

See Dropbox for Due Date
T. Kinsman

Homework is to be programmed in Python, R, Matlab, or Java.

When coding, assume that the grader has no knowledge of the language or API calls but can read comments. Use prolific comments before each section of code, or function call to explain what the code does, and why you are using it.

Do not use single letter variable names. That is for theory classes. This is an applied class.

Hand in your results, and the commented code, in the associated dropbox. Again, put all files in a directory with your name on it, HWNN_LastName_Firstname. Then zip up the entire directory and submit the zip file.

Inside the directory we should find two files:

- A. HWNN_<LASTNAME>_<Firstname>_results.pdf – your write up of what you did.
- B. HWNN_<LASTNAME>_<Firstname>_program.ext – your program.

Substitute the homework number for NN.

Feel free to look over each other's shoulders, at each other's work, but do your own work.
Let me know whom you worked with. Do not hand in copies of each other's code.

Background Ridiculousness from a Fictional Universe:

Abominable snowfolk have been observed in the Himalayan Mountains. After decades of careful and expensive observations, the associated data was collected. What most of us do not know is that when they mature, the hair on the shoulders of most abominable snowfolk turns gray. This is like the silver-backed great apes in other parts of the world.

It is now known that stress causes mammal's hairs to lose color and turn gray. The mechanism was recently discovered. (Not kidding: <https://thenextweb.com/science/2020/01/23/scientists-figured-out-why-stress-turns-your-hair-gray/>)

The data collected includes information about each individual observed. The curated data available to you now includes the age at which the individual started having gray hair on their shoulders, and the individual's approximate height in cm.

Overview:

It is believed that there are two sub-species of these snowfolks. Cluster them into two types, based on the age at which they matured.

Details follow:

(continued)

1. **Exploratory Data Analysis:** (2 pts)

You are also provided with some mystery data, in the file MysteryData.txt. It consists of two underlying groups. This data is pre-quantized to the nearest unit.

- a. Compute the average and standard deviation of this data? Use a package. **What value did you get?**
- b. Remove the last value from the data (the last row), and then re-compute the average value. How did average values change? Why do you think you observed these changes? What caused this amount of change?

Hint: we know that outliers pull the computation for the mean value off. A point that is very far from the average will change the average. Look at how much the computation for your average value changed. Then look at how much of an outlier the data point you removed was. Make your own conclusions here.

Also, remember the rule of physics that says that the center of mass is a good representative for an object. What happens if you add or subject weight at the center of mass? I don't like having to give away the punch line, but make sure you understand this point.

2. **1D Clustering using Otsu's method on the age. (5 pts total)**

Implement Otsu's method from scratch. (1 pts)

Use Otsu's method to find the best age to form two clusters with.

Make the decision based on "age \leq this_threshold".

What is the best threshold you found?

- We are going to do some noise removal. Quantizing into bins removes small variations in the data.
- Quantize the snowfolks age into bins that are 2 years interval using the ranges the ranges [0-2), [2,4), ... so the first bin is age zero up-to (but not including) age 2.
`quantized_data = floor(raw_data / BIN_SIZE_for_AGE) * BIN_SIZE_for_AGE`
- Implement Otsu's method to separate the snowfolk's population into two clusters. That is, we are using Otsu's method to binarize the data. There are other methods. We are quantizing the data into two groups.

Answer these Questions:

- a. What age should we use to best separate the two clusters? (1 pt.)
- b. What is the minimum mixed variance that resulted? (1 pt.)
- c. Breaking Ties: How would your program handle a situation where the minimum mixed variance occurred twice? Does this situation happen? (1 pt.)
- d. What other methods could we have used to quantize the data? (1 pt.)

(continued)

3. (2 pt) Exploring Regularization:

You are going to replace Otsu's method from the previous routine with a cost function.

Let: $\text{Cost_Function} = \text{Objective_Function} + \text{Regularization}$.

The objective function, the thing you really care about, is Otsu's method.

Let the mixed variance be the objective function, as previously defined. The objective function is Otsu's method, as you just computed. So, instead of finding the best threshold for Otsu's method, you will now find the best threshold for the Cost Function. In other words, you will find the best threshold for

$\text{Cost} = (\text{Otsu's method}) + (\text{Regularization Term})$

We want to add a regularization term to help make the two clusters the same size, however, it should not overwhelm the objective function. For understanding the next equation, imagine alpha is 1 for a second.

$\text{Regularization} = \alpha * \text{abs}((\text{Number of Points in First Group}) - (\text{Number of Points in Second Group})) / \text{NormFactor}$

Set the NormFactor to 100. This is about half the size of the data.

Explore different values of alpha around:

[100, 1, 1/5, 1/10, 1/20, to 1/25, to 1/50, and 1/100, 1/1000].

Which value of alpha causes the "best" splitting point to change? Do you notice anything?

There is not necessarily a change.

4. (2 pt) Graphing

Ignoring regularization, plot a graph of the mixed variance for the snowfolk's data that is given to you versus the quantized age.

Add a circular point indicating the value used to segment the data into two clusters.

Clearly label all the axes.

You should see that the function goes down, and then goes back up again.

If you run out of time, use Excel to do the plot. However, I highly recommend you learn matplotlib. You want to add this to your resume. AND, it helps to not rely on some company's product.

AND it will help you the rest of the semester.

5. (1 pt) Conclusion and Discussion:

Write a conclusion that describes what you learned in this homework.

Feel free to repeat answers from the previous questions.

How do these learning relate to life and physics as we know it?

How do your results relate to life, the universe, and everything?

This is college. Write at least a five-sentence paragraph here.