

## LISÄTEHTÄVÄT viikko 35

Jenni Ylisirniö

### Tehtävä 1

Tallenna oheinen diabetes.csv ja lataa se pandasin dataframeen. Tulosta muutamia tunnuslukuja mm. count, mean, min, max, std.

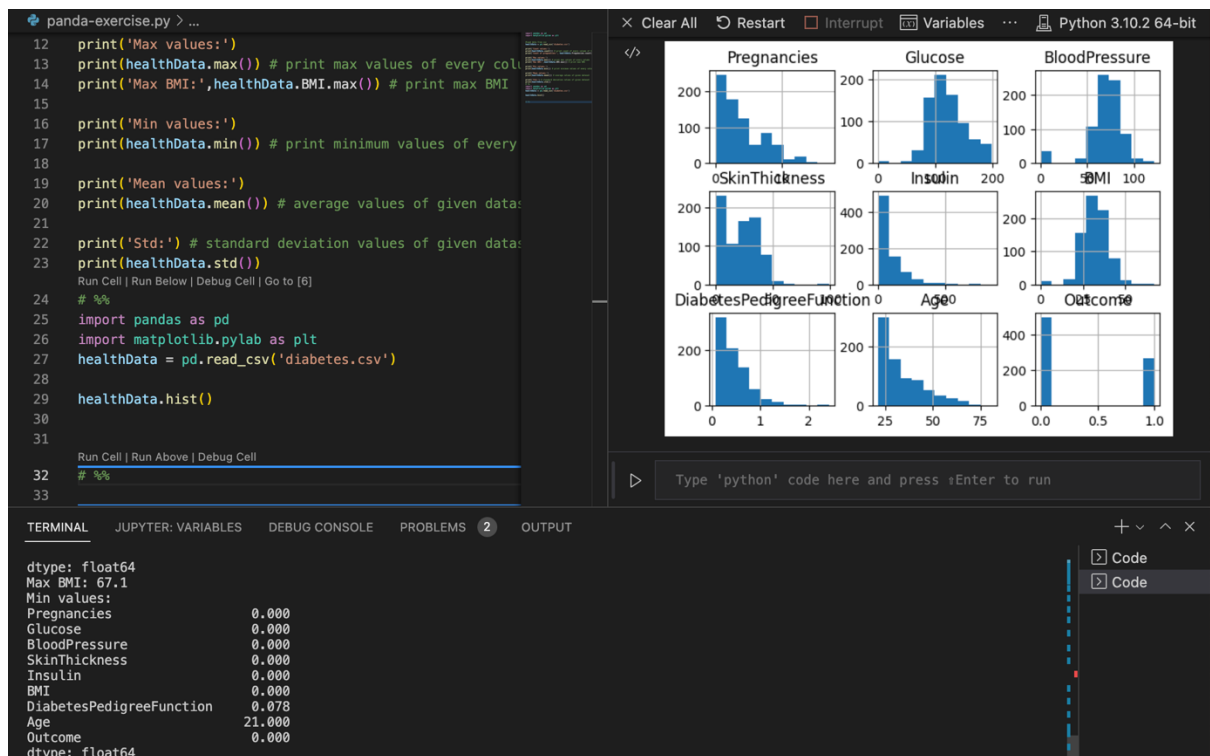
```
panda-exercise.py > ...
1  import pandas as pd
2
3  #load data from csv
4  healthData = pd.read_csv('diabetes.csv')
5
6  print('Count values:')
7  print(healthData.count()) # print count of every column of the dataset
8  print('Count of pregnancies:', healthData.Pregnancies.count()) # print count of pregnancies
9
10 print('Max values:')
11 print(healthData.max()) # print max values of every column
12 print('Max BMI:', healthData.BMI.max()) # print max BMI
13
14 print('Min values:')
15 print(healthData.min()) # print minimum values of every column
16
17 print('Mean values:')
18 print(healthData.mean()) # average values of given dataset
19
20 print('Std:') # standard deviation values of given dataset
21 print(healthData.std())
```

TERMINAL	JUPYTER	DEBUG CONSOLE	PROBLEMS	OUTPUT
Age		33.240885		
Outcome		0.348958		
dtype: float64				
Std:				
Pregnancies		3.369578		
Glucose		31.972618		
BloodPressure		19.368155		
SkinThickness		15.952218		
Insulin		115.244002		
BMI		7.884160		
DiabetesPedigreeFunction		0.331329		
Age		11.760232		
Outcome		0.476951		

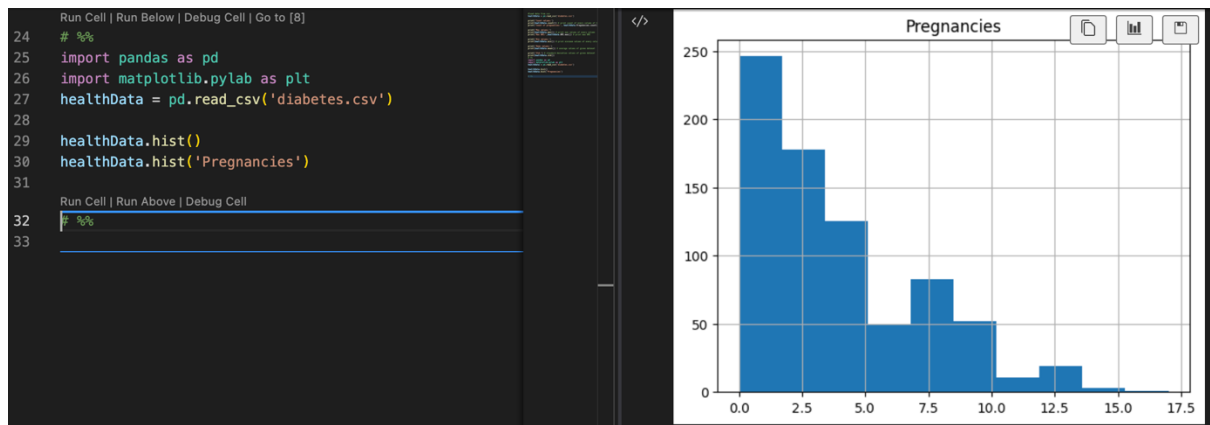
Screenshotissa näkyy vain pieni osa ohjelman ajosta, mutta ajo toimii täydellisesti.

## Tehtävä 2

Piirrä histogrammi kuvaaja datasta:



Yksi histogrammi näyttää vähän siistimmältä:



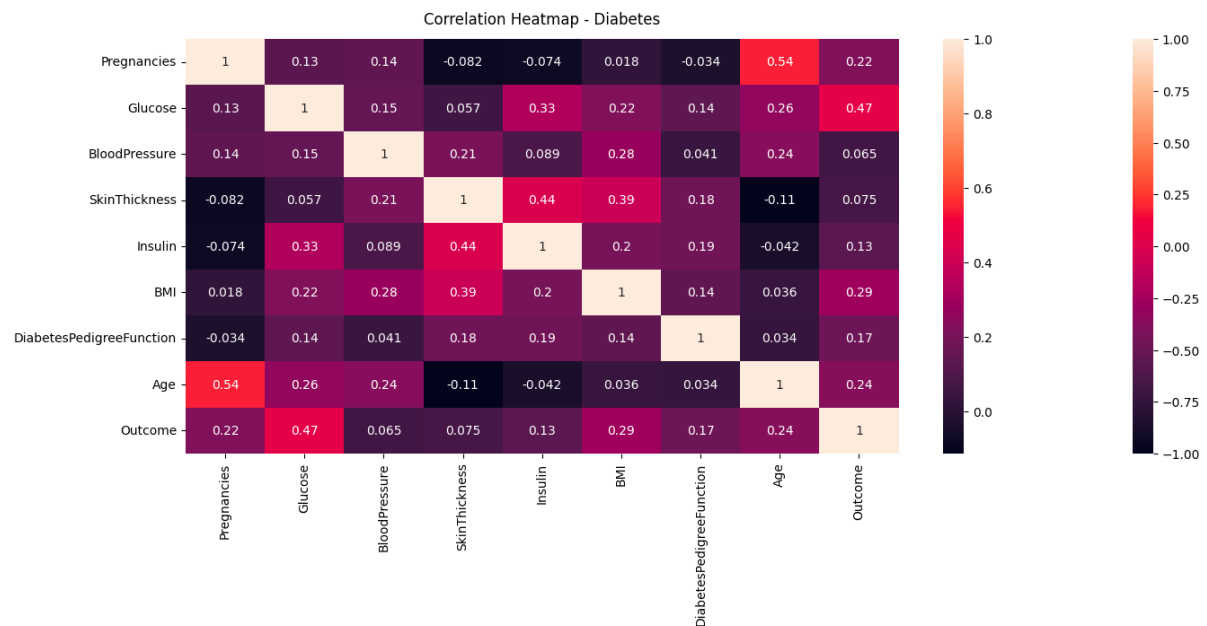
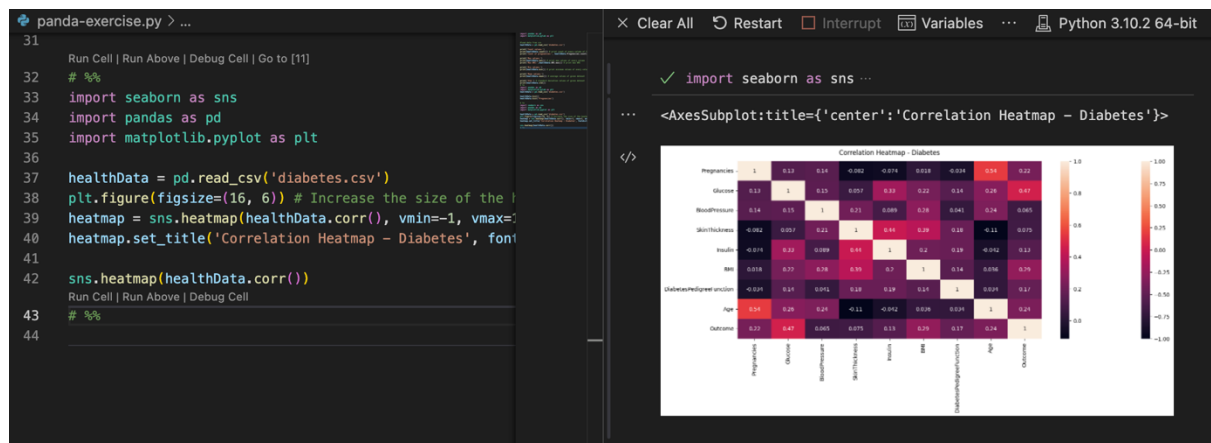
## Tehtävä 3

Piirrä korrelaatioheatmap datasta:

```
import seaborn as sns
import pandas as pd
import matplotlib.pyplot as plt

healthData = pd.read_csv('diabetes.csv')
plt.figure(figsize=(16, 6)) # Increase the size of the heatmap.
heatmap = sns.heatmap(healthData.corr(), vmin=-1, vmax=1, annot=True)
heatmap.set_title('Correlation Heatmap - Diabetes', fontdict={'fontsize':12},
pad=12)

sns.heatmap(healthData.corr())
```



<https://medium.com/@szabo.bibor/how-to-create-a-seaborn-correlation-heatmap-in-python-834c0686b88e>

## Tehtävä 4

Laske potilaiden lukumäärä iän mukaan siten, että suurin lukumäärä on ensin.

```
44 import pandas as pd
45
46 healthData = pd.read_csv('diabetes.csv')
47 dups = healthData.pivot_table(index = ['Age'], aggfunc = 'size')
48 print(dups)
49
Run Cell | Run Above | Debug Cell
50 # %%
51
```

Output exceeds the [size limit](#). Open the full output [data in a text editor](#)

Age	
21	63
22	72
23	38
24	46
25	48
26	33
27	32
28	35
29	29
30	21
31	24
32	16
33	17
34	14

Osasin laskea potilaiden lukumäärän iän mukaan, mutta järjestystä en osannut vaihtaa niin, että suurin lukumäärä olisi ensin. Yritin `sort_values`, mutta ei jostain syystä onnistunut.

Tulosta myös montako diabetestapausta (1) ja ei diabetesta (0) on aineistossa.

```
Run Cell | Run Above | Debug Cell | Go to [34]
50 # %%
51 import pandas as pd
52
53 healthData = pd.read_csv('diabetes.csv')
54 dups = healthData.pivot_table(index = ['Outcome'], aggfunc = 'size')
55 print(dups)
Run Cell | Run Above | Debug Cell
56 # %%
57
```

dtype: int64

✓ import pandas as pd ...

... Outcome

0	500
1	268

dtype: int64

Eli 500 ei ole (0) diabetestä ja 268 on (1) diabetes.

## Tehtävä 5

Onko aineistossa nan arvoja? Jos on, missä sarakkeissa ja montako?

```
panda-exercise.py > ...
56 # %%
57 import pandas as pd
58
59 healthData = pd.read_csv('diabetes.csv')
60 count_of_nan = healthData.isna().sum()
61 print('COUNT NAN VALUES')
62 print(count_of_nan)
63 healthData.notnull()
64
65
```

✓ import pandas as pd ...

... COUNT NAN VALUES

Pregnancies	0
Glucose	0
BloodPressure	1
SkinThickness	0
Insulin	0
BMI	0
DiabetesPedigreeFunction	0
Age	0
Outcome	0

dtype: int64

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Diabe
0	True	True	False	True	True	True	
1	True	True	True	True	True	True	
2	True	True	True	True	True	True	

Aineistossa on yksi NaN (not a number) arvo. Se löytyy Blood Pressure -sarakeesta, riviltä yksi.