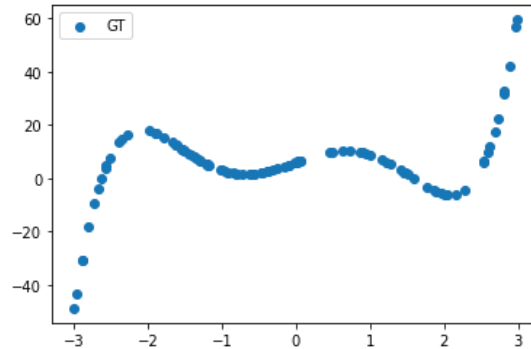


**Assignment 1**  
**CS 4783/5783**  
**Due: 09/08/2020 11:59 am**

**[Question 1]**

**[30 points]**

Suppose that you are conducting a scientific experiment where you are observing the effects of one variable ( $x_{train.npy}$  and  $x_{test.npy}$ ) on the output ( $y_{train.npy}$  and  $y_{test.npy}$ ). On visualizing the relationship between the variables, you see the following plot:



Your goal is to come up with a linear regression model that can take the training data ( $x_{train.npy}$  and  $y_{train.npy}$ ) and model the relationship between the variables  $x$  and  $y$ . You should implement your own version of linear regression either using gradient descent or normal equations. You **SHOULD NOT** use any pre-packaged library such as Sci-Kit Learn. However, you may use these libraries to verify your output if you wish.

Here are some things to keep in mind for tackling this problem:

1. Try to plot this relationship on your own using [matplotlib](#). You can also visualize the test data to see if it gives you any clues about the underlying relationship between the variables.
2. Use your knowledge gleaned from the previous step to answer the following questions:
  - a. Is the relationship linear?
  - b. Do I need feature engineering to add any non-linearity?
    - i. If so, how can I engineer these features? *[Hint: Basis Functions!]*
    - ii. What are some functions that I can try?
      1. Plot each of them individually to verify!

You will need to write a short report detailing your thought process, the code you wrote in Python to implement the linear regression model and the equation that models the relationship between  $x$  and  $y$  that you found. You should provide evidence that corroborates your final statement such as plots, prediction errors, etc.

**[Question 2]**

**[20 points]**

Imagine that you are a realtor in Hogsmeade. You have data points that correspond to the recent sales of different houses in and around Hogsmeade. Your goal is to help Hagrid estimate the prices of houses that he can use to sell or buy listings. Can you use your knowledge of linear regression to find the best regression model? Use your implementation from Question 1 (without any basis functions) to answer the following questions.

1. What is the average least squares error for the given data using your simple linear regression model?
2. Which factor has the most effect on the final value? How do you know this? Can you use only this feature to predict the price?

3. Which factor has the least effect on the final value? How do you know this? What effect does removing this feature have on the performance?

**Submission Requirements:**

You will need to submit the following as a single ZIP file:

1. A short report detailing your work and answers to the questions presented above.
2. Your code as a IPython notebook that can be run on Google Colab.
3. A README file on any dependencies that are required to run your code.

**Note:**

1. If your code does not run on Colab, you will not get any credit for the code segment. We will only grade what is in your report.
  - a. This includes any syntax errors due to indentation, unnamed/unknown libraries that were not listed in the README file, etc.
2. Please submit code only in Python and in the IPython notebook format. You can write your answers as part of the notebook if you do not want a separate report file, but it must be comprehensive.
  - a. Any code not in Python will not be graded at all.