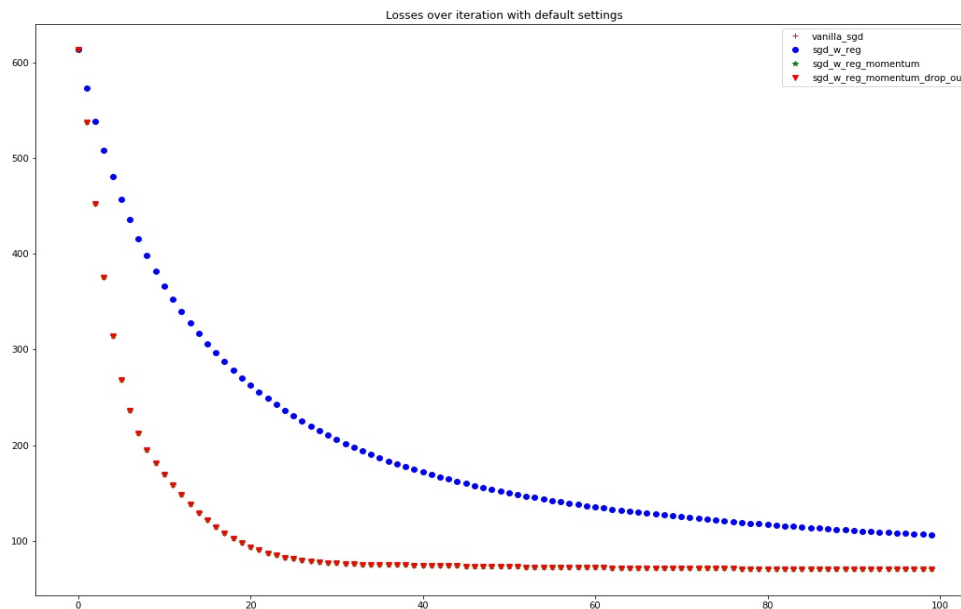# Question 5: Building a Neural Networks

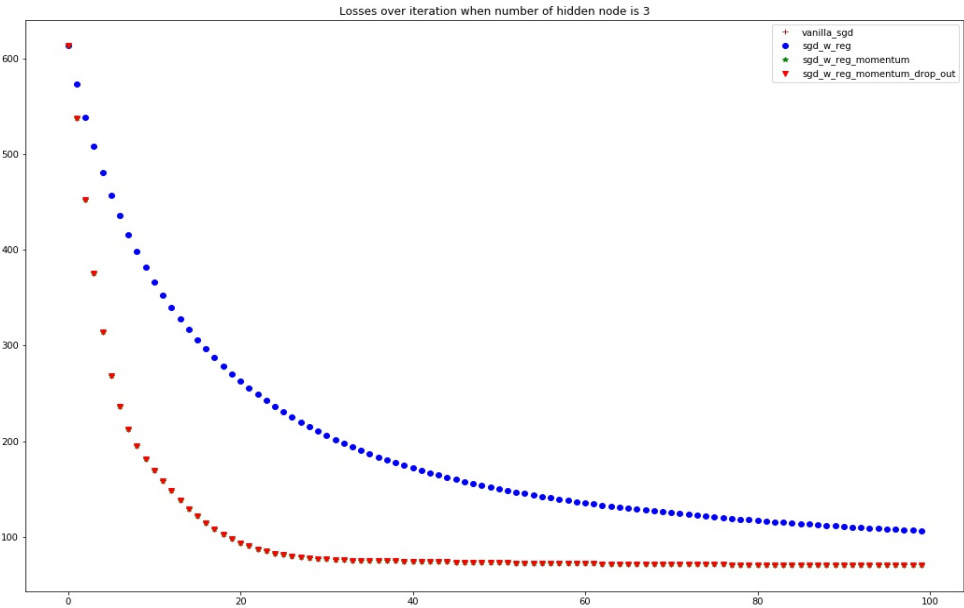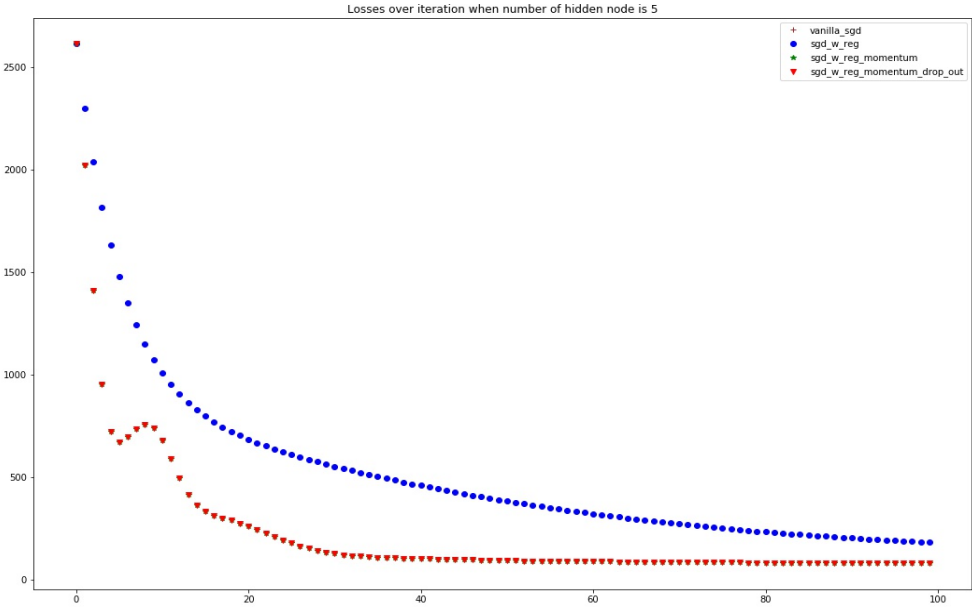## Part 1: For a ReLu network:

**A**



(a) From the figure above, we can see vanilla-sgd and sgd-w-reg have relative similar performance, while sgd-w-momentum and sgd-w-reg-momentum-drop-out have similar performance. The later two perform better than the previous two methods. We can also see that after 100 iterations, sgd-w-momentum and sgd-w-reg-momentum-dropout achieved same loss.

(b) The momentum methods works obviously better since momentum helps accelerate gradients in the right direction during training process. Therefore, we can see loss drops faster in momentum methods with same number of iterations. However, both regularization and dropout does little effect on the performance, and which suggest that the model may be underfit under 100 iterations.
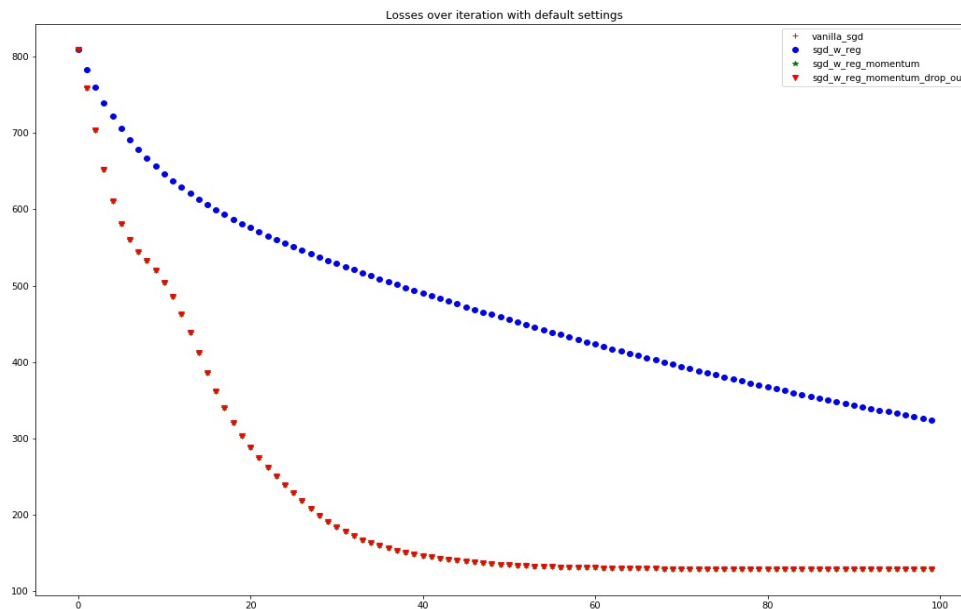
**B**

Losses over iteration when number of hidden node is 5



Losses over iteration when number of hidden node is 3

(a) Both sgd-w-momentum and sgd-w-reg-momentum-drop-out achieved same loss after 100 iterations. They perform better than the first two methods.

(b) The performance is not better than when we use 3 hidden nodes. The results suggest that the model starts to overfit when number of hidden node is 5.

**C**

| Model | Time |
|---|---|
| vanilla-sgd | 1.4158 |
| sgd-w-reg | 1.1950 |
| sgd-w-momentum | 1.0863 |
| sgd-w-reg-momentum-drop-out | 0.9850 |

# Part 2: For a Sigmoid network

**A**



From the figure above, we can see vanilla-sgd and sgd-w-reg have relative similar performance, while sgd-w-momentum and sgd-w-reg-momentum-drop-out have similar performance. The

later two perform better than the previous two methods. We can also see that after 100 iterations, sgd-w-momentum and sgd-w-reg-momentum-drop-out achieved same loss.

The momentum methods works obviously better since momentum helps accelerate gradients in the right direction during training process. Therefore, we can see loss drops faster in momentum methods with same number of iterations. However, both regularization and dropout does little effect on the performance, and which suggest that the model may be underfit under 100 iterations.

**B**



Losses over iteration with default settings

Both sgd-w-momentum and sgd-w-reg-momentum-drop-out achieved same loss after 100 iterations. They perform better than the first two methods. This is better than when number of hidden nodes is 3, which means that we need more complex models to learn better from the data.

**C**

| Model | Time |
|---|---|
| vanilla-sgd | 1.1793 |
| sgd-w-reg | 1.1569 |
| sgd-w-momentum | 1.0662 |
| sgd-w-reg-momentum-drop-out | 1.2817 |

**D**

(1) Relu is better.

(2) RuLu is better.

| Model | 3-node-Relu Loss | 3-node-Sigmoid Loss |
|:---:|:---:|:---:|
| vanilla-sgd | 106.0580 | 324.2820 |
| sgd-w-reg | 106.0685 | 324.3093 |
| sgd-w-momentum | 69.9962 | 128.6403 |
| sgd-w-reg-momentum-drop-out | 69.9962 | 128.6404 |

| Model | 5-node-Relu Loss | 5-node-Sigmoid Loss |
|:---:|:---:|:---:|
| vanilla-sgd | 182.3663 | 247.1844 |
| sgd-w-reg | 182.3803 | 247.2096 |
| sgd-w-momentum | 79.4184 | 91.3771 |
| sgd-w-reg-momentum-drop-out | 79.4184 | 91.3771 |

# Question 6: Activation Functions

We create a python script to simulate different model structures and activation functions, with random data. Since there are numerous plots for each of structure and method (more than a thousand each), therefore I only plot the first 50 different shape and placed them in folder so that you can check. Here I only randomly choose some plots to illustrate.

Function:

$$y = V \cdot ReLu(wx + b) + v_b$$
$$y = V \cdot Sigmoid(wx + b) + v_b$$

## One Layer One Node

## One Layer Two Nodes

### ReLu

Yunjia Zeng, collaborate with Shaoyu Feng, Mengtong Zhang
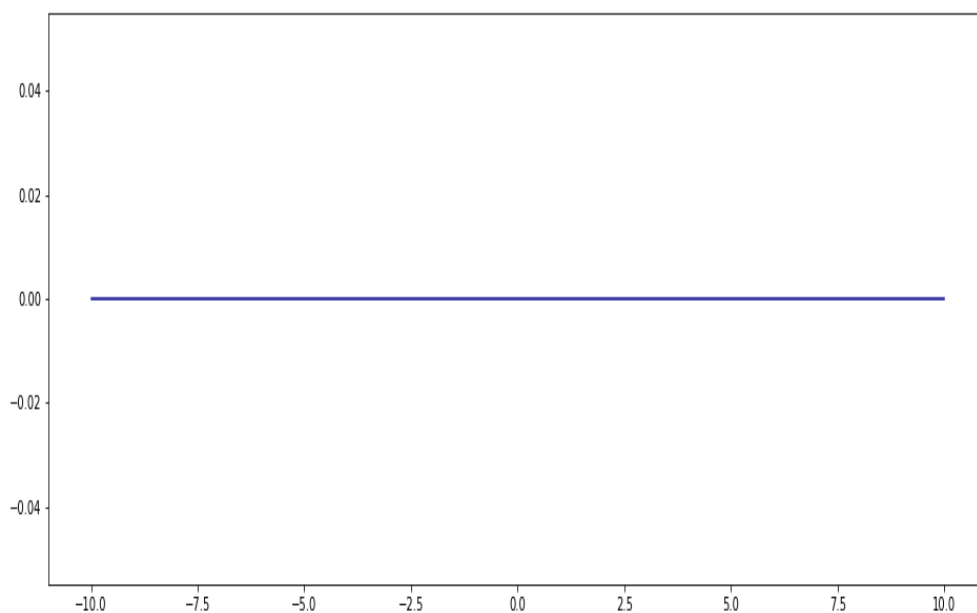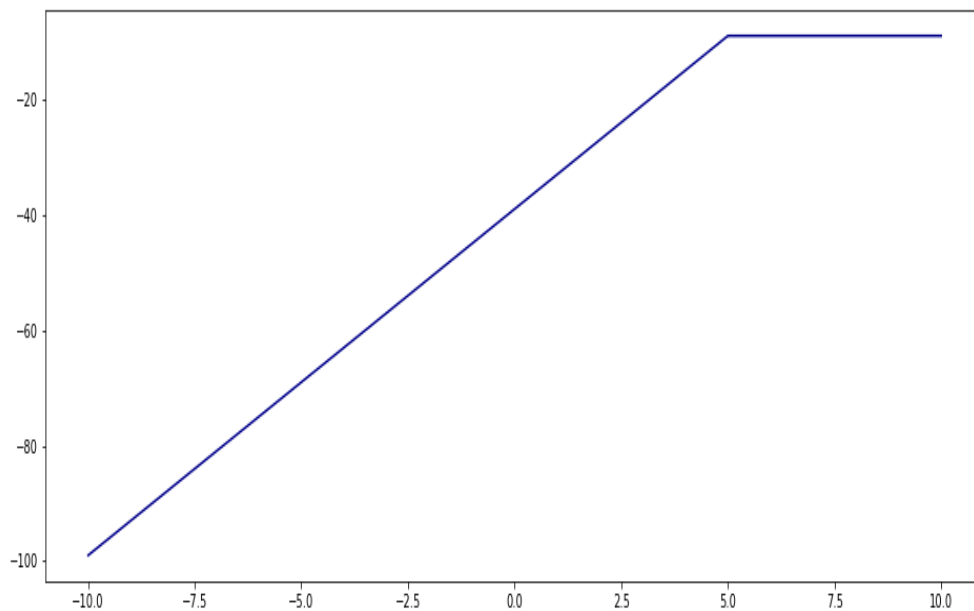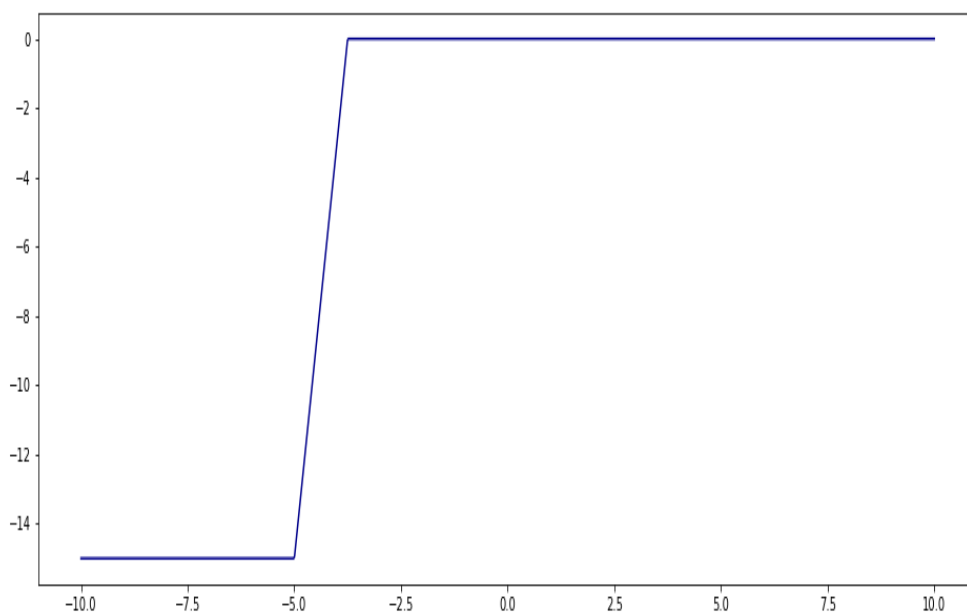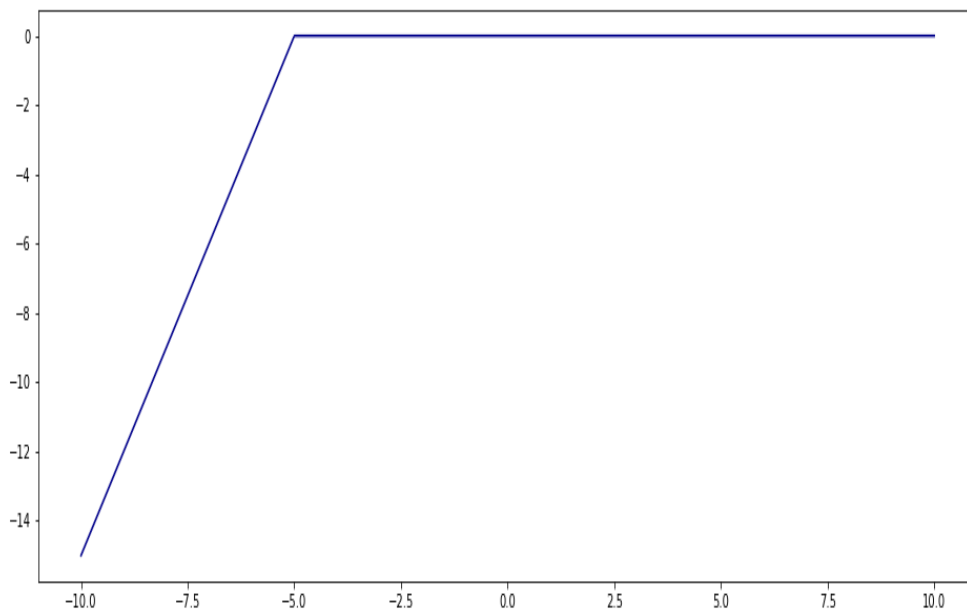yz682            **Assignment 2**            ANLY 601
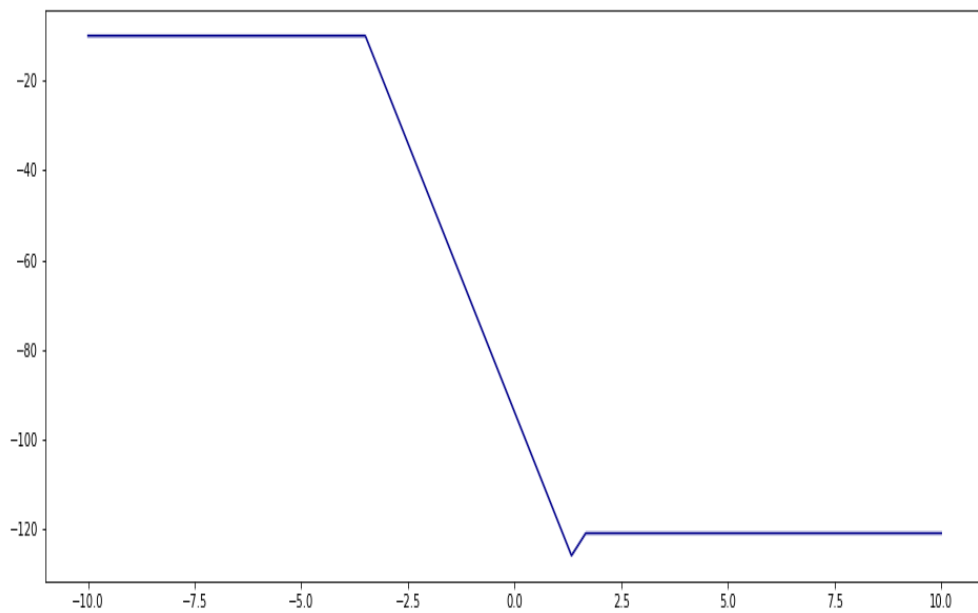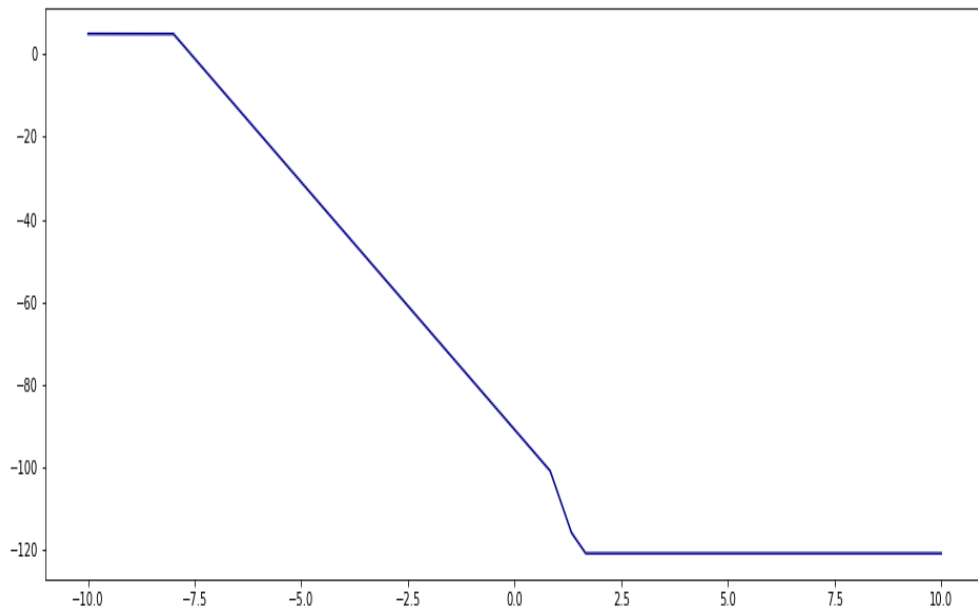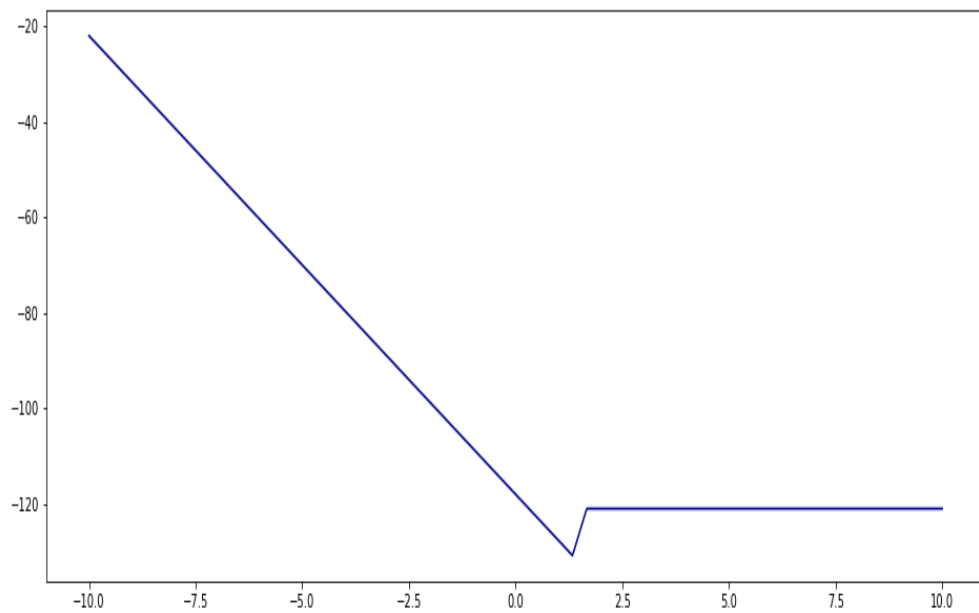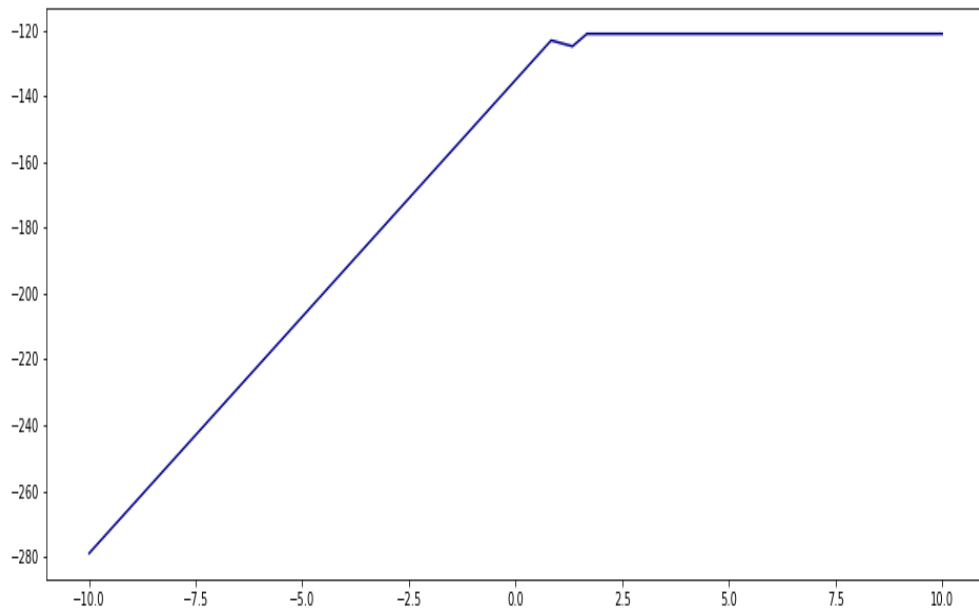April 16, 2020

**Sigmoid**

## Two Layers One Node

### ReLu

**Sigmoid**

## Two Layers Two Nodes

### ReLu

**Sigmoid**