

# Programming Assignment 2

## CSE 253: Neural Networks for Pattern Recognition

Winter 2019

### Instructions

Due on Saturday, February 2nd, 2019

1. Please submit your assignment on Gradescope. There are two components to this assignment: written homework (Problems 1a-c), and a programming part. You will be writing a report in a conference paper format for the programming part of this assignment, reporting your findings. **The report should be written using L<sup>A</sup>T<sub>E</sub>X or Word in NIPS format** (NIPS is the top machine learning conference, and it is now dominated by deep nets - it will be good practice for you to write in that format!). The templates, both in **Word** and **L<sup>A</sup>T<sub>E</sub>X** are available from the [2015 NIPS format site](#).
2. For the programming part, please work in pairs. In extraordinary circumstances, we will allow you to do it on your own. Please discuss your circumstances with your TA, who will then present your case to me. Again, **don't forget to include a paragraph for each team member in your report that describes what each team member contributed to the project.**
3. You need to submit all of the source codes files and a *readme.txt* file that includes detailed instructions on how to run your code.  
  
You should write clean code with consistent format, as well as explanatory comments, as this code may be reused in the future. Do not submit any of your output plot files or .pyc files, just the .py files and a readme.
4. Using PyTorch, or any off-the-shelf code is strictly prohibited.
5. Any form of copying, plagiarizing, grabbing code from the web, having someone else write your code for you, etc., is cheating. We expect you all to do your own work, and when you are on a team, to pull your weight. Team members who do not contribute will not receive the same scores as those who do. Discussions of course materials and homework solutions are encouraged, but you should write the final solutions to the written part alone. Books, notes, and Internet resources can be consulted, but not copied from. Working together on homework must follow the spirit of the **Gilligan's Island Rule** (Dymond, 1986): No notes can be made (or recording of any kind) during a discussion, and you must watch one hour of Gilligan's Island or something equally insipid before writing anything down. Suspected cheating has been and will be reported to the UCSD Academic Integrity office.

### Multi-layer Neural Networks

In this assignment, we will be classifying handwritten digits from Yann LeCun's MNIST Database. In Assignment 1, we classified the faces dataset using a single-layer neural network with different output activation functions. (Logistic and Softmax regression). In this assignment, we are going to classify the MNIST dataset using multi-layer neural networks with softmax outputs.

## Part I

# Homework problems to be solved individually, and turned in individually

For this part we will not be accepting handwritten reports. Please use latex or word for your report. MathType is a handy tool for equations in Word. The free version (MathType Lite) has everything you need. This should be done individually, and each team member should turn in his or her own work separately.

1. (15pts) **For multiclass classification on the MNIST dataset, we will use the cross-entropy error function and softmax as the output layer.** In our network, we will have a hidden layer between the input and output, that consists of  $J$  units with the tanh activation function. So this network has three layers: an input layer, a hidden layer and a softmax output layer.

*Notation:* We use index  $k$  to represent a node in output layer and index  $j$  to represent a node in hidden layer and index  $i$  to represent a node in the input layer. Additionally, the weight from node  $i$  in the input layer to node  $j$  in the hidden layer is  $w_{ij}$ . Similarly, the weight from node  $j$  in the hidden layer to node  $k$  in the output layer is  $w_{jk}$ .

- (a) (10pts) **Derivation.** In the following discussion,  $n$  denotes the  $n$ th input pattern. Derive the expression for  $\delta$  for both the units of output layer ( $\delta_k^n$ ) and the hidden layer ( $\delta_j^n$ ). Recall that the definition of  $\delta$  is  $\delta_i^n = -\frac{\partial E^n}{\partial a_i^n}$ , where  $a_i^n$  is the weighted sum of the inputs to unit  $i$ . There are two “hard parts” to this: 1) taking the derivative of the softmax; and 2) figuring out how to apply the chain rule to get the hidden deltas. Bishop and Chapter 8 of the PDP books both have good hints on the latter, and Bishop on the former. However, crucial steps have been left out of the Bishop derivation (Chapter 6). Our main hint here is: break it up into two parts (see equation 6.161 in Bishop), when  $k = k'$  and when it doesn't. Note that Bishop (Equation 4.31) defines  $\delta_j^n$  without a minus sign, which is different than we defined it above, and differently than the PDP book chapter 8.
- (b) (2pts) **Update rule.** Write the update rule for  $w$ 's in terms of the  $\delta$ 's you derived above using learning rate  $\alpha$ , starting with the gradient descent rule:

$$w_{ij} = w_{ij} - \alpha \frac{\partial E}{\partial w_{ij}} \quad (1)$$

where

$$\frac{\partial E}{\partial w_{ij}} = \sum_n \frac{\partial E^n}{\partial w_{ij}} \quad (2)$$

You have to write both the update rules, the hidden to output layer ( $w_{jk}$ ) update rule and the input to hidden ( $w_{ij}$ ) update rule in a generalized form. (Hint: you will have to use chain rule for differentiation.)

$$\frac{\partial E^n}{\partial w_{ij}} = \frac{\partial E^n}{\partial a_j^n} \frac{\partial a_j^n}{\partial w_{ij}} \quad (3)$$

- (c) (3pts) **Vectorize computation.** The computation is much faster when you update all  $w_{ij}$ s and  $w_{jk}$ s at the same time, using matrix multiplications rather than **for** loops. Please show the update rule for the weight matrix from the hidden layer to output layer and the matrix from input layer to hidden layer, using matrix/vector notation.

## Part II

# Team Programming Assignment

2. **Classification.** Classification on the MNIST database. Refer to your derivations from Problem 2.

- (a) (0pts) Read in the MNIST data using the “load\_data” function provided in the code.
- (b) (5pts) Check your code for computing the gradient using a small subset of data. You can compute the slope with respect to one weight using the numerical approximation:

$$\frac{d}{dw}E(w) \approx \frac{E(w + \epsilon) - E(w - \epsilon)}{2\epsilon}$$

where  $\epsilon$  is a small constant, e.g.,  $10^{-2}$ . Compare the gradient computed using numerical approximation with the one computed as in backpropagation. The difference of the gradients should be within big-O of  $\epsilon^2$ , so if you used  $10^{-2}$ , your gradients should agree within  $10^{-4}$ . (See section 4.8.4 in Bishop for more details). Note that  $w$  here is *one* weight in the network. You can only check one weight at a time this way - every other weight must stay the same.

Choose one output bias weight, one hidden bias weight, and two hidden to output weights and two input to hidden weights, and show that the gradient obtained for that weight after backpropagation is within ( $O(\epsilon^2)$ ) of the gradient obtained by numerical approximation. For each selected weight  $w$ , first increment the weight by small value  $\epsilon$ , do a forward pass for one training example, and compute the loss. This value is  $E(w + \epsilon)$ . Then reduce  $w$  by the same amount  $\epsilon$ , do a forward pass for the same training example and compute the loss  $E(w - \epsilon)$ . Then compute the gradient using equation mentioned above and compare this with gradient obtained by backpropagation. Report the results in a Table.

- (c) (10pts) Using the vectorized update rule you obtained from 1(c), perform gradient descent to learn a classifier that maps each input data to one of the labels  $t \in \{0, \dots, 9\}$ , using a one-hot encoding. Use 50 hidden units. ***For this programming assignment, use mini-batch stochastic gradient descent throughout, in all problems.***

You should use momentum in your update rule, i.e., include a momentum term weighted by  $\gamma$ , and set  $\gamma$  to 0.9. You should use cross-validation for early stopping of your training: Stop training when the error on the validation set goes up. Use the following criteria - If the validation error goes up for some *threshold* number of epochs, stop training and save the weights which resulted in minimum validation error. The validation set error should go up at some point, although this didn't always happen with the faces.

Describe your training procedure. Plot your training and validation accuracy (i.e., percent correct) vs. number of training epochs, as well as training and validation loss vs. number of training epochs. Report accuracy on test set using the best weights obtained through early stopping.

You may experiment with different learning rates, but you only need to report your results and plots on the best learning rate you find.

- (d) (5pts) **Experiment with Regularization.** Starting with the network you used for part c, with new initial random weights, add weight decay to the update rule. (You will have to decide the amount of regularization, i.e.,  $\lambda$ , a factor multiplied times the weight decay penalty. Experiment with 0.001 and 0.0001) Again, plot training and validation loss, training and validation accuracy, and report final test accuracy. For this problem, train about 10% more epochs than you found in part c (i.e., if you found that 100 epochs were best, train for 110 for this problem). Comment on the change of performance, if any.
- (e) (5pts) **Experiment with Activations.** Starting with the network of part c, try using different activation functions for the hidden units. You are already using tanh, try the other two below. Note that the derivative changes when the activation rule changes!!

- i. Sigmoid.  $f(z) = \frac{1}{1+e^{-z}}$

- ii. ReLU.  $f(z) = \max(0, z)$

The weight update rule is exactly the same for each activation function. The only thing that changes is the derivative of the activation function when computing the hidden unit  $\delta$ s. For each activation function you try, plot training and validation loss on one graph, training and validation accuracy on another, and report final test accuracy. Comment on the change of performance.

- (f) (5pts) **Experiment with Network Topology.** Starting with the network from part c, consider how the topology of the neural network changes the performance.
  - i. Try halving and doubling the number of hidden units. Plot training and validation loss, training and validation accuracy, and report final test accuracy. How does performance change? Explain your results.
  - ii. Change the number of hidden layers. Use two hidden layers instead of one. Create a new architecture that uses two hidden layers of equal size and has approximately the same number of parameters, as the previous network with one hidden layer of 50 units. By that, we mean it should have roughly the same total number of weights and biases. Again, plot training and validation loss, training and validation accuracy, and report final test accuracy. How did the performance change?

## Instructions for Programming Assignment

The MNIST dataset has been randomly shuffled, split into training, validation and testing data and uploaded on Resources page and the github repo in pickle format(zipped).

You need to edit the **neuralnet\_starter.py** file to complete the assignment. This file is a skeleton code that is designed to guide you to build and implement your neural net in an efficient and modular fashion, and this will give you a feel for what developing models in PyTorch will be like.

A **config** dictionary is provided which has all the information necessary to build the model. The purpose of each flag is indicated in the comment next to it. Use this dictionary to decide the architecture, activation functions, etc. **Please do not add additional keys to it.**

The class **Activation** includes the definitions for all activation functions and their gradients, which you need to fill in. The definitions of 'forward\_pass' and 'backward\_pass' have been implemented for you in this class. The code is structured in such a way that each activation function is treated as an additional layer on top of a linear layer that computes the net input ( $a$ ) to the unit. To add an activation layer after a fully-connected or linear layer, a new object of this class needs to be instantiated and added to the model.

The **Layer** class denotes a standard fully-connected/ linear layer. The 'forward\_pass' and 'backward\_pass' functions need to be implemented by you. As the name suggests, 'forward\_pass' takes in an input vector 'x' and outputs the variable 'a'. Do not apply the activation function on the computed weighted sum of inputs since the activation function is implemented as a separate layer, as mentioned above. The function 'backward\_pass' takes the weighted sum of the deltas from the layer above it as input, computes the gradient for its weights (to be saved in 'd\_w') and biases (to be saved in 'd\_b'). If there is another layer below that (multiple hidden layers), it also passes the weighted sum of the deltas back to the previous layer. Otherwise, if the previous layer is the input layer, it stops there.

The **Neuralnetwork** class defines the entire network. The '\_\_init\_\_' function has been implemented for you which uses the 'config' specifications to generate the network. Make sure to understand this function very carefully since good understanding of this will be needed while implementing 'forward\_pass' and 'backward\_pass' for this class. The function 'forward\_pass' takes in the input dataset 'x' and targets (in one hot encoded form) as input, performs a forward pass on the data 'x' and returns the loss and predictions. The 'backward\_pass' function computes the error signal from saved predictions and targets and performs a backward pass through all the layers by calling backward pass for each layer of the network, until it reaches the first hidden layer above the input layer (usually there will only be one hidden layer for this project, but when there are more, there will be more backward passes). The 'loss\_func' function computes cross-entropy loss by taking in the logits (a fancy term for prediction  $y$ ) and targets and returns this loss.

Additionally, you need to implement the **softmax**, **load\_data**, **trainer** and **test** functions. The requirements for these functions and all other functions are given in the code.

Furthermore, a couple of things to take care of:

- **Do not** add additional keys in **config** dictionary.
- **Do not** modify the main function in the code.
- **Do not** modify the **checker.py** file. This is the file that has test cases for your code. You can download it and use it to verify your results to ensure your implementation is correct.
- You are allowed to write additional functions if you feel the need to do so.
- As such, the code is solvable using numpy and pickle libraries only. However, if you feel the need to use additional libraries, you can do so as long as they don't have implemented functions for backprop, etc., and you mention these dependencies in the Readme file. That said, make sure to include clear instructions in the Readme file to run your code. If you haven't done so and if we are not able to run your code using the instructions you provide, you lose points.