

## CSC 604, Fall 2024, Assignments 8 and 9

handed out Nov. 11, due Nov. 18 at 11:59pm on Canvas

In this assignment, we are moving over to the area of Data Analysis. In particular, we will analyze more deeply the data set you encountered in Assignment 5. This time, we will use the full data set rather than the portion we used in Assignment 5. The dataset contains the medical records of patients who suffered from heart failure, collected during their follow-up period, where each patient profile has 13 clinical features. The database allows researchers to explore the incidence of death during the follow-up period of patients who suffered heart failure. In particular, it can help researchers associate clinical or demographic features with the chance of death after heart failure (True, this is not a cheerful data set! I could have found something a bit more fun!). This is a 2-part assignment covering both the data science concepts discussed in class such as the PANDA concept of a dataframe and its manipulation and the visualization of your findings using Matplot Lib. In particular, you will familiarize yourself with three types of graphs in the Matplot library of Python: Pie Charts, Bar Plots, and simple Bar Plots.

You are asked to work using Jupyter Notebooks instead of PyCharm or Spyder since I would like you to familiarize yourself with the type of IDE most commonly used when using Python in Data Science.

### **Here is the description of the assignment:**

We will be working on the full Heart Failure Clinical Records data set already partially used in Assignment 5: <https://archive.ics.uci.edu/dataset/519/heart+failure+clinical+records>. Don't lose the link we will be also using the data set in the next assignment!

The data set contains the records of patients represented according to the following variables:

- age: age of the patient (years)
- anaemia: decrease of red blood cells or hemoglobin (boolean)
- creatinine phosphokinase (CPK): level of the CPK enzyme in the blood (mcg/L)
- diabetes: if the patient has diabetes (boolean)
- ejection fraction: percentage of blood leaving the heart at each contraction (percentage)
- high blood pressure: if the patient has hypertension (boolean)
- platelets: platelets in the blood (kiloplatelets/mL) - sex: woman or man (binary)
- serum creatinine: level of serum creatinine in the blood (mg/dL)
- serum sodium: level of serum sodium in the blood (mEq/L)
- smoking: if the patient smokes or not (boolean)
- time: follow-up period (days)
- death event: if the patient died during the follow-up period (boolean)

### **Part I: PANDAS**

Your first task will be to read the file and display it.

The screenshot shows a JupyterLab environment with a data table displayed in a code cell. The table has 9 rows and 13 columns. The columns are: age, anaemia, creatinine\_phosphokinase, diabetes, ejection\_fraction, high\_blood\_pressure, platelets, serum\_creatinine, and serum\_cholesterol. The data is as follows:

	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelets	serum_creatinine	serum_cholesterol
0	75.0	0	582	0	20	1	265000.00	1.9	
1	55.0	0	7861	0	38	0	263358.03	1.1	
2	65.0	0	146	0	20	0	162000.00	1.3	
3	50.0	1	111	0	20	0	210000.00	1.9	
4	65.0	1	160	1	20	0	327000.00	2.7	
...	...	...	...	...	...	...	...	...	...
94	62.0	0	61	1	38	1	155000.00	1.1	
95	55.0	0	1820	0	38	0	270000.00	1.2	
96	45.0	0	2060	1	60	0	742000.00	0.8	
97	45.0	0	2413	0	38	0	140000.00	1.4	
98	50.0	0	196	0	45	0	395000.00	1.6	

The interface also shows a sidebar with installed extensions: @jupyter-widgets/jupyterlab-manager, jupyterlab-plotly, and @pyviz/jupyterlab\_pyviz. The bottom status bar indicates the system time as 8:32 AM on 11/3/2024.

Please also display general information/statistics about each column of the data set in the following two ways, a) and b):

a)

The screenshot shows a JupyterLab environment with a summary statistics table displayed in a code cell. The table has 8 rows and 9 columns. The columns are: age, anaemia, creatinine\_phosphokinase, diabetes, ejection\_fraction, high\_blood\_pressure, platelets, serum\_creatinine, and serum\_cholesterol. The data is as follows:

	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelets	serum_creatinine	serum_cholesterol
count	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	
mean	60.833893	0.431438	581.839465	0.418060	38.083612	0.351171	263358.029264		
std	11.894809	0.496107	970.287881	0.494067	11.834841	0.478136	97804.236869		
min	40.000000	0.000000	23.000000	0.000000	14.000000	0.000000	25100.000000		
25%	51.000000	0.000000	116.500000	0.000000	30.000000	0.000000	212500.000000		
50%	60.000000	0.000000	250.000000	0.000000	38.000000	0.000000	262000.000000		
75%	70.000000	1.000000	582.000000	1.000000	45.000000	1.000000	303500.000000		
max	95.000000	1.000000	7861.000000	1.000000	80.000000	1.000000	850000.000000		

The interface also shows a sidebar with installed extensions: @jupyter-widgets/jupyterlab-manager, jupyterlab-plotly, and @pyviz/jupyterlab\_pyviz. The bottom status bar indicates the system time as 9:33 AM on 11/3/2024.

b)

The screenshot shows a JupyterLab interface with a file browser on the left and a code editor in the center. The code editor displays a data table with the following columns: age, anaemia, creatinine\_phosphokinase, diabetes, ejection\_fraction, high\_blood\_pressure, platelets, and serum. The table contains three rows of data: mean, Std.Dev, and Var. The values are as follows:

	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelets	serum
mean	60.833893	0.431438	581.839465	0.418060	38.083612	0.351171	2.633580e+05	
Std.Dev	11.894809	0.496107	970.287881	0.494067	11.834841	0.478136	9.780424e+04	
Var	141.486483	0.246122	941458.571457	0.244102	140.063455	0.228614	9.565669e+09	

You must then find out and print the following quantities:

The number of people 45 or less in the database is A.

The number of people older than 45 in the database is B.

The youngest person in the database is C years old.

The oldest person in the database is D years old.

The number of people in the database who died following a heart failure event is E out of F patients.

The number of people in the database who survived following a heart failure event: G out of F patients.

In other words, the percentage of people who died following a heart failure event is H%.

The percentage of young people ( $\leq 45$ ) who died after a heart failure event is I%.

The percentage of old people ( $> 45$ ) who died after a heart failure event is J%.

The percentage of smokers who died after a heart failure event is K%.

The number of men in the database is L.

The number of men survivors in the database is M or N%.

The number of women in the database is O.

The number of women survivors in the database is P or Q%.

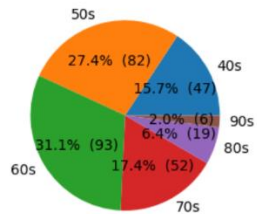
Letters A-Q above need to be replaced by the correct numbers. If the number is a decimal, please display two decimal places. If it is an integer, please display it as such, without decimals.

In addition, if you are interested in additional statistics about this data set, please go ahead and include them.

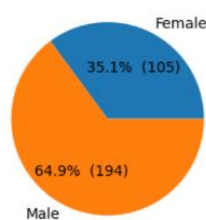
## PART II: Matplotlib Lib

Your first task will be to show pie charts to capture the proportions of age, sex, and survivors in the Heart Failure data set. The age groups considered should be 40s (< 50), 50s (>= 50; < 60), 60s (>= 60; < 70), 70s (>= 70; < 80), 80s (>= 80; < 90), 90s (>= 90).

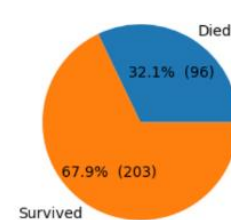
Age of Patients who suffered heart failure



Proportion of Male/Female who experienced Heart Failure

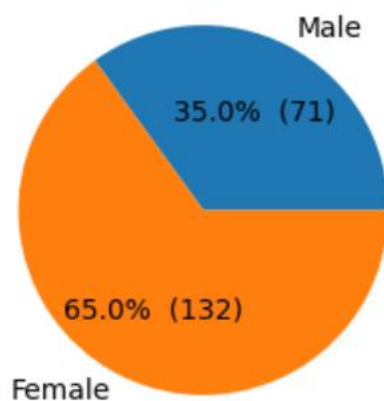


Heart Failure Survivors

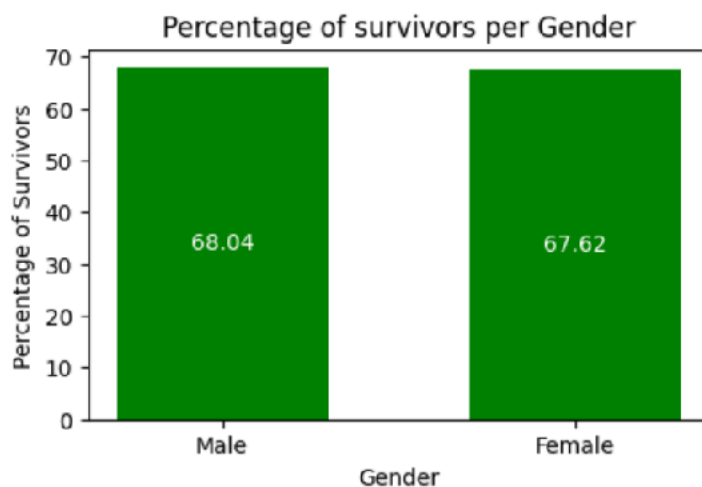


Another Pie chart will capture the proportion of male to female Heart Failure survivors:

Proportion of Male/Female Heart Failure Survivors



Next, you are asked to draw a bar graph showing the percentage of survivors per gender:



Finally, you will plot a graph tracing the percentage of Heart Failure survivors per age group.

