

# MATH3670 Project Proposal

February 3, 2025

Nicholas Stone  
903887158  
MATH-3670-N

## I. Dataset #1: Scatterplot

### I.a. Data and Description

The restraints for Dataset #1 are  $2000 > n > 100$ . The provided Dataset contains 193 entries, as can be assessed by running the following in **dataset1code.py**:

```
print(f"Total number of entries (rows) in the dataset: {len(df_avg_iq)}")
```

Or by checking the total number of rows in the dataset, which can be found [here](#) as a .txt file in my Github repository for this project. Please utilize the earlier link to access the full dataset. A snippet of the data can also be found through Reference 8 in the Appendix.

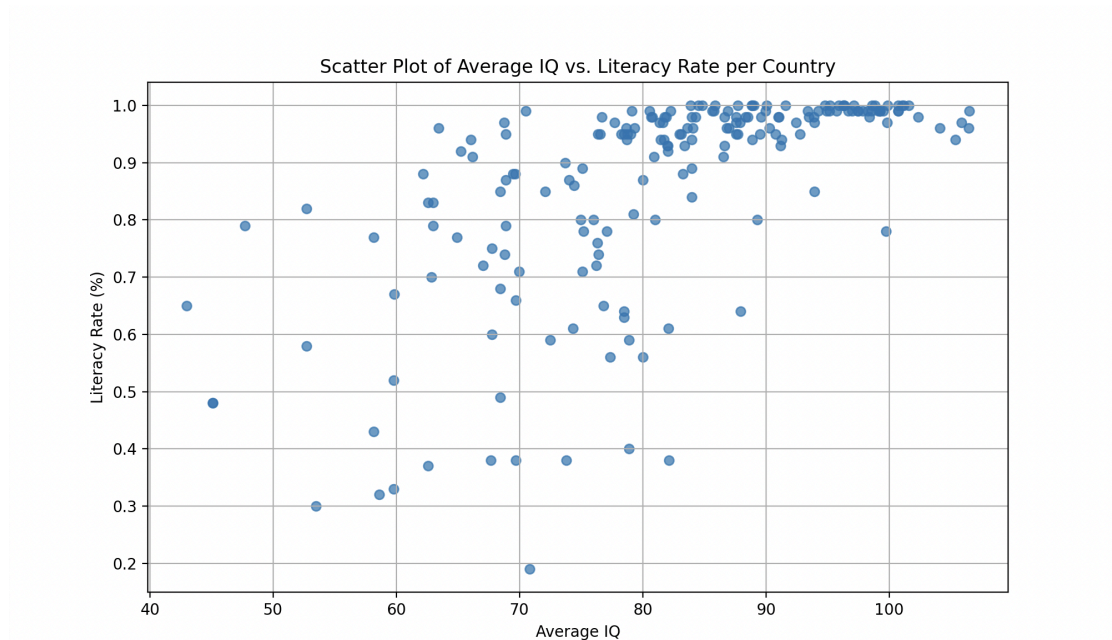
### I.b. Data Source

The source to this data can be found [here](#). The data can be downloaded locally by running **dataset1.py** in your terminal. It will provide you with a file path, which is relative to your computer and must be utilized in **dataset1code.py** or else it will not recognize the file path.

Average IQ will be  $X_i$ , where Literacy Rate will be  $Y_i$ . As you can see throughout the scatter plot in I.c, As a country's average IQ increases, the country's literacy rate tends to as well, indicating a general upward trajectory up to a literacy rate of 1.0, which flattens around an average IQ of 70, where most countries with IQs of 70 or above tend to maintain high literacy rates from  $\sim .9$  to 1.0. The data points which are below an IQ of 75 are much more sparse in comparison to those countries which have higher average IQs of 75.

### I.c. Scatter Plot of Average IQ vs. Literacy Rate per Country

*Provided on the next page.*



## II. Dataset #2: Histogram

### II.a. Data and Description

The restraints for Dataset #2 are  $5000 > n > 100$ . The provided Dataset contains 270 entries, as can be assessed by checking the total number of rows in the dataset, which can be found [here](#) as a .txt file in my Github repository for this project. Please utilize the earlier link to access the full dataset. A snippet of the data can also be found through Reference 9 in the Appendix

The source to this data can be found [here](#). The data can be downloaded locally by running **dataset2.py** in your terminal. It will provide you with a file path, which is relative to your computer and must be utilized in **dataset2code.py** or else it will not recognize the file path.

What you will observe is that the x-axis represents Age, while the y-axis represents the frequency of heart disease. What you will notice is that  $k = 10$  bins were used, effectively showcasing the frequency of heart disease amongst decades of the people who were assessed. The graph showcases a large concentration of heart disease being reported towards late 50's and early 60's, while it tapers off towards younger ages and even drops off towards older ages. The initial frequency histogram shows that the height of cases per age group is approaching 60, and the relative frequency histogram helps put that into perspective by denoting that as a little over .200.

### II.b. Sample Mean/Variance for the data set above:

#### Sample Mean

The following formula for Sample Mean was used:

$$\bar{x} = \frac{\sum x_i}{n}$$

Its code implementation can be seen [here](#). The calculated Sample Mean for this dataset was 54.43333333333333. You can run this python file by cd'ing into **project\_part1/calculations/** and running **python3 calculate\_mean.py**.

### Sample Variance

The following formula for Sample Variance was used:

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

Its code implementation can be seen [here](#). The calculated Sample Variance for this dataset was 82.97509293680297. You can run this python file by cd'ing into **project\_part1/calculations/** and running **python3 calculate\_variance.py**.

### II.c. Quartiles

The following three values separate the four quartiles of the data.

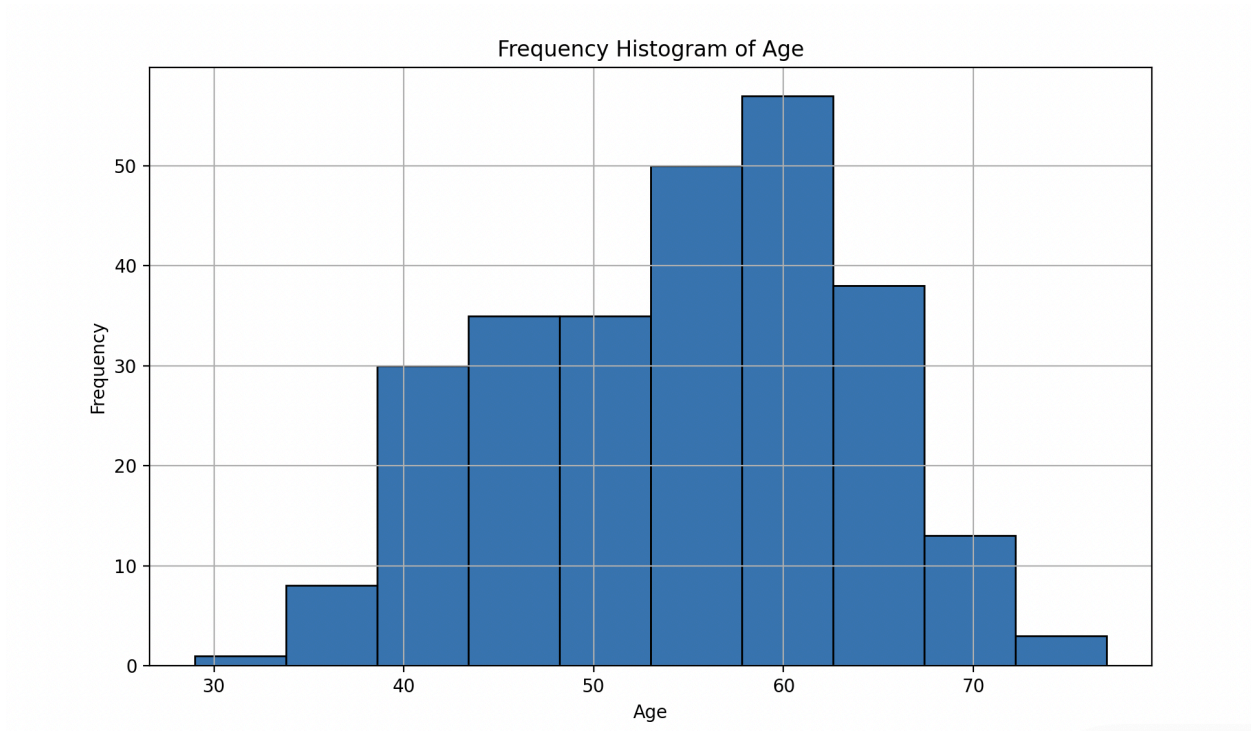
Q1: 48.0

Q2: 55.0

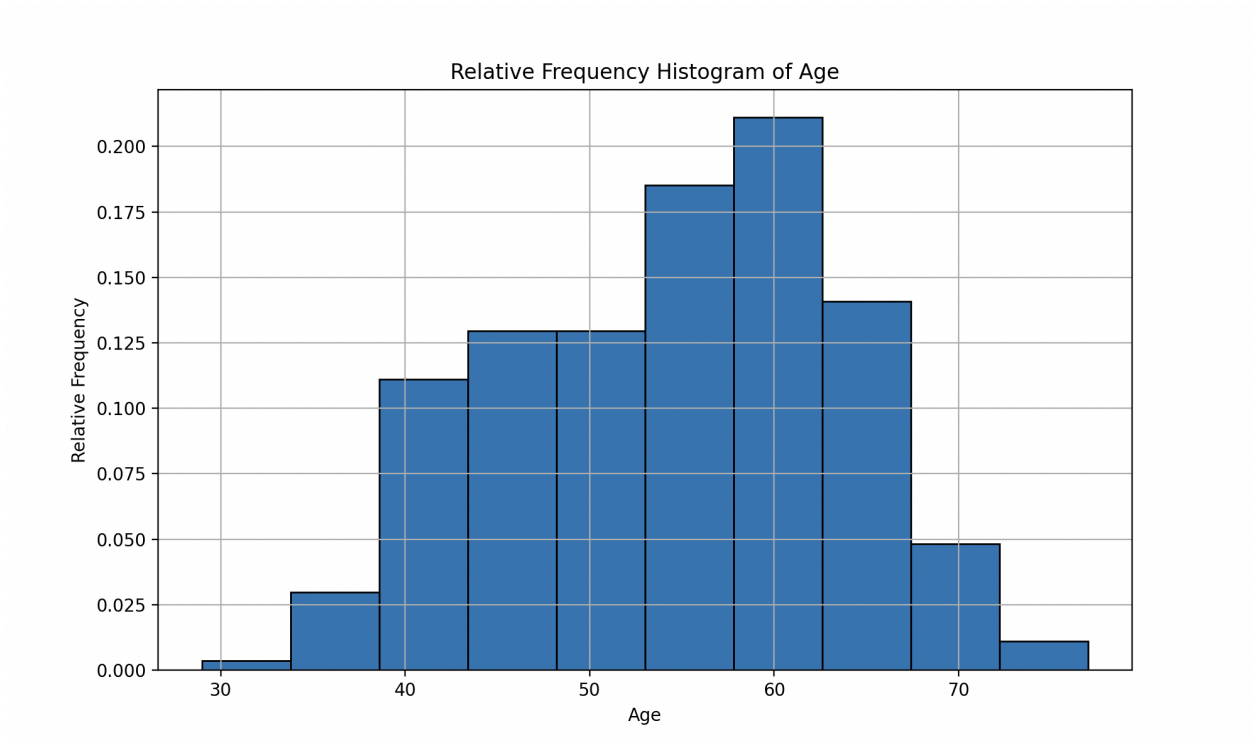
Q3: 61.0

Its code implementation can be seen [here](#). You can run this python file by cd'ing into **project\_part1/calculations/** and running **python3 calculate\_quartiles.py**.

II.d. Frequency Histogram



II.e. Relative Frequency Histogram



C should take the value 270.0 (Logic can be observed in the code).

### III. Code and Documentation

The code and associated documentation can be found [here](#). This link to my Github repository contains the below file structure and project history for your convenience. Else, code is provided in the appendix below.

```
/project_part1
├── /dataset1
│   ├── dataset1.py
│   └── dataset1code.py
├── /dataset2
│   ├── dataset2.py
│   └── dataset2code.py
├── /misc
│   ├── /calculations
│   │   ├── calculate_mean.py
│   │   ├── calculate_quartiles.py
│   │   └── calculate_variance.py
│   └── /datasets_txt
│       ├── dataset1info.txt
│       └── dataset2info.txt
```

## Appendix

### Reference 1: dataset1.py

```
1  #Taken from Kaggle, run file to download data locally
2
3  import kagglehub
4
5  path = kagglehub.dataset_download("mlippo/average-global-iq-per-country-with-other-stats")
6
7  print(["Path to dataset files:", path])
```

### Reference 2: dataset1code.py (File path is cut off in photo but can be observed in Github repo)

```
1  # Imports
2  import pandas as pd
3  import matplotlib.pyplot as plt
4
5  # File Path (relative)
6  file_path = "/Users/nicholasstone/.cache/kagglehub/datasets/mlippo/average-global-iq-per-country-with-other-stats/versions/3/avgIQp
7
8  # Interpreting data as Pandas dataframe
9  df_avg_iq = pd.read_csv(file_path)
10
11 # Looking to see if Average IQ and Literacy Rates occur in the data
12 if 'Average IQ' in df_avg_iq.columns and 'Literacy Rate' in df_avg_iq.columns:
13
14     # Creates the scatter plot
15     plt.figure(figsize=(10, 6))
16     plt.scatter(df_avg_iq['Average IQ'], df_avg_iq['Literacy Rate'], alpha=0.7)
17
18     # Labels and Title
19     # Label for x-axis
20     plt.xlabel('Average IQ')
21     # Label for y-axis
22     plt.ylabel('Literacy Rate (%)')
23     # Title
24     plt.title('Scatter Plot of Average IQ vs. Literacy Rate per Country')
25
26     # Background grid initialized for visual purposes
27     plt.grid(True)
28
```

```
28
29     # Initialize the scatterplot in a separate frame
30     plt.show()
31 ~ !lse:
32     # 'Error' statement
33     print("Ensure the dataset has 'Average IQ' and 'Literacy Rate' columns.")
34
```

### Reference 3: dataset2.py

```
1 #Taken from Kaggle, run file to download data locally
2
3 import kagglehub
4
5 path = kagglehub.dataset_download("luvhariishkhathi/heart-disease-patients-details")
6
7 print("Path to dataset files:", path)
8
9 heart_disease.csv
```

### Reference 4: dataset2code.py

```
1 #Imports
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import numpy as np
5
6 # File Path (relative)
7 file_path = "/Users/nicholasstone/.cache/kagglehub/datasets/luvhariishkhathi/heart-disease-patients-details/versions/1/heart_disease.csv"
8
9 # Interprets data as Pandas dataframe
10 df_heart_disease = pd.read_csv(file_path)
11
12 # Looks for Age column in data
13 if 'age' in df_heart_disease.columns:
14     data = df_heart_disease['age'].dropna()
15
16 # Defines k bins
17 k = 10 # Can be adjusted, but 10 is ideal number in my opinion, essentially separates by decades,
18         # which makes sense in the context of heart disease
19
20 # Creating the Frequency Histogram
21 plt.figure(figsize=(10, 6))
22 counts, bins, patches = plt.hist(data, bins=k, edgecolor='black')
23
24 # Title initialized
25 plt.title('Frequency Histogram of Age')
26
27 # X-axis titled for Age
28 plt.xlabel('Age')
29
30 # Y-axis titled for Frequency
31 plt.ylabel('Frequency')
32
33 # Initialized grid for visual pruposes
34 plt.grid(True)
35
36 # Initialized the histogram in a separate frame
37 plt.show()
38
39 # Defining C variable
40 C = counts.sum() # Total number of data points
41 relative_frequencies = counts / C
42
43 # Creating the Relative Frequency Histogram
44 plt.figure(figsize=(10, 6))
45 plt.bar(bins[:-1], relative_frequencies, width=np.diff(bins), edgecolor='black', align='edge')
46
47 # Title initialized
48 plt.title('Relative Frequency Histogram of Age')
49
50 # X-axis titled for age
51 plt.xlabel('Age')
52
53 # Y-axis titled for relative frequency
54 plt.ylabel('Relative Frequency')
55
56 # Initialized grid for visual purposes
57 plt.grid(True)
58
```

```

59 | # Initialized the relative histogram in a separate frame, keep in mind that the first histogram must be x'd
60 | #out of to see this second histogram
61 | plt.show()
62 |
63 | # Display C, must have x'd out of both histograms to see in terinal (i.e. the program has terminated)
64 | print(f"Value of C (normalizing constant): {C}")
65 | else:
66 |     #'Error' statement
67 |     print("Age not present in data.")

```

### Reference 5: calculate\_mean.py

```

1 | # Imports
2 | import pandas as pd
3 |
4 | #Function
5 | def calculate_mean(data):
6 |
7 |     #Return statement for sample mean
8 |     return sum(data) / len(data) if data else None
9 |
10 | #Main code execution
11 | if __name__ == "__main__":
12 |
13 |     #File path (relative)
14 |     file_path = "/Users/nicholasstone/.cache/kagglehub/datasets/luvharishkhathi/heart-disease-patients-details/versions/1/heart_disease.csv"
15 |
16 |     #Reads to Pandas dataframe
17 |     df = pd.read_csv(file_path)
18 |
19 |     #Looking for age in the dataframe
20 |     if "age" in df.columns:
21 |
22 |         #Converts age to a list
23 |         data = df["age"].dropna().tolist()
24 |
25 |         #Calculates sample mean by referring to python file
26 |         sample_mean_value = calculate_mean(data)
27 |
28 |         #Prints to terminal
29 |         print(f"Sample Mean is {sample_mean_value}")
30 |     else:
31 |
32 |         #'Error' message
33 |         print("Age not present in data.")

```



## Reference 6: calculate\_quartiles.py

```
1  #Imports
2  import pandas as pd
3
4  #Function to calculate the quartiles
5  def calculate_quartiles(data):
6
7      #Listing the quartiles and calculating them
8      quartiles = {
9          #1st quartile
10         "25 percentile": data.quantile(0.25),
11
12         #2nd quartile
13         "50 percentile": data.quantile(0.5),
14
15         #3rd quartile
16         "75 percentile": data.quantile(0.75)
17     }
18
19     #Return statement
20     return quartiles
21
22 #Main code execution
23 if __name__ == "__main__":
24
25     #File Path (relative)
26     file_path = "/Users/nicholasstone/.cache/kagglehub/datasets/luvhari/hkhati/heart-disease-patients-details/versions/1/heart_disease.csv"
27
28     #Uses Pandas to read dataframe
29     df = pd.read_csv(file_path)
30
```

```
31     #Checking for age data
32     if "age" in df.columns:
33         data = df["age"].dropna()
34         quartile_values = calculate_quartiles(data)
35         print("Quartile Values:")
36
37         #Individually prints quartiles to terminal
38         for quartiled, num in quartile_values.items():
39             print(f"{quartiled}: {num}")
40     else:
41         #'Error' statement
42         print("Age not present in data.")
43
```

## Reference 7: calculate\_variance.py

```
1  #Imports
2  import pandas as pd
3
4  #Function to calculate variance
5  def calculate_variance(data):
6      n = len(data)
7      if n < 2:
8          return None
9
10     #Sample variance calculation (See pdf for rationale)
11     mean_value = sum(data) / n
12     variance = sum((x - mean_value) ** 2 for x in data) / (n - 1)
13     return variance
14
15 #Main code execution
16 if __name__ == "__main__":
17
18     #File Path (relative)
19     file_path = "/Users/nicholasstone/.cache/kagglehub/datasets/luvharishkhathi/heart-disease-patients-details/versions/1/heart_disease.csv"
20
21     #Reads as a Pandas dataframe
22     df = pd.read_csv(file_path)
23
24     #Checks to see if age data exists
25     if "age" in df.columns:
26         data = df["age"].dropna().tolist()
27         variance_value = calculate_variance(data)
28
29         #Prints out sample variance final number
30         print(f"Sample Variance: {variance_value}")
31
32     else:
33         # 'Error' statement
34         print("Age not present in data.")
```

## Reference 8: Example of dataset1info.txt

Data in Rows and Columns:								
Rank	Country	Average IQ	Continent	Literacy Rate	Nobel Prices	HDI (2021)	Mean years of schooling - 2021	
1	Japan	106.48	Asia	0.99	29	0.925	13.4	
2	Taiwan	106.47	Asia	0.96	4	NaN	NaN	
3	Singapore	105.89	Asia	0.97	0	0.939	11.9	
4	Hong Kong	105.37	Asia	0.94	1	0.952	12.2	
5	China	104.10	Asia	0.96	8	0.768	7.6	
6	South Korea	102.35	Asia	0.98	0	0.925	12.5	
7	Belarus	101.60	Europe	1.00	2	0.808	12.1	
8	Finland	101.20	Europe	1.00	5	0.940	12.9	
9	Liechtenstein	101.07	Europe	1.00	0	0.935	12.5	
10	Germany	100.74	Europe	0.99	111	0.942	14.1	
11	Netherlands	100.74	Europe	0.99	22	0.941	12.6	
12	Estonia	100.72	Europe	1.00	0	0.890	13.5	
13	Luxembourg	99.87	Europe	1.00	2	0.930	13.0	
14	Macao	99.82	Asia	0.97	0	NaN	NaN	
15	Cambodia	99.75	Asia	0.78	0	0.593	5.1	
16	Canada	99.52	North America	0.99	28	0.936	13.8	
17	Australia	99.24	Oceania	0.99	12	0.951	12.7	
18	Hungary	99.24	Europe	0.99	13	0.846	12.2	
19	Switzerland	99.24	Europe	0.99	27	0.962	13.9	
20	United Kingdom	99.12	Europe	0.99	137	0.929	13.4	
21	North Korea	98.82	Asia	1.00	0	NaN	NaN	
22	Slovenia	98.60	Europe	1.00	1	0.918	12.8	
23	New Zealand	98.57	Oceania	0.99	3	0.937	12.9	
24	Austria	98.38	Europe	0.98	22	0.916	12.3	
25	Iceland	98.26	Europe	0.99	1	0.959	13.8	
26	Denmark	97.83	Europe	0.99	13	0.948	13.0	
27	Belgium	97.49	Europe	0.99	11	0.937	12.4	
28	United States	97.43	North America	0.99	400	0.921	13.7	

*Reference 9: Example of dataset2info.txt*

1	age	sex	chest	resting_blood_pressure				serum_cholesterol				fasting_blood_sugar	resting_electrocardiographic_results	maximum_heart_rate_
2	70	1	4	130	322	0	2	109	0	2.4	2	3	3	1
3	67	0	3	115	564	0	2	160	0	1.6	2	0	7	0
4	57	1	2	124	261	0	0	141	0	0.3	1	0	7	1
5	64	1	4	128	263	0	0	105	1	0.2	2	1	7	0
6	74	0	2	120	269	0	2	121	1	0.2	1	1	3	0
7	65	1	4	120	177	0	0	140	0	0.4	1	0	7	0
8	56	1	3	130	256	1	2	142	1	0.6	2	1	6	1
9	59	1	4	110	239	0	2	142	1	1.2	2	1	7	1
10	60	1	4	140	293	0	2	170	0	1.2	2	2	7	1
11	63	0	4	150	407	0	2	154	0	4.0	2	3	7	1
12	59	1	4	135	234	0	0	161	0	0.5	2	0	7	0
13	53	1	4	142	226	0	2	111	1	0.0	1	0	7	0
14	44	1	3	140	235	0	2	180	0	0.0	1	0	3	0
15	61	1	1	134	234	0	0	145	0	2.6	2	2	3	1
16	57	0	4	128	303	0	2	159	0	0.0	1	1	3	0
17	71	0	4	112	149	0	0	125	0	1.6	2	0	3	0
18	46	1	4	140	311	0	0	120	1	1.8	2	2	7	1
19	53	1	4	140	203	1	2	155	1	3.1	3	0	7	1
20	64	1	1	110	211	0	2	144	1	1.8	2	0	3	0
21	40	1	1	140	199	0	0	178	1	1.4	1	0	7	0
22	67	1	4	120	229	0	2	129	1	2.6	2	2	7	1
23	48	1	2	130	245	0	2	180	0	0.2	2	0	3	0
24	43	1	4	115	303	0	0	181	0	1.2	2	0	3	0
25	47	1	4	112	204	0	0	143	0	0.1	1	0	3	0
26	54	0	2	132	288	1	2	159	1	0.0	1	1	3	0
27	48	0	3	130	275	0	0	139	0	0.2	1	0	3	0
28	46	0	4	138	243	0	2	152	1	0.0	2	0	3	0
29	51	0	3	120	295	0	2	157	0	0.6	1	0	3	0
30	58	1	3	112	230	0	2	165	0	2.5	2	1	7	1