# DAT565/ Assignment 2 / Group 96

David Duong
thuan@chalmers.se

Savinjith Walisadeera
savwal@chalmers.se

September 16, 2024

## Problem 2

### i)

Most expensive house in Kungälv municipality during the year 2022 was 10 500 000 and the cheapest 1 650 000. The median price was 5 000 000. First quartile was 4 012 500 and third quartile was 5 795 000.

### ii)

We chose the number of bins from 'Square Root Rule' which states number of bins $= \sqrt{n}$. In our case we had 200 data points giving us $\sqrt{190} \approx 14$.
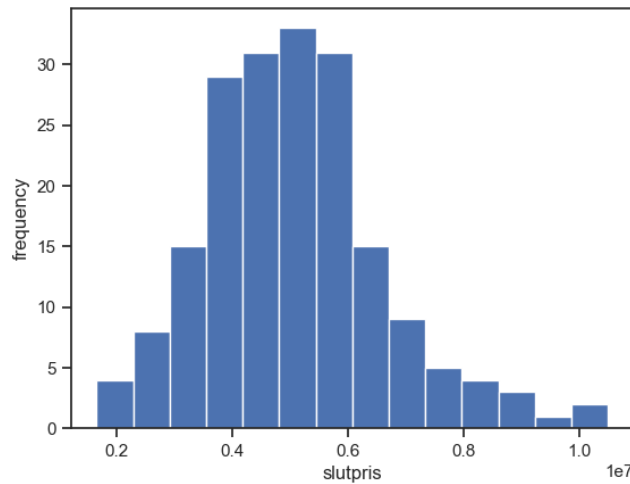


Figure 1: Histogram of closing prices and frequency in Kungälv municipality.
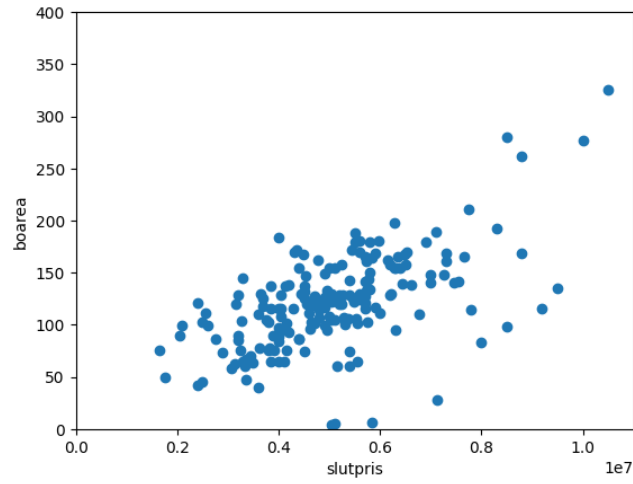
**iii)**



Figure 2: Scatter plot of *boarea* to closing prices.
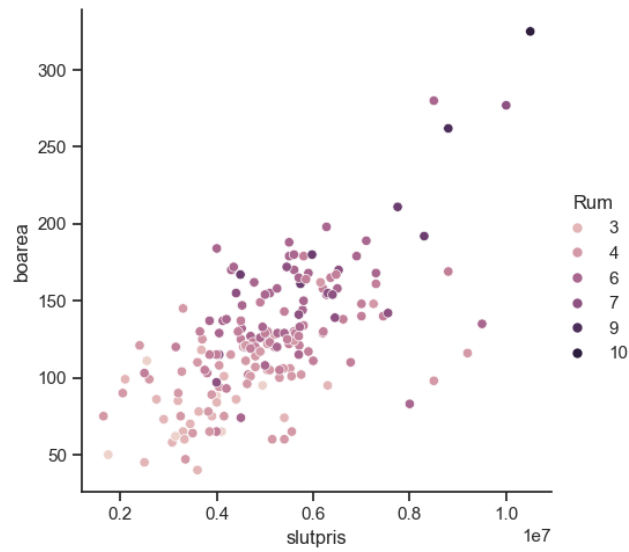
**iv)**



Figure 3: Scatter plot of *boarea* to closing prices with the number of rooms of the house defining the color of the point.

**v)**

From looking at the scatter plots in *Figure 2* and *Figure 3* we can see that *boarea* grows linearly to the closing price. As we look at growing closing prices in *Figure 3* we see that the colors of the points get darker as *boarea* rises. This makes sense as the number of rooms usually increase with *boarea*. Therefore we can conclude that the closing price is dependent on *boarea*.

# A  Code

```python
1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  from bs4 import BeautifulSoup
5  import csv
6  import re
7  import seaborn as sns
8  import matplotlib
9  import matplotlib.colors
10
11 import pandas as pd
12 data = list()
13
14 months = {'januari':'01', 'februari':'02', 'mars':'03'
       , 'april':'04','maj':'05','juni':'06', 'juli':'07',
        'augusti':'08', 'september':'09', 'oktober':'10','
       november':'11','december':'12'}
15 for j in range(1,5):
16     with open("kungalv_slutpriser\
           kungalv_slutpris_page_{}{}.html".format(j,0),
           encoding = 'utf-8') as fp:
17         soup = BeautifulSoup(fp, 'html.parser')
18     sellings = soup.find_all(attrs={'class':'sold-
           results__normal-hit'})
19     for sold in sellings:
20         sold_at = sold.find('span', class_='hcl-label␣
               hcl-label--state␣hcl-label--sold-at')
21         sold_at = sold_at.text.strip()[5:]
22         for month in months:
23             sold_at = sold_at.replace(month, months[
                   month])
24         sold_at = sold_at.split()
25
26         address = sold.find('h2', class_='sold-
               property-listing__heading␣qa-selling-price-
               title␣hcl-card__title')
27         address_text = address.text.strip()
28         location = address.parent.div
29         location.span.clear()
```

```python
30              location = list(location.stripped_strings)
31              location = location[0].split(',')
32              for i in range(len(location)):
33                  location[i] = location[i].strip()
34
35              try:
36                  tomt = sold.find('div', class_="sold-
                        property-listing__land-area").text.
                        strip()[:-8]
37              except:
38                  tomt = pd.NA
39
40              price = sold.find('span', class_='hcl-text hcl
                    -text--medium').text.strip()[9:-3]
41              price = int(re.sub('[^0-9]+', '', price))
42              rooms = sold.find('div', class_='sold-property
                    -listing__subheading sold-property-
                    listing__area')
43              rooms = rooms.text.replace(' ', '').replace('\
                    n', ' ')
44              rooms_val = re.sub('[^0-9 ]+', '', rooms).
                    split()
45
46              boarea = pd.NA
47              room_amount = pd.NA
48              biarea = pd.NA
49              if len(rooms_val) == 2:
50                  if bool(re.search('(rum)',rooms)):
51                      boarea = int(rooms_val[0])
52                      room_amount = int(rooms_val[1])
53                  else:
54                      boarea = int(rooms_val[0])
55                      biarea = int(rooms_val[1])
56              elif len(rooms_val) == 1:
57                  boarea = int(rooms_val[0])
58              elif len(rooms_val) == 3:
59                  boarea = int(rooms_val[0])
60                  biarea = int(rooms_val[1])
61                  room_amount = int(rooms_val[2])
62              elif len(rooms_val) == 0:
63                  boarea = pd.NA
64                  biarea = pd.NA
65                  room_amount = pd.NA
66
67
68              row = {
69                      'Year of sale':sold_at[2],
70                      'Month of sale':sold_at[1],
71                      'Day of sale':sold_at[0],
72                      'Address': address_text,
```

```python
73                    'Ort': location[0],
74                    'Kommun': location[1],
75                    'boarea': boarea,
76                    'biarea': biarea,
77                    'Rum': room_amount,
78                    'total_area': boarea+biarea,
79                    'plot_area': tomt,
80                    'slutpris': price
81                        }
82            data.append(row)
83
84  for i in range(0,4):
85      for j in range(1,10):
86          with open("kungalv_slutpriser\
               kungalv_slutpris_page_{}{}.html".format(i,j
               ), encoding = 'utf-8') as fp:
87               soup = BeautifulSoup(fp, 'html.parser')
88          sellings = soup.find_all(attrs={'class':'sold-
               results__normal-hit'})
89          for sold in sellings:
90              sold_at = sold.find('span', class_='hcl-
                   label hcl-label--state hcl-label--sold-
                   at')
91              sold_at = sold_at.text.strip()[5:]
92              for month in months:
93                  sold_at = sold_at.replace(month,
                       months[month])
94              sold_at = sold_at.split()
95              address = sold.find('h2', class_='sold-
                   property-listing__heading qa-selling-
                   price-title hcl-card__title')
96              address_text = address.text.strip()
97              location = address.parent.div
98              location.span.clear()
99              location = list(location.stripped_strings)
100             location = location[0].split(',')
101             for ind in range(len(location)):
102                 location[ind] = location[ind].strip()
103
104             try:
105                 tomt = sold.find('div', class_="sold-
                       property-listing__land-area").text.
                       strip()[:-8]
106             except:
107                 tomt = pd.NA
108
109             price = sold.find('span', class_='hcl-text
                    hcl-text--medium').text.strip()[9:-3]
110             price = int(re.sub('[^0-9]+', '', price))
```

```
111              rooms = sold.find('div', class_='sold-
                    property-listing__subheading␣sold-
                    property-listing__area')
112              rooms = rooms.text.replace('␣', '').
                    replace('\n', '␣')
113              rooms_val = re.sub('[^0-9␣]+', '', rooms).
                    split()
114
115              boarea = pd.NA
116              room_amount = pd.NA
117              biarea = pd.NA
118              if len(rooms_val) == 2:
119                  if bool(re.search('(rum)',rooms)):
120                      boarea = int(rooms_val[0])
121                      room_amount = int(rooms_val[1])
122                  else:
123                      boarea = int(rooms_val[0])
124                      biarea = int(rooms_val[1])
125              elif len(rooms_val) == 3:
126                  boarea = int(rooms_val[0])
127                  biarea = int(rooms_val[1])
128                  room_amount = int(rooms_val[2])
129              elif len(rooms_val) == 1:
130                  boarea = int(rooms_val[0])
131              elif len(rooms_val) == 0:
132                  boarea = pd.NA
133                  biarea = pd.NA
134                  room_amount = pd.NA
135
136
137              row = {
138                  'Year␣of␣sale':sold_at[2],
139                  'Month␣of␣sale':sold_at[1],
140                  'Day␣of␣sale':sold_at[0],
141                  'Address': address_text,
142                  'Ort': location[0],
143                  'Kommun': location[1],
144                  'boarea': boarea,
145                  'biarea': biarea,
146                  'Rum': room_amount,
147                  'total_area': boarea+biarea,
148                  'plot_area': tomt,
149                  'slutpris': price
150                      }
151              data.append(row)
152  data = pd.DataFrame(data)
153
154  data.to_csv('kungalv_prices.csv',index=None)
155
156  sellings_2022 = data[data['Year␣of␣sale'] == '2022']
```

```python
157 print(sellings_2022['slutpris'].max())
158 print(sellings_2022['slutpris'].min())
159 print(sellings_2022['slutpris'].median())
160 print(np.quantile(sellings_2022['slutpris'], 0.25))
161 print(np.quantile(sellings_2022['slutpris'], 0.75))
162 plt.hist(sellings_2022['slutpris'], int(np.ceil(np.
        sqrt(190))))
163
164 sellings_2022_boarea = sellings_2022.dropna(subset=['
        boarea'], inplace=False)
165
166 fig, ax = plt.subplots()
167 ax.scatter(sellings_2022_boarea['slutpris'],
        sellings_2022_boarea['boarea'])
168 plt.xlim(0,1.1*10**7)
169 plt.ylim(0,400)
170
171 sns.set_theme(style='ticks')
172 sns.relplot(data=sellings_2022_boarea, x='slutpris', y
        ='boarea', hue='Rum')
173 plt.show()
```