

## COMP 5600/6600

### Assignment 2 – Basic Machine Learning

**Due:** 10/09/2023 11:59 pm

**Submission:** On Canvas

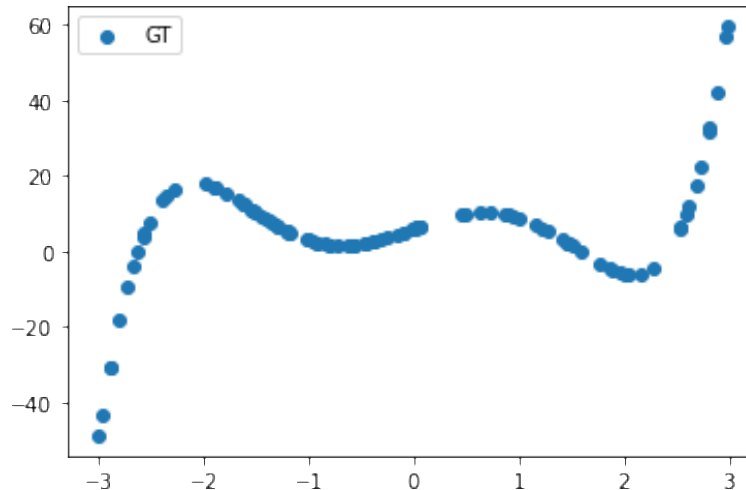
#### Overview:

In this assignment, you will be implementing and analyzing three different machine learning algorithms. Each algorithm will focus on a different problem scenario that can be solved by different learning methods. You must implement all algorithms from scratch using Python and should not use any external libraries other than NumPy. You can use other libraries for data loading and pre-processing as necessary.

#### [Problem 1]

[25 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 using gradient descent. Things to consider:

1. Try to plot this relationship on your own using matplotlib. You can also visualize the test data (`x_test.npy` and `y_test.npy`) 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 you need feature engineering to add any non-linearity?
  - c. If so, how can you engineer these features?
  - d. What are some functions that you can try?
  - e. 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.

### **[Problem 2]**

**[25 Points]**

See the attached housing data (Assignment2\_NB\_Data.xlsx). Each tab in the Excel file contains training and test splits. Your goal is to construct a Naïve Bayes classifier for this data.

1. Compute and show the conditional probability distribution for each feature. Explain how you got these values and show your work. Note: You are expected to do this part of the question by hand. Explain how you got the probability distribution for at least two features in detail.
2. Using your conditional probability table, write a Python code that will compute the probabilities for each example in the test data. Your program should output the probabilities of each class as well as the final classification based on the MAP rule. Note: You should hard-code the conditional probabilities from the previous step into your code.
3. Analyze the performance of the code using different metrics (accuracy, precision, recall, etc.) and briefly discuss your insights about the performance of the approach (good or bad) and the merits of using one metric over the other. You can use plots and other mechanisms to support your conclusions.

### **[Problem 3]**

**[25 Points]**

In this question, you will be using k-means to perform image compression. Implement a naïve version of the k-means algorithm based on your understanding. Your code must take the number of clusters  $k$  as input and perform k-means clustering on the given image (test\_image.png). Once the algorithm finishes running, the cluster centroids represent the top- $k$  common colors in the image. Iterate through each pixel in the image and assign the closest color to each pixel. Save and visualize

the resulting image. For reading and writing images, you can use OpenCV, which is an open-source computer vision toolkit. The following code will load the image into a NumPy array. You can use this as input to your K-Means algorithm.

```
import cv2
img = cv2.imread('test_image.png')
height, width, channels = np.shape(img)
for i in width:
    for j in height:
        pixel = img[j][i] # Read the pixel at location (i,j)
        img[j][i] = newValue # Assign a new value to the pixel
```

Experiment with different values of  $k$  and briefly describe your thoughts about which value works best for this problem. You can use plots, error bars, etc. to support your conclusions.

### **Deliverables:**

A single IPython Notebook that contains your code and report. You can use the text cells to write your report and embed any plots, illustrations, and/or images that you need to support your claims.