

HOMEWORK 1 ASSIGNMENT

Assigned: Thursday, August 27, 2015 **Deadline:** Tuesday, September 1 by 5:00pm

The purpose of this first homework is to get us exploring data using iPython notebook and pandas/numpy. The only way to “fail” this one is to not try.

In this assignment, we will gain practice using:

- iPython notebook
- Relevant Python packages

DATA & CONTEXT

In this assignment, we will explore the passenger list of the Titanic, as provided in a well-known Kaggle competition. For this assignment, we are concerned only with initial exploration. We may build a predictive model later, but not as part of this assignment. The focus of the assignment is to answer the specific questions listed below in the section “Homework Questions.”

The dataset is a list of passengers. The second column of the dataset is a “label” for each person indicating whether that person survived (1) or did not survive (0). Here is the Kaggle page with more information on the dataset (data dictionary):

<http://www.kaggle.com/c/titanic-gettingStarted/data>

Don’t worry about downloading the data from Kaggle; we have provided the HW1 dataset for you. Download it from Github and Move this data file to the directory where you like to do your iPython notebook work.

SUBMITTING YOUR WORK

Submit your work by placing the iPython Notebook in the folder in Google Drive with the file name \$YOUR_LAST_NAME_HW1

HOMEWORK QUESTIONS

Please answer the following questions about your data exploration in the iPython notebook. Feel free to explore further. These questions are a guide and a minimum, not a limit ;-)

1. How many passengers are in our passenger list? From here forward, we'll assume our dataset represents the full passenger list for the Titanic.
2. What is the overall survival rate?
3. How many male passengers were onboard?
4. How many female passengers were onboard?
5. What is the overall survival rate of male passengers?
6. What is the overall survival rate of female passengers?
7. What is the average age of all passengers onboard?
 - a. How did you calculate this average age?
 - b. Note that some of the passengers do not have an age value. How did you deal with this? What are some other ways of dealing with this?
8. What is the average age of passengers who survived?
9. What is the average age of passengers who did not survive?
10. At this (early) point in our analysis, what might you infer about any patterns you are seeing?
11. How many passengers are in each of the three classes of service (e.g. First, Second, Third?)
12. What is the survival rate for passengers in each of the three classes of service?
13. What else might you conclude?
14. Last, if we were to build a predictive model, which features in the data do you think we should include in the model and which can we leave out? Why?