know...
170            monthdict =     {
171                "januari": "Jan",
172                "februari": "Feb",
173                "mars": "Mar",
174                "april": "Apr",
175                "maj": "May",
176                "juni": "Jun",
177                "juli": "Jul",
178                "augusti": "Aug",
179                "september": "Sep",
180                "oktober": "Oct",
181                "november": "Nov",
182                "december": "Dec",
183            }
184
185            monthlist = [
```

```
186                "januari",
187                "februari",
188                "mars",
189                "april",
190                "maj",
191                "juni",
192                "juli",
193                "augusti",
194                "september",
195                "oktober",
196                "november",
197                "december"
198            ]
199            for month in monthlist:
200                if date_of_sale.find(month)>0:
201                    date_of_sale = date_of_sale.replace(
202                        month, monthdict[month])
203            #####################

205            ## Closing price ####
206            tag = t.find('span', class_="hcl-text hcl-text--medium")
207            lbl = re.findall(r".+", tag.get_text(strip=True))
208            Closing_price = re.sub(
209                r"\W*|[x][a][0]|[S]{1}[l,u,t,p,r,i,s]+|[K|k][R|r]",
210                "",
211                str(lbl))

213            ## Put the data in the csv file
214            myList = (
215                datetime.datetime.strptime(date_of_sale, r"%d %b %Y").date(),
216                PropAddress,
217                Location_of_the_estate,
218                Area_boarea,
219                Area_biarea,
220                The_number_of_rooms,
221                Area_of_the_plot,
222                Closing_price)
223            writer.writerow(myList)
224 ## ####################################################################

226 ## Close the file ####################################################
227 csvfile.close()
228 ## ####################################################################

231 ## Problem 2 ########################################################

233 try:
234     df = pd.read_csv(csv_file_path, encoding='utf-8')
235     print("UTF-8")
236 except UnicodeDecodeError:
237     df = pd.read_csv(csv_file_path, encoding='latin1')
238     print("Latin 1")

240 df['date_of_sale'] = pd.to_datetime(
241     df['date_of_sale'], format=r'%Y-%m-%d', errors='coerce')

243 df = df[df['date_of_sale'].dt.year == 2022].reset_index(drop=True)
244 five_number_summary = df.describe(include='all').transpose()[
245     ['min', '25%', '50%', '75%', 'max']
246     ]
247 print(five_number_summary)
248
```

```
249   ## Graphs ############################################################
250   ##
251   ## Distribution of the closing prices in 2022
252   plt.figure(figsize=(8, 6))
253   bin_edges = [1000000, 2000000, 3000000, 4000000, 5000000, 6000000,
254                7000000, 8000000, 9000000, 10000000, 11000000]
255   plt.hist(df['slutpris'], bins=bin_edges,
256            edgecolor='black', color='#00664b', alpha=0.7)
257   plt.grid(axis='y', linestyle='-', alpha=0.3)
258   plt.xlabel('Closing␣prices␣\n(in␣tens␣of␣milions␣of␣SEK)')
259   plt.ylabel('Frequency\n\n')
260   plt.savefig('histogram_closing_prices.pdf')
261   # plt.show()
262
263   ## Scatter plot of the boarea and the closing prices
264   plt.figure(figsize=(8, 6))
265   plt.scatter(
266       df['boarea'], df['slutpris'], color='#00664b', alpha=0.6, s=50)
267   plt.grid(True, linestyle='-', alpha=0.3)
268   plt.ylabel('Closing␣prices␣\n(in␣tens␣of␣milions␣of␣SEK)\n')
269   plt.xlabel('\nBoarea␣(in␣sqr␣meters)')
270   plt.savefig('scatter_plot_closing_prices.pdf')
271   # plt.show()
272
273   ## Scatter plot of the boarea and the closing prices
274   ## differenciated by number of rooms
275   plt.figure(figsize=(10, 6))
276   plt.scatter(df['boarea'], df['slutpris'], c=df['rum'], alpha=0.7, s=50)
277   plt.colorbar(label='Number␣of␣rooms')
278   plt.grid(True, linestyle='-', alpha=0.3)
279   plt.ylabel('Closing␣prices␣\n(in␣tens␣of␣milions␣of␣SEK)\n')
280   plt.xlabel('\nBoarea␣(in␣sqr␣meters)')
281   plt.savefig('scatter_plot_closing_prices_with_number_of_rooms.pdf')
282   # plt.show()
```