# DAT565/DIT407 Assignment 2

Louis PRUNIER
prunier.louis@icloud.com

Marco SPEZIALE
speziale@student.chalmers.se

2024-04-10

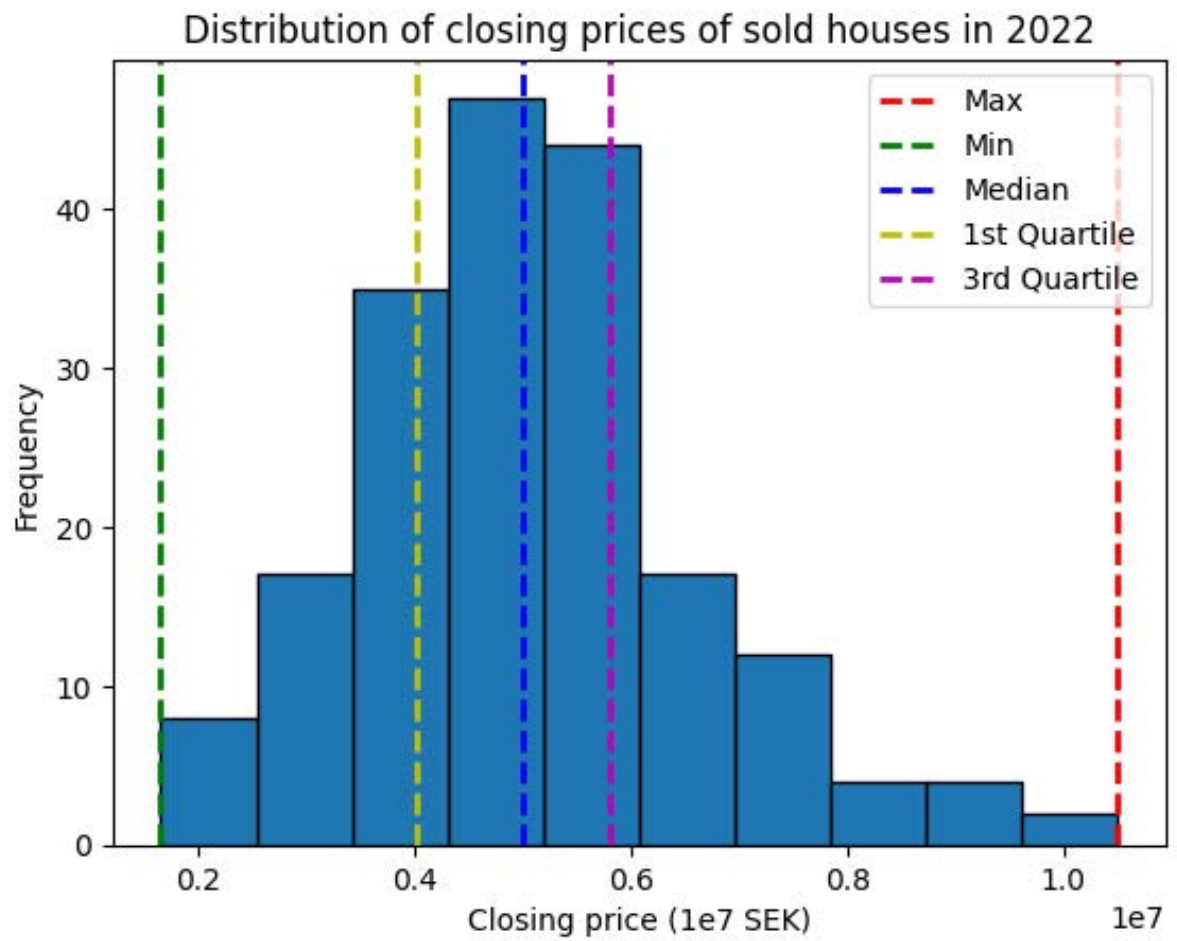# 1 Figure 1



Figure 1: Distribution of Closing Prices in 2022

## 2 Figure 2

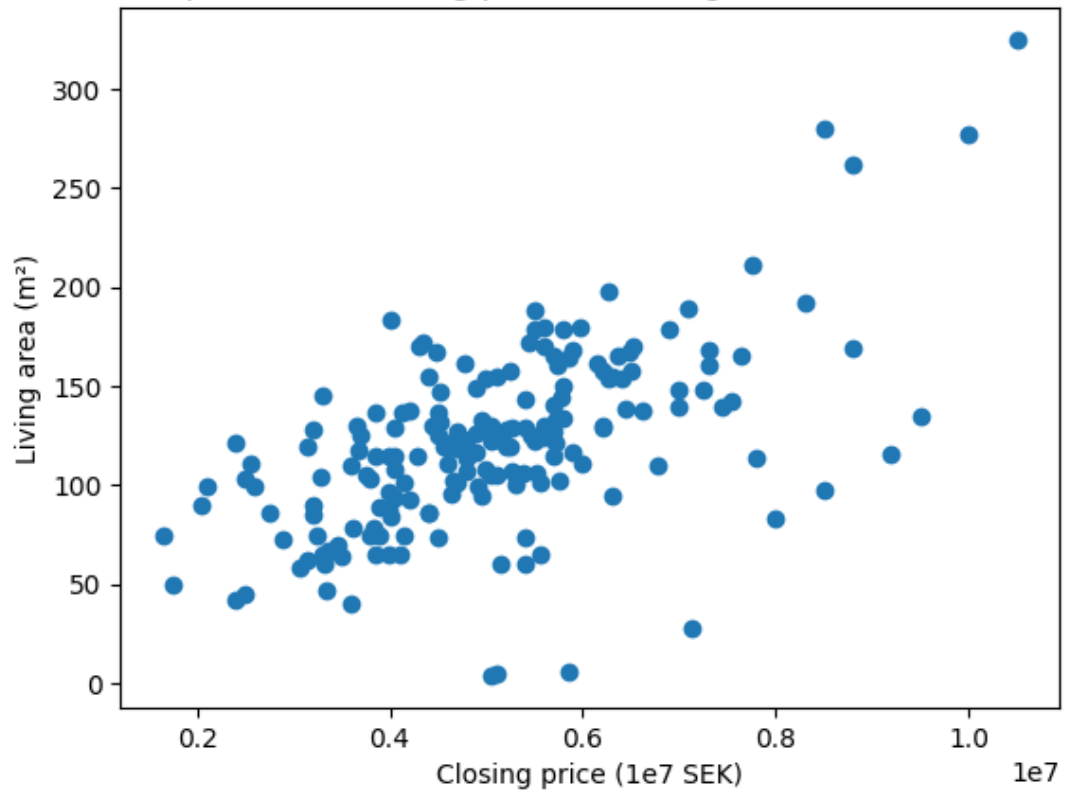Relationship between closing price and living area of sold houses in 2022

Figure 2: Relationship between house Boarea and closing price (2022)
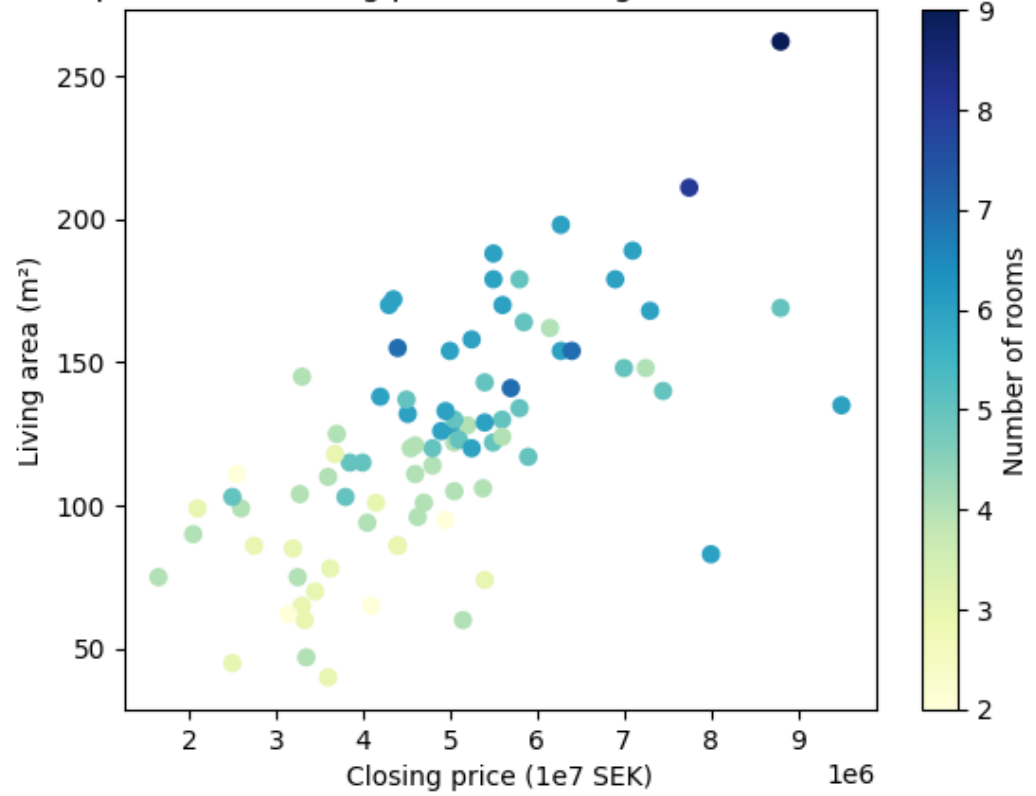
# 3 Figure 3



Figure 3: Relation between house boarea and closing price by number of rooms

# 4 Discussion

## 4.1 Figure 1

Figure 1 is a bar chart illustrating the distribution of house closing prices in 2022 in Kungälv. On this histogram, 5 key statistical concepts are added: the minimum value, the maximum value, the median, the first quartile and the third quartile. On this histogram, we can see that house sale prices in 2022 in Kungälv range from SEK 1,650,000 to a maximum of SEK 10,500,000. The median price is SEK 5,000,000. 25% of closing prices are below SEK 4,000,000 and conversely only 25% of closing prices are above SEK 5,795,000. We chose 10 as the number of bins to show the distribution as the prices range from approximately 1,000,000 to 10,000,000 SEK, and each bin can therefore represent (approximately) a step in millions of SEK. 10 is also enough bins to make each line corresponding to the five-number summary fit into one unique bin each, clearly showing the frequency associated with the given statistical concept.

## 4.2 Figure 2

Figure 2 is a scatter-plot showing the relationship between living areas (boarea) and the closing price of the house. We can see from the trend line that prices increase in a linear fashion as living area also increases, reflecting a fairly logical correlation between living area and house price. Figure 2 also shows that the high concentration of points (houses) is mainly between 100 and 150 m2 at a price of SEK 5,000,000 (median).

## 4.3 Figure 3

Finally, Figure 3 is identical to number 2, except that the number of bedrooms in the house is also indicated (colored dots). On this figure, we can see that the smallest houses (between 50 and 100 m2) have the fewest bedrooms (between 0 and 2 bedrooms), and conversely the largest houses have the most bedrooms (over 250 m2 and 8 to 10 bedrooms). The vast majority of houses (between 100 and 150 m2) have around 4 bedrooms (between 2 and 6 bedrooms to be broader).

## 4.4 Conclusion

So, the conclusions we can draw are that the closing price of a house in Kungälv is strongly correlated to the variable 'living area' (BoArea), even if there are disparities with some outliers. But also, and quite logically, the largest houses have the most bedrooms and vice versa for the smallest, which conversely doesn't really influence the closing price but depends more on the size of the house's living area.

# 5  Appendix

```
1  #%%
2  import os
3  from bs4 import BeautifulSoup
4  import pandas as pd
5
6
7  ## PART 1 ##
8
9  # methods to clean up the extracted information into
       useful formats
10
11 def cleanup_date(date_info):
12     split_date = date_info.split()
13     pure_date = split_date[1:]
14     return '␣'.join(pure_date)
15
16 def cleanup_address(address_info):
17     return address_info.strip()
18
19 def cleanup_location(location_info):
20     location_info_split = location_info.split()
21     return '␣'.join(location_info_split)
22
23 def cleanup_bo_area(bo_area_info):
24     bo_area = bo_area_info.strip()
25     bo_area = bo_area.replace(',','.')
26     if bo_area.isnumeric():
27         return float(bo_area)
28     else:
29         return None
30
31 def cleanup_bi_area(bi_area_info):
32     bi_area_split = bi_area_info.split()
33     bi_area_pure = bi_area_split[1:2]
34     bi_area = ''.join(bi_area_pure)
35     bi_area = bi_area.replace(',','.')
36     if bi_area.isnumeric():
37         return float(bi_area)
38     else:
39         return None
40
41 def cleanup_rooms(nr_rooms_info):
42     nr_rooms_split = nr_rooms_info.split()
```

```python
43        nr_rooms_pure = nr_rooms_split[:1]
44        nr_rooms = ''.join(nr_rooms_pure)
45        if nr_rooms.isnumeric():
46            nr_rooms = int(nr_rooms)
47        else:
48            nr_rooms = None
49
50   def cleanup_area_room(area_and_room_info):
51        area_room = []
52        area_and_room_info_split = area_and_room_info.
             split()
53        area_pure = area_and_room_info_split[:1]
54        room_pure = area_and_room_info_split[2:3]
55
56        area = ''.join(area_pure)
57        area = area.replace(',','.')
58        if area.isnumeric():
59            area = float(area)
60        else:
61            area = None
62        area_room.append(area)
63
64        nr_rooms = ''.join(room_pure)
65        if nr_rooms.isnumeric():
66            nr_rooms = int(nr_rooms)
67        else:
68            nr_rooms = None
69        area_room.append(nr_rooms)
70        return area_room
71
72   def cleanup_plotarea(plot_area_info):
73        plot_area = ''
74        for char in plot_area_info:
75            if char.isnumeric():
76                plot_area += char
77        plot_area = plot_area[:-1]
78        plot_area.replace(',','.')
79        return float(plot_area)
80
81   def cleanup_price(price_info):
82        price = ''
83        for char in price_info:
84            if char.isnumeric():
85                price += char
86        return float(price)
87
```

```
88
89  # directory with html files
90  directory = 'kungalv_slutpriser'
91
92  data = []
93
94  # for each file in directory, extract the sought-out
        information
95  for filepath in os.listdir(directory):
96      with open(os.path.join(directory, filepath),
            encoding='utf-8') as fp:
97          soup = BeautifulSoup(fp, 'html.parser')
98
99      result = soup.find_all('li', {'class': 'sold-
            results__normal-hit'})
100
101     # find the relevant pieces of information
102     for element in result:
103         date_info = element.find('span', {'class': '
                hcl-label␣hcl-label--state␣hcl-label--sold-
                at'}).text
104         date = cleanup_date(date_info)
105
106         address_info = element.find('h2', {'class': '
                sold-property-listing__heading␣qa-selling-
                price-title␣hcl-card__title'}).text
107         address = cleanup_address(address_info)
108
109         location_info = element.find('div', {'class':
                'sold-property-listing__location'}).
                contents[3].contents[2]
110         location = cleanup_location(location_info)
111
112         if len(element.find('div', {'class': 'sold-
                property-listing__subheading␣sold-property-
                listing__area'}).contents) > 1:
113             bo_area_info = element.find('div', {'class
                    ': 'sold-property-listing__subheading␣
                    sold-property-listing__area'}).contents
                    [0]
114             bo_area = cleanup_bo_area(bo_area_info)
115
116             bi_area_info = element.find('div', {'class
                    ': 'sold-property-listing__subheading␣
                    sold-property-listing__area'}).contents
                    [1].text
```

```python
117                    bi_area = cleanup_bi_area(bi_area_info)
118
119                    if (bo_area and bi_area):
120                        area = bo_area + bi_area
121                    elif (bo_area):
122                        area = bo_area
123
124                    nr_rooms_info = element.find('div', {'
                           class': 'sold-property-
                           listing__subheading␣sold-property-
                           listing__area'}).contents[2]
125                    nr_rooms = cleanup_rooms(nr_rooms_info)
126                else:
127                    area_and_room_info = element.find('div', {
                           'class': 'sold-property-
                           listing__subheading␣sold-property-
                           listing__area'}).contents[0]
128                    area_and_room = cleanup_area_room(
                           area_and_room_info)
129
130                    area = area_and_room[0]
131                    bi_area = None
132                    bo_area = area
133
134                    nr_rooms = area_and_room[1]
135
136            if element.find('div', {'class': 'sold-
                   property-listing__land-area'}):
137                    plot_area_info = element.find('div', {'
                           class': 'sold-property-listing__land-
                           area'}).text
138                    plot_area = cleanup_plotarea(
                           plot_area_info)
139            else:
140                    plot_area = None
141
142            price_info = element.find('span', {'class': '
                   hcl-text␣hcl-text--medium'}).text
143            price = cleanup_price(price_info)
144
145            data.append([date, address, location, bo_area,
                   bi_area, area, nr_rooms, plot_area, price
                   ])
146
147 # turn data into dataframe
```

```python
148  df = pd.DataFrame(data, columns=['Date', 'Address', '
         Location', 'Bo-area', 'Bi-area', 'Total␣area', '
         Rooms', 'Plot', 'Price'])
149
150  # turn dataframe into csv file
151  csv = df.to_csv('houseprices.csv')
152
153  #%%
154  ## PART 2 ##
155  import pandas as pd
156  import matplotlib.pyplot as plt
157
158  df = pd.read_csv('houseprices.csv')
159  df = df.drop('Unnamed:␣0', axis=1)
160
161  sold_2022 = df[df['Date'].str.contains('2022')]
162
163  closing_prices_2022 = sold_2022['Price']
164  closing_prices_2022
165
166  min_2022 = closing_prices_2022.min()
167  max_2022 = closing_prices_2022.max()
168  median_2022 = closing_prices_2022.median()
169  first_quartile_2022 = closing_prices_2022.quantile
         (0.25)
170  third_quartile_2022 = closing_prices_2022.quantile
         (0.75)
171
172  # histogram
173  plt.hist(closing_prices_2022, bins=10, edgecolor='
         black')
174
175  plt.axvline(max_2022, color='r', linestyle='dashed',
         linewidth=2, label='Max')
176  plt.axvline(min_2022, color='g', linestyle='dashed',
         linewidth=2, label='Min')
177  plt.axvline(median_2022, color='b', linestyle='dashed'
         , linewidth=2, label='Median')
178  plt.axvline(first_quartile_2022, color='y', linestyle=
         'dashed', linewidth=2, label='1st␣Quartile')
179  plt.axvline(third_quartile_2022, color='m', linestyle=
         'dashed', linewidth=2, label='3rd␣Quartile')
180
181  plt.xlabel('Closing␣price␣(1e7␣SEK)')
182  plt.ylabel('Frequency')
```

```
183  plt.title('Distribution␣of␣closing␣prices␣of␣sold␣
         houses␣in␣2022')
184
185  plt.legend()
186  plt.show()
187
188  # scatter plot
189  bo_areas_2022 = sold_2022['Bo-area']
190  bo_areas_2022
191
192  plt.scatter(closing_prices_2022, bo_areas_2022)
193  plt.ylabel('Living␣area␣(m^2)')
194  plt.xlabel('Closing␣price␣(1e7␣SEK)')
195  plt.title('Relationship␣between␣closing␣price␣and␣
         living␣area␣of␣sold␣houses␣in␣2022')
196
197  plt.show()
198
199  # colored version
200  nr_rooms_2022 = sold_2022['Rooms']
201
202  scatter = plt.scatter(closing_prices_2022,
         bo_areas_2022, c=nr_rooms_2022, cmap='YlGnBu')
203  plt.ylabel('Living␣area␣(m^2)')
204  plt.xlabel('Closing␣price␣(1e7␣SEK)')
205  plt.title('Relationship␣between␣closing␣price␣and␣
         living␣area␣of␣sold␣houses␣in␣2022')
206
207  plt.colorbar(scatter, label='Number␣of␣rooms')
208  plt.show()
```