

# CIS581: Assignment 4: Deep Learning Basics

## Due: Dec. 4, 2018 at 11:59 am

### Instructions

- This is an **individual** assignment. “**Individual**” means each student must hand in their **own** answers, and each student must write their **own** code in the homework. It is admissible for students to collaborate in solving problems. To help you actually learn the material, what you write down must be your own work, not copied from any other individual. You must also list the names of students (maximum two) you collaborated with.
- There is no single answer to most problems in deep learning, therefore the questions will often be underspecified. You need to fill in the blanks and submit a solution that solves the (practical) problem. Document the choices (hyperparameters, features, neural network architectures, etc.) you made in the write-up.
- All the code should be written in Python. You should use PyTorch Only to complete this homework.
- Four datasets are needed in this homework, namely, “random dataset”, “line dataset”, and “detection dataset”. They can be downloaded from the course wiki website. Course wiki also has a template code for the first question.
- You must submit your solutions online on **Canvas**. Compress your files into a ZIP file named “1\_<penn\_key>.zip”, which should contain **1 PDF report** and **4 Python files**, each of which contains the Python code for each part. Note that you should include all the figures in your report. There will be one submission for the code (submit your zip file) and one for the report (only the pdf).

## Overview

This homework aims at investigating the basics of neural networks, i.e., activation and loss functions (non-linearity and objectives). As discussed in class, changing activation and loss functions will highly influence the efficiency of learning through back-propagation. And even sometimes, they are coupled together. There are 4 parts of the homework:

1. Build visualization tools to plot activations, losses, and gradients. It is the first step of understanding a function through its energy landscape of outputs and gradients. From the gradient response, it is possible to predict how it will work in a deep network.
2. Experiment with a fully connected network on a random (toy) dataset. In this part, we move one step forward: from theory to experiment. To start with, we use a rather simple network to experiment with different activation and loss functions.
3. Experiment with a convolutional network on two toy datasets, namely, line dataset and detection dataset. In this part, we experiment how the network can learn to interpret image patterns. Also, with little revision, the network can also predict attributes of an object.
4. To understand how the neural network is trained, it is helpful to visualize all of the aspects of the process. Re-produce the experiment that was shown in class. We will visualize the space into which the first layer of the network transforms the data to understand how the non-linear warping helps classification

### 1 Plot Loss and Gradient (20%)

In this part, you will write code to plot the output and gradient for a single neuron with Sigmoid activation and two different loss functions. As shown in Figure 1, You should implement a single neuron with input 1, and calculate different losses and corresponding error.

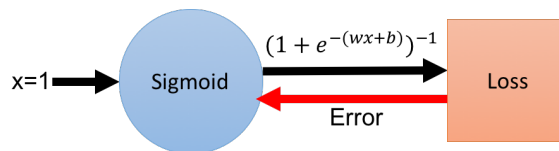


Figure 1: Network diagram for part 1.

All the figures plotted in this part should have the same range of x-axis and y-axis. The range should be centered at 0 but the extend should be picked so as to see the difference clearly.

A set of example plots are provided in Figure 2. Here we use ReLU (instead of Sigmoid) activation and L2 loss as an example.

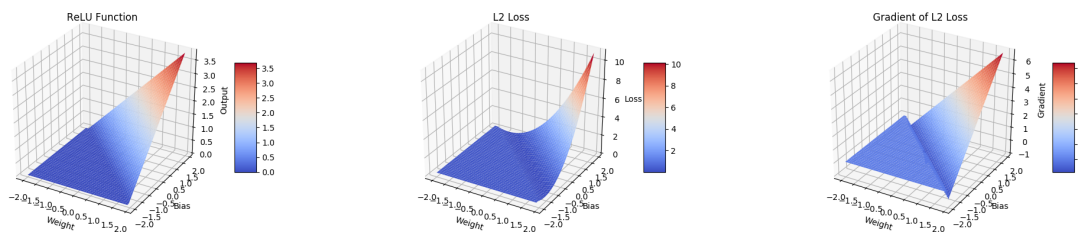


Figure 2: Example plots with ReLU activation and L2 loss. Left: Output of ReLU function. Middle: Loss plot with L2 loss. Right: Gradient plot.

1. (3%) Plot a 3D figure showing the relations of output of Sigmoid function and weight/bias. To be specific, x-axis is weight, y-axis is bias, and z-axis is the output.

Hint: Use the Python package *matplotlib* and the function *plot\_surface* from *mpl\_toolkits.mplot3d* to draw 3D figures.

2. (3%) Experiment with L2 loss. The L2 loss is defined as  $\mathcal{L}_{L2} = (\hat{y} - y)^2$ , where  $y$  is the ground truth and  $\hat{y}$  is the prediction. Let  $y = 0.5$  and plot a 3D figure showing the relations of L2 loss and weight/bias. To be specific, x-axis is weight, y-axis is bias, and z-axis is the L2 loss.
3. (4%) Experiment with back-propagation with L2 loss. Compute  $\frac{\partial \mathcal{L}_{L2}}{\partial \text{weight}}$  and plot 3D figure showing the relations of gradient and weight/bias. To be specific, x-axis is weight, y-axis is bias, and z-axis is the gradient w.r.t. weight.
4. (3%) Experiment with cross-entropy loss. The cross-entropy loss is defined as  $\mathcal{L}_{CE} = -(y \log \hat{y} + (1 - y) \log(1 - \hat{y}))$ , where  $y$  is the ground truth probability and  $\hat{y}$  is the predicted probability. Let  $y = 0.5$  and plot a 3D figure showing the relations of cross-entropy loss and weight/bias. To be specific, x-axis is weight, y-axis is bias, and z-axis is the cross-entropy loss.

5. (4%) Experiment with back-propagation with cross-entropy loss. Compute  $\frac{\partial \mathcal{L}_{CE}}{\partial \text{weight}}$  and plot 3D figure showing the relations of gradient and weight/bias. To be specific, x-axis is weight, y-axis is bias, and z-axis is the gradient w.r.t. weight.
6. (3%) Explain what you observed from the above 5 plots. The explanation should include: 1) What's the difference between cross-entropy loss and L2 loss? 2) What's the difference between the gradients from cross-entropy loss and L2 loss? and 3) Predict how these differences will influence the efficiency of learning.

## 2 Experiment with Fully Connected Network (20%)

Starting from this question, you should use a deep learning framework. In this part, you will extend the theoretical analysis in the first part to quick experiments on a random toy dataset. By doing experiments, you will experience how different activation functions and loss functions influence the learning (or convergence) of a simple neural network. Construct a simple neural network as shown in Figure 3.

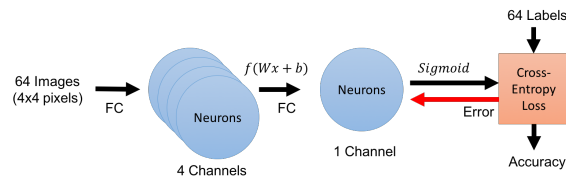


Figure 3: Network diagram for part2

The random toy dataset contains 64 images, each of which is of resolution  $4 \times 4$ . The labels is  $\in \{0, 1\}$ . For all the experiments in this part, use simple gradient descent method with learning rate 0.1. Weights are initialized with truncated normal distribution with mean 0 and standard deviation 0.1. Biases are initialized with constant of 0.1. Batch size is set to 64.

Note: The maximum number of training iterations is 10,000.

1. (4%) Use Sigmoid function as neuron activation and L2 loss for the network. Plot two figures showing 1) loss vs training iterations, and 2) accuracy vs training iterations. Stop the training when accuracy reaches 100%.
2. (4%) Use Sigmoid function as neuron activation and cross-entropy loss for the network. Plot two figures showing 1) loss vs training iterations, and 2) accuracy vs

training iterations. Stop the training when accuracy reaches 100%.

3. (4%) Use ReLU (Rectified Linear Units) function as neuron activation and L2 loss for the network. (Change the activation function of the first fully connected layer to ReLU.) layer Plot two figures showing 1) loss vs training iterations, and 2) accuracy vs training iterations. Stop the training when accuracy reaches 100%.
4. (4%) Use ReLU function as neuron activation and cross-entropy loss for the network. Plot two figures showing 1) loss vs training iterations, and 2) accuracy vs training iterations. Stop the training when accuracy reaches 100%.
5. (4%) Sort the settings according to the training iterations to reach 100% accuracy, and explain the reasons of different convergence rates.

### 3 Solving XOR with a 2-layer Perceptron (20%)

In this question you are asked to build and visualize a 2-layer perceptron that computes the XOR function. The network architecture is shown in Figure 4. The MLP has 1 hidden layer with 2 neurons. The activation function used for the hidden layer is the hyperbolic tangent function. Since we aim to model a boolean function the output of the last layer is passed through a sigmoid activation function to constrain it between 0 and 1.

1. (5%) Formulate the XOR approximation as an optimization problem using the cross entropy loss. *Hint: Your dataset consists of just 4 points,  $\mathbf{x}_1 = (0, 0)$ ,  $\mathbf{x}_2 = (0, 1)$ ,  $\mathbf{x}_3 = (1, 0)$  and  $\mathbf{x}_4 = (1, 1)$  with ground truth labels 0, 1, 1 and 0 respectively.*
2. (10%) Use gradient descent to learn the network weights that optimize the loss. Intuitively, the 2 layer perceptron first performs a nonlinear mapping from  $(x_1, x_2) \rightarrow (h_1, h_2)$  and then learns a linear classifier in the  $(h_1, h_2)$  plane. For different steps during training visualize the image of each input point  $\mathbf{x}_i$  in the  $(h_1, h_2)$  plane as well as the decision boundary (separating line) of the classifier.
3. (5%) What will happen if we don't use an activation function in the hidden layer? Is the network be able to learn the XOR function? Justify your answer.

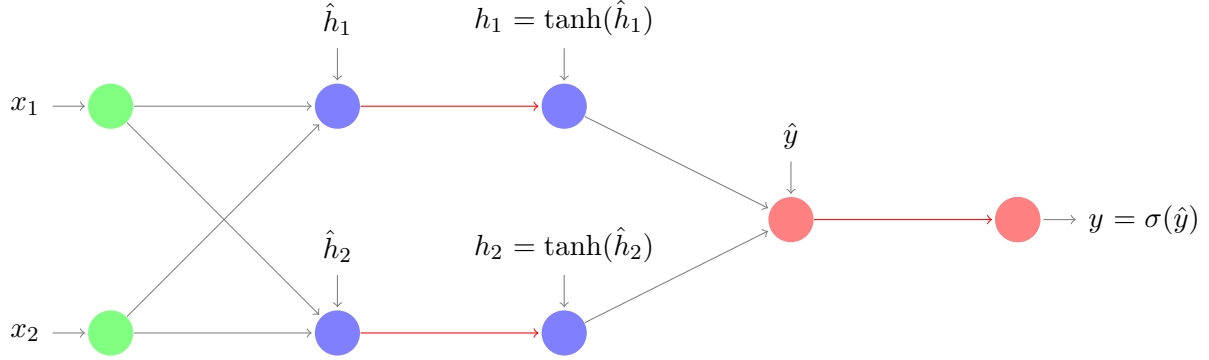


Figure 4: Graphical representation of the 2-layer Perceptron

## 4 Experiment with Convolutional Network (15%)

In this part, you will explore the convolutional networks. By doing experiments, you will experience how convolutional networks learn to interpret images and capture meaningful structure. Construct a convolutional network as shown in Figure 5.

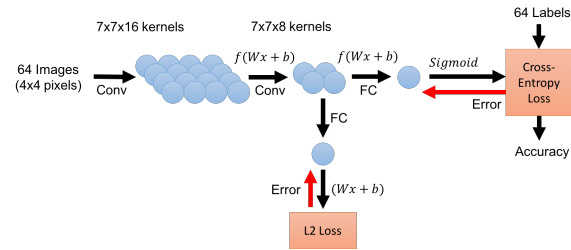


Figure 5: Network diagram for part3

The experiments will be done on two toy datasets (line dataset and detection dataset), each of which contains 64 images of resolution  $8 \times 8$ . The labels is  $\in \{0, 1\}$ . For all the experiments in this part, use simple gradient descent method with learning rate 0.1. Weights are initialized with truncated normal distribution with mean 0 and standard deviation 0.1. Biases are initialized with constant of 0.1. Batch size is set to 64.

Note: The maximum number of training iterations is 10,000.

1. (10%) Experiment with the line dataset (Figure 6). You will construct a network with two convolutional layers and one fully connected layer, concatenated with a cross-entropy loss. (A standard classification network.)

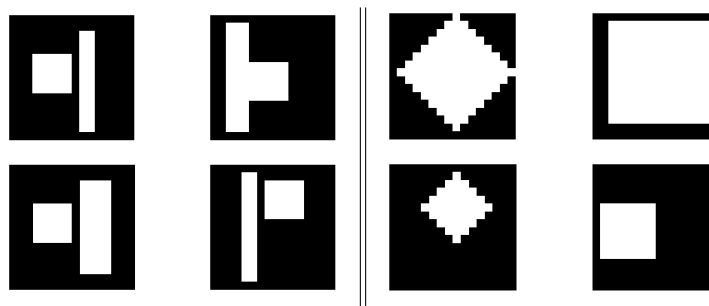


Figure 6: **Left:** 4 sample images from the line dataset. The left column corresponds to label 0, as the right column 1. **Right:** 4 sample images from the detection dataset. The left column corresponds to label 0, as the right column 1. The width of each sample is 7, 6, 4, 3.

Use ReLU as the activation functions of convolutional layers. Plot two figures showing 1) loss vs training iterations, and 2) accuracy vs training iterations. Stop the training when accuracy reaches 100%. Compare the results with part 2.

2. (5%) Experiment with the detection dataset (Figure 6). Use the same network architecture from the previous questions. In addition, add one more fully connected layer on top of the convolutional layer. This fully connected layer (with linear activation) acts as a regressor, followed by an L2 loss. It is used to predict the width of the object. The learning rate for L2 loss is set to 0.001. We define the regression prediction is correct if the predicted width is within 0.5 of the ground truth width.

Use ReLU as the activation functions of convolutional layers. Plot four figures showing 1) cross-entropy loss vs training iterations, 2) classification accuracy vs training iterations, 3) L2 loss vs training iterations, and 4) regression accuracy vs training iterations. Stop the training when regression accuracy reaches 100%.

Compare the results with the previous question.