

# DAT565/DIT407 Assignment 2

Saif Sayed  
gussayedfa@student.gu.se

Gona Ibrahim Abdulrahman  
gusibrigo@student.gu.se

2024-04-10

This is a report for assignment 2 for the course *Introduction to Data Science & AI* from Chalmers and Gothenburg University.

## Problem 1: Scraping house prices

In this task, the goal was to scrape house prices from a tarball file containing HTML files and extract relevant information from each advertisement of a closing price. The data was scraped from Hemnet in October 2023. The extracted information included the date of sale, address, location, bo-area (habitable area), bi-area (non-habitable area), total-area (combined area), rooms, area of the plot, and the closing price.

To solve the problem, the following steps were followed:

**Importing the necessary libraries:** The code began by importing the required libraries, including `os`, `tarfile`, `pandas`, `datetime`, `re`, and `BeautifulSoup`.

**Extracting the HTML files:** The tarball file was opened, and the HTML files were extracted to a designated directory. The list of HTML files was obtained using the `os` module.

**Parsing HTML files:** The code then looped through each HTML file, opened it, and read its contents. `Beautiful Soup` was used to parse the HTML content.

**Extracting information from advertisements:** Within each HTML file, the code used `Beautiful Soup` to find all advertisements. For each advertisement, the relevant information was extracted, including the date of sale, address, location, bo-area, bi-area, total-area, rooms, area of the plot, and the closing price. Discrepancies in how the information was presented were taken into account.

**Storing the extracted data:** The extracted information was stored in a list as individual entries. The information for each advertisement was appended to the list.

**Converting data to a DataFrame and saving as CSV:** The data list was converted to a `pandas DataFrame`. Finally, the `DataFrame` was saved as a CSV file named `"house_prices.csv"` with the specified columns.

The resulting CSV file contains the extracted information for each entry in the dataset, including the address, location, bo-area, bi-area, total-area, rooms, area of the plot, date of sale, and the sale price.

By following these steps, the problem of scraping house prices and extracting relevant information from the HTML files was successfully solved. The resulting CSV file can be further analyzed and used for various purposes like data exploration, visualization, or modeling.

Our source code can be found in Appendix A of this document.

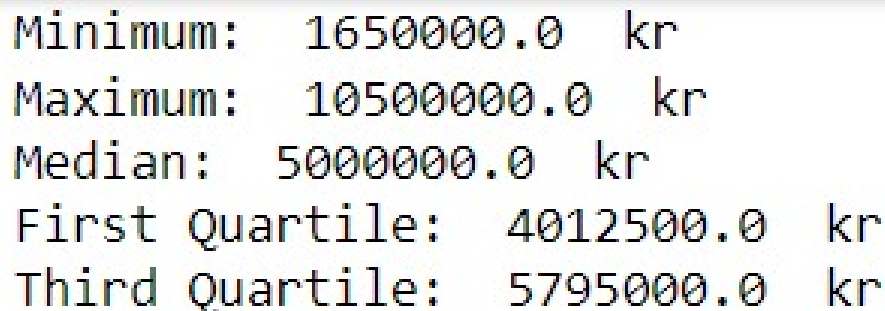
## Problem 2: Analyzing 2022 house sales

In this task, we analyzed the data of houses sold in 2022 and produced several plots to gain insights into the closing prices and their relationship with other variables. The steps taken to analyze the data and create the plots are described in the following sub-sections.

### 1 Five Number Summary

**Filtering the data:** We created a new DataFrame called *houses\_sold\_2022* by selecting houses that were sold in 2022 from the original dataset.

**Computing the five-number summary:** We computed the five-number summary (minimum, maximum, median, first quartile, and third quartile) of the closing prices for houses sold in 2022 using the `describe()` function on the 'Sale Price (kr)' column (see Figure 1).



```
Minimum: 1650000.0 kr
Maximum: 10500000.0 kr
Median: 5000000.0 kr
First Quartile: 4012500.0 kr
Third Quartile: 5795000.0 kr
```

Figure 1: Five number summary for houses sold in 2022

## 2 Histogram

**Histogram of closing prices:** We constructed a histogram to visualize the distribution of closing prices. The number of bins was determined using the square root choice method, where the number of bins is the square root of the total number of observations (see Figure 2).

## 3 Scatter Plot

**Scatter plot of closing price vs. bo-area:** We created a scatter plot to examine the relationship between the closing price and the boarea (habitable area) of the houses sold in 2022 (see Figure 3).

## 4 Scatter Plot (grouped by Number of Rooms)

**Scatter plot with color-coded rooms:** We repeated the scatter plot but colorized the observations based on the number of rooms in each house. Each point in the scatter plot was assigned a color based on the number of rooms, using a colormap (see Figure 4).

## 5 Discussion

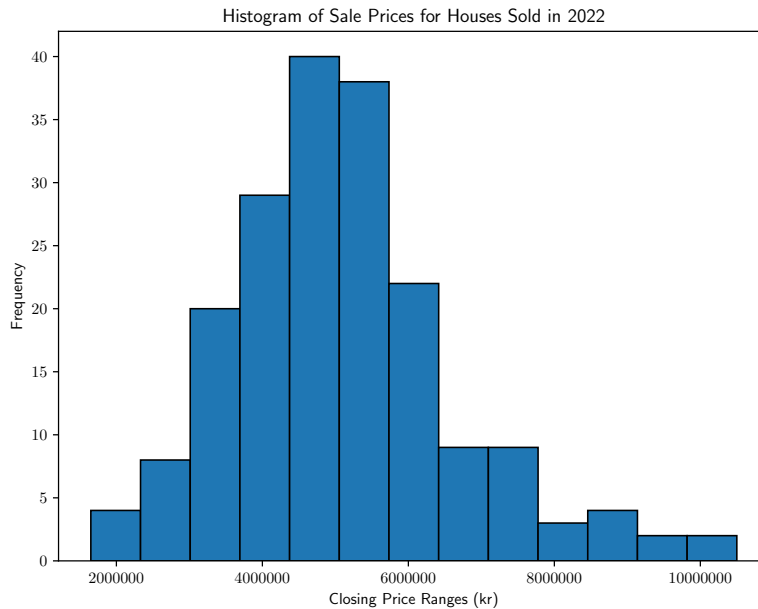


Figure 2: Frequency of Sale Prices for Houses Sold in 2022

**Figure 2:** The histogram in the image represents the distribution of sale prices for houses sold in 2022.

The histogram shows the distribution of house sale prices in 2022. Most houses were sold in the price range of 400 000 kr to 600 000 kr, with a peak frequency around 35. As prices increased beyond 600 000 kr, the frequency of sales declined. Below 400 000 kr, there was a sharp drop in frequency. This suggests a price sensitivity among buyers. However, further analysis would be needed to understand other factors influencing house prices.

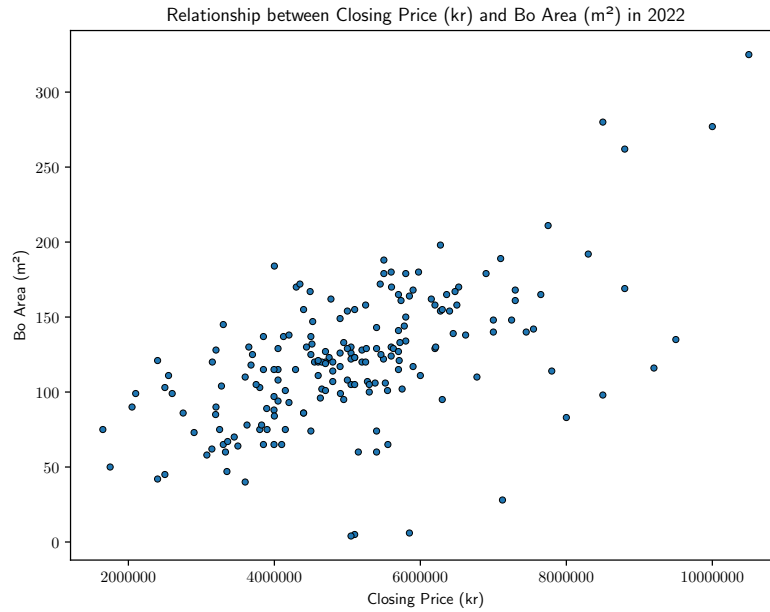


Figure 3: Relationship between Closing Price (kr) and Bo Area (m²) in 2022

**Figure 3 and Figure 4:** The scatter plots illustrate the relationship between Closing Price (kr) and Bo Area (m²) for houses sold in 2022:

The scatter plot shows that larger Bo Area (m²) tends to be associated with higher Closing Prices (kr) for houses sold in 2022. Buyers are willing to pay more for spacious properties. However, there is price variability even for similar area sizes. Other factors like location and amenities also play a role. Overall, the positive correlation emphasizes that size matters in house prices.

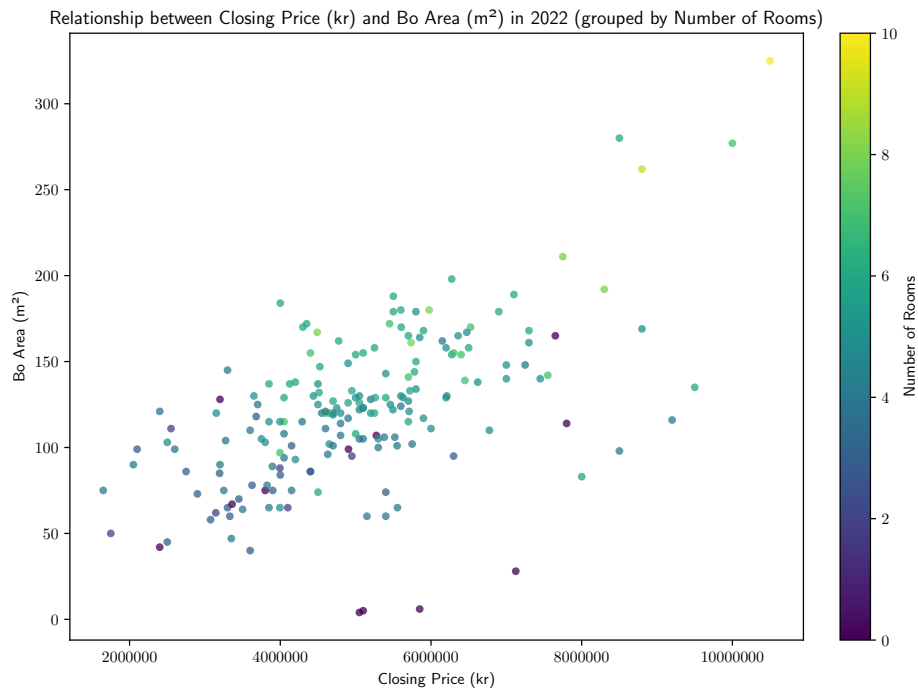


Figure 4: Relationship between Closing Price (kr) and Bo Area (m<sup>2</sup>) in 2022 (grouped by Number of Rooms)

## A Python code

This is the code we used to extract relevant information, perform data cleaning and plot our figures.

```

1  #import modules
2  import re
3  import csv
4  from bs4 import BeautifulSoup
5  import os
6  import tarfile
7  import pandas as pd
8  import locale
9  from datetime import datetime
10 import numpy as np
11 import matplotlib.pyplot as plt
12
13 #set local datetime to swedish to interpret swedish
    months
14 locale.setlocale(locale.LC_TIME, "sv_SE")
15 '''
16 Problem 1
17 '''

```

```

18 # Get current working directory of source file
19 root_dir = os.getcwd()
20
21 # Path to the input tar.gz file
22 tar_file_path = os.path.join(root_dir, '
    kungolv_slutpriser.tar.gz')
23 # Path where the tarbile file is extracted to
24 extract_dir = root_dir
25 # Path to all the html files
26 html_dir = os.path.join(root_dir, 'kungolv_slutpriser'
    )
27
28 # open tarbile for reading
29 with tarfile.open(tar_file_path, 'r', encoding='utf-8'
    ) as tar:
30     tar.extractall(extract_dir) # extract all files
31     file_list = os.listdir(html_dir) # list all html
        files in directory
32     data = [] # create placeholder to store data
33     # loop through all the html files
34     for filename in file_list :
35         # Check if the extracted file is an HTML file
36         if filename.endswith('.html'):
37             file_path = os.path.join(html_dir,
                filename)
38             #open a specific html file in the above
                file_path
39             with open(file_path, 'r', encoding='utf-8'
                ) as file:
40                 # Read the contents of the HTML file
41                 html_content = file.read()
42                 # Parse html file using beautifulsoup
43                 soup = BeautifulSoup(html_content, 'html.
                    parser')
44                 # find all ads by html id/class
45                 ads = soup.find_all('li', class_='sold-
                    results__normal-hit')
46                 # loop through all the ads
47                 for ad in ads:
48                     # Initialize variables for each piece of
                        information
49                     date_of_sale = ''
50                     address = ''
51                     location = ''
52                     bo_area = ''
53                     bi_area = ''
54                     total_area = ''
55                     area_plot = ''
56                     sale_price = ''
57                     rooms = ''

```

```

58
59     # Extract the date of sale, find
        element by id/class
60     date_elem = ad.find('span', class_='
        hcl-label_hcl-label--state_hcl-
        label--sold-at')
61     if date_elem:
62         date_of_sale = date_elem.text.
            strip().replace("S 1d", "")
63         date_of_sale = datetime.strptime(
            date_of_sale, '%d_%B_%Y')
64         date_of_sale = pd.to_datetime(
            date_of_sale, format = '%Y-%m-%
            d') # converts to datetime
            format
65
66     # Extract the address and location,
        find element by id/class
67     address_elem = ad.find('h2', class_='
        sold-property-listing__heading_qa-
        selling-price-title_hcl-card_title
        ')
68     if address_elem:
69         address = address_elem.text.strip
            ()
70
71     parent_location = ad.find('span',
        class_='property-icon_property-icon
        --result').parent
72     if parent_location:
73         parent_location.span.decompose()
74         location_strip = parent_location.
            text.strip().replace("\n", "")
            # remove unnecessary string
            literals
75         location_split = location_strip.
            split(',')
76         locations = list(map(lambda x: x.
            strip(), location_split)) #
            remove leading white spaces
77         location = ", ".join(locations)
78
79     # Extract the bi area, find element by
        id/class
80     bi_elem = ad.find('span', class_='
        listing-card__attribute--normal-
        weight')
81     if bi_elem:
82         bi_text = bi_elem.text.strip()

```

```

83         match = re.search(r'(\d+)',
84                             bi_text) # collect only the
85                                         digits using regex in the
86                                         extracted data
87         if match:
88             bi_area = match.group(1)
89         else:
90             bi_area= ""
91
92         # Extract the bo area, bi area, rooms
93         # and total area; find element by
94         # id/class
95     area_elem = ad.find('div', class_='
96                         sold-property-listing__subheading_
97                         sold-property-listing__area')
98     if area_elem:
99         area_elem.span.decompose if
100             area_elem.span else area_elem
101         area_text = area_elem.text.strip()
102         match = re.findall(r'\d+',
103                             area_text) # collect only the
104                                         digits using regex in the
105                                         extracted data
106         if len(match) > 2:
107             # if 3 pieces of info.
108                 available for area and room
109                 , then store all three data
110                 and compute total area.
111             bo_area = int(match[0])
112             bi_area = int(match[1])
113             rooms = int(match[2])
114             total_area = int(match[0]) +
115                             int(match[1])
116         elif len(match) == 2 and bi_area
117             != "":
118             # if 2 pieces of info.
119                 available which includes bi
120                 area, then set room to 0,
121                 calculate total area
122             bo_area = int(match[0])
123             bi_area = int(match[1])
124             rooms = 0
125             total_area = int(match[0]) +
126                             int(match[1])
127         elif len(match) == 1:
128             # if only 1 piece of info.
129                 available, then set room
130                 and bi area to 0, calculate
131                 total area
132             bo_area = int(match[0])

```



```

110         bi_area = 0
111         rooms = 0
112         total_area = int(match[0])
113     elif len(match) == 2 and bi_area
114         == "":
115         # if 2 pieces of info.
116         # available which doesn't
117         # include bi area, then set
118         # bi area to 0
119         bo_area = int(match[0])
120         bi_area = 0
121         rooms = int(match[1])
122         total_area = int(match[0])
123
124     # Extract the area plot, find element
125     # by id/class
126     plot_elem = ad.find('div', class_='
127     sold-property-listing__land-area')
128     if plot_elem:
129         plot_text = plot_elem.text.strip()
130         match = re.findall(r'\d+',
131         plot_text) # collect only the
132         # digits using regex in the
133         # extracted data
134         if match:
135             area_plot = str(match[0]) + "
136             m " + "tmt"
137
138     # Extract the sale price, find element
139     # by id/class
140     price_elem = ad.find('span', class_='
141     hcl-text_hcl-text--medium')
142     if price_elem:
143         price = price_elem.text.strip()
144         match = re.findall(r'\d+', price)
145         # collect only the digits using
146         # regex in the extracted data
147         if match:
148             sale_price = ""
149             for value in match:
150                 sale_price += str(value)
151             sale_price = int(sale_price)
152     # Append the extracted data to the
153     # list
154     data.append([address, location,
155     bo_area, bi_area, total_area, rooms
156     , area_plot, date_of_sale,
157     sale_price])
158
159 #store the data in a panda dataframe

```

```

141 df = pd.DataFrame(data, columns=['Address', 'Location',
    , 'Bo□Area□(m )', 'Bi□Area□(m )', 'Total□Area□(
    m )', 'Rooms', 'Area□Plot', 'Date□of□Sale', 'Sale□
    Price□(kr)'])
142 df.to_csv('house_prices.csv', index=False, encoding='
    utf-8') # convert to csv
143
144 '''
145 Problem 2
146 '''
147 houses_sold_2022 = df[df['Date□of□Sale'].dt.year ==
    2022] #create and store dataframe only for the year
    2022
148 sale_price_stats = houses_sold_2022['Sale□Price□(kr)'
    ].describe()
149 #collect statistical data about sales price on the
    dataframe
150 minimum = sale_price_stats['min']
151 maximum = sale_price_stats['max']
152 median = sale_price_stats['50%']
153 first_quartile = sale_price_stats['25%']
154 third_quartile = sale_price_stats['75%']
155
156 print("Minimum:□", minimum, "□kr")
157 print("Maximum:□", maximum, "□kr")
158 print("Median:□", median, "□kr")
159 print("First□Quartile:□", first_quartile, "□kr")
160 print("Third□Quartile:□", third_quartile, "□kr")
161
162 '''Histogram of Sale Prices for Houses Sold in 2022'''
163
164 sale_prices = np.asarray(houses_sold_2022['Sale□Price□
    (kr)'].dropna(), float)
165 fig1, ax1 = plt.subplots(figsize=(8, 6))
166 # the number of bins was determined using the square
    root choice method
167 num_bins = int(len(houses_sold_2022['Sale□Price□(kr)'
    ]) ** 0.5)
168 ax1.hist(sale_prices, bins=num_bins, edgecolor='black',
    , linewidth=1, label='Closing□prices')
169 ax1.set_xlabel('Closing□Price□Ranges□(kr)')
170 ax1.set_ylabel('Frequency')
171 ax1.set_title('Histogram□of□Sale□Prices□for□Houses□
    Sold□in□2022')
172 ax1.ticklabel_format(useOffset=1, style='plain', axis=
    'x') # change the tick format on the x-axis to be
    more comprehensive
173 fig1.savefig('histogram.pdf')
174

```

```

175  """Relationship between Closing Price (kr) and Bo Area
      (m ) in 2022"""
176
177  fig2, ax2 = plt.subplots(figsize=(8, 6))
178  bo_area = pd.to_numeric(houses_sold_2022['Bo_Area_(m
      )']).dropna(), errors = 'coerce').values
179  ax2.scatter(sale_prices, bo_area, s= 15, edgecolor='
      black', linewidth=0.5)
180  ax2.set_title("Relationship between Closing Price (kr)
      and Bo Area (m ) in 2022")
181  ax2.set_xlabel('Closing Price (kr)')
182  ax2.set_ylabel('Bo Area (m )')
183  ax2.ticklabel_format(useOffset=1, style='plain', axis=
      'x') # change the tick format on the x-axis to be
      more comprehensive
184  fig2.savefig('scatter_plot.pdf')
185
186  '''Relationship between Closing Price (kr) and Bo Area
      (m ) in 2022 (grouped by Number of Rooms)'''
187
188  fig3, ax3 = plt.subplots(figsize=(8, 6),
      constrained_layout=True)
189  rooms = pd.to_numeric(houses_sold_2022['Rooms']).dropna
      (), errors = 'coerce').values
190  ax3.scatter(sale_prices, bo_area , c= rooms, cmap='
      viridis', s=15, alpha= 0.75)
191  ax3.set_title("Relationship between Closing Price (kr)
      and Bo Area (m ) in 2022 (grouped by Number of
      Rooms)")
192  ax3.set_xlabel('Closing Price (kr)')
193  ax3.set_ylabel('Bo Area (m )')
194  sm = plt.cm.ScalarMappable(cmap='viridis') # mapping
      scalar values to the colormap based on the values
      stored in 'rooms'
195  sm.set_array(rooms)
196  fig3.colorbar(sm, label='Number of Rooms', ax=ax3) #
      create colorbar
197  ax3.ticklabel_format(useOffset=1, style='plain', axis=
      'x') # change the tick format on the x-axis to be
      more comprehensive
198  fig3.savefig('scatter_plot_room.pdf')

```