

Homework #3

CSEP 546: Machine Learning

Prof. Byron Boots

Due: **Wednesday** Dec 1, 2021 11:59pm Pacific Time

Please review all homework guidance posted on the website before submitting to GradeScope. Reminders:

- Please provide succinct answers and supporting reasoning for each question. Similarly, when discussing experimental results, concisely create tables and/or figures when appropriate to organize the experimental results. All explanations, tables, and figures for any particular part of a question must be grouped together.
- Please typeset your submission in a pdf file. Failure to do so may result in a points deduction.
- For every problem involving generating plots, please include the plots as part of your PDF submission.
- When submitting to Gradescope, please link each question from the homework in Gradescope to the location of its answer in your homework PDF. Failure to do so may result in deductions of up to *[5 points]*. For instructions, see https://www.gradescope.com/get_started#student-submission.
- If you collaborate on this homework with others, you must indicate who you worked with on your homework. Failure to do so may result in accusations of plagiarism.
- For every problem involving code, please include the code as part of your PDF for the PDF submission *in addition to* submitting your code to the separate assignment on Gradescope created for code. Not submitting all code files will lead to a deduction of *[1 point]*.
- Please indicate your final answer to each question by placing a box around the main result(s). To do this in \LaTeX , one option is using the `boxed` command.

Not adhering to these reminders may result in point deductions.

Conceptual Questions

1. The answers to these questions should be answerable without referring to external materials. Briefly justify your answers with a few words.
 - a. [2 points] Consider training kernel ridge regression with a Gaussian RBF kernel ($K(u, v) = \exp\left(-\frac{\|u-v\|_2^2}{2\sigma^2}\right)$). It seems to underfit the training set: should you increase or decrease σ ?
 - b. [2 points] True or False: Training deep neural networks requires minimizing a non-convex loss function, and therefore gradient descent might not reach the globally-optimal solution.
 - c. [2 points] True or False: It is a good practice to initialize all weights to zero when training a deep neural network.
 - d. [2 points] True or False: We use non-linear activation functions in a neural network's hidden layers so that the network learns non-linear decision boundaries.
 - e. [2 points] True or False: Given a neural network, the time complexity of the backward pass step in the backpropagation algorithm can be prohibitively larger compared to the relatively low time complexity of the forward pass step.

What to Submit:

- **Writeup:** For each part a-e, 1-2 sentences containing your answer.

Coding

Introduction to PyTorch

2. PyTorch is a great tool for developing, deploying and researching neural networks and other gradient-based algorithms. In this problem we will explore how this package is built and re-implement some of its core components. First start by reading `README.md` file provided in `intro_pytorch` subfolder. A lot of problem statements will overlap between here, readmes and comments in functions.
 - a. [10 points] You will start by implementing components of our own PyTorch modules. You can find these in folders: `layers`, `losses` and `optimizers`. Almost each file there should contain at least one problem function, including exact directions for what to achieve in this problem. Finally, you should implement functions in `train.py` file.
 - b. [5 points] Next we will use the above module to perform hyperparameter search. Here we will also treat loss function as a hyper-parameter. However, because cross-entropy and MSE are different, we are going to use two different files: `crossentropy_search.py` and `mean_squared_error_search.py`. For each you will need to build and train (in the provided order) 5 models:
 - Linear neural network (Single layer, no activation function)
 - NN with one hidden layer (2 units) and sigmoid activation function after the hidden layer
 - NN with one hidden layer (2 units) and ReLU activation function after the hidden layer
 - NN with two hidden layer (each with 2 units) and Sigmoid, ReLU activation functions after first and second hidden layers, respectively
 - NN with two hidden layer (each with 2 units) and ReLU, Sigmoid activation functions after first and second hidden layers, respectively

For each loss function, submit a plot of losses from training and validation sets. All models should be on the same plot (10 lines per plot), with two plots total (1 for MSE, 1 for cross-entropy).

- c. [5 points] For each loss function, report the best performing architecture (best performing is defined here as achieving the lowest validation loss at any point during the training), and plot it's guesses on test set. You should use function `plot_model_guesses` from `train.py` file. Finally, report accuracy of that model on a test set.
- d. [3 points] Is there a big gap in performance between between MSE and Cross-Entropy models? If so, explain why it occurred? If not explain why different loss functions achieve similar performance? Answer in 2-4 sentences.

What to Submit:

- **Part b:** 2 plots (one per loss function), with 10 lines each, showing both training and validation loss of each model. Make sure plots are titled, and have proper legends.
- **Part c:** Names of best performing models (i.e. descriptions of their architectures), and their accuracy on test set.
- **Part c:** 2 scatter plots (one per loss function), with predictions of best performing models on test set.
- **Part d:** 2-4 sentence written reponse to provided questions.
- **Code** on Gradescope through coding submission

Resources

For the next question you will use a lot of PyTorch. Please feel free to reference [PyTorch Documentation](#), when needed.

If you do not have access to GPU, you might find [Google Colaboratory](#) useful. It allows you to use a cloud GPU for free. To enable it make sure: “Runtime” → “Change runtime type” → “Hardware accelerator” is set to “GPU”. When submitting please download and submit a `.py` version of your notebook.

Neural Networks for MNIST

3. In previous homeworks, we used ridge regression for training a classifier for the MNIST data set. Similarly in previous homework, we used logistic regression to distinguish between the digits 2 and 7.

In this problem, we will use PyTorch to build a simple neural network classifier for MNIST to further improve our accuracy.

We will implement two different architectures: a shallow but wide network, and a narrow but deeper network. For both architectures, we use d to refer to the number of input features (in MNIST, $d = 28^2 = 784$), h_i to refer to the dimension of the i -th hidden layer and k for the number of target classes (in MNIST, $k = 10$). For the non-linear activation, use ReLU. Recall from lecture that

$$\text{ReLU}(x) = \begin{cases} x, & x \geq 0 \\ 0, & x < 0. \end{cases}$$

Weight Initialization

Consider a weight matrix $W \in \mathbb{R}^{n \times m}$ and $b \in \mathbb{R}^n$. Note that here m refers to the input dimension and n to the output dimension of the transformation $x \mapsto Wx + b$. Define $\alpha = \frac{1}{\sqrt{m}}$. Initialize all your weight matrices and biases according to $\text{Unif}(-\alpha, \alpha)$.

Training

For this assignment, use the Adam optimizer from `torch.optim`. Adam is a more advanced form of gradient descent that combines momentum and learning rate scaling. It often converges faster than regular gradient descent in practice. You can use either Gradient Descent or any form of Stochastic Gradient Descent. Note that you are still using Adam, but might pass either the full data, a single datapoint or a batch of data to it. Use cross entropy for the loss function and ReLU for the non-linearity.

Implementing the Neural Networks

- a. [10 points] Let $W_0 \in \mathbb{R}^{h \times d}$, $b_0 \in \mathbb{R}^h$, $W_1 \in \mathbb{R}^{k \times h}$, $b_1 \in \mathbb{R}^k$ and $\sigma(z): \mathbb{R} \rightarrow \mathbb{R}$ some non-linear activation function applied element-wise. Given some $x \in \mathbb{R}^d$, the forward pass of the wide, shallow network can be formulated as:

$$\mathcal{F}_1(x) := W_1 \sigma(W_0 x + b_0) + b_1$$

Use $h = 64$ for the number of hidden units and choose an appropriate learning rate. Train the network until it reaches 99% accuracy on the training data and provide a training plot (loss vs. epoch). Finally evaluate the model on the test data and report both the accuracy and the loss.

- b. [10 points] Let $W_0 \in \mathbb{R}^{h_0 \times d}$, $b_0 \in \mathbb{R}^{h_0}$, $W_1 \in \mathbb{R}^{h_1 \times h_0}$, $b_1 \in \mathbb{R}^{h_1}$, $W_2 \in \mathbb{R}^{k \times h_1}$, $b_2 \in \mathbb{R}^k$ and $\sigma(z): \mathbb{R} \rightarrow \mathbb{R}$ some non-linear activation function. Given some $x \in \mathbb{R}^d$, the forward pass of the network can be formulated as:

$$\mathcal{F}_2(x) := W_2 \sigma(W_1 \sigma(W_0 x + b_0) + b_1) + b_2$$

Use $h_0 = h_1 = 32$ and perform the same steps as in part a.

- c. [5 points] Compute the total number of parameters of each network and report them. Then compare the number of parameters as well as the test accuracies the networks achieved. Is one of the approaches (wide, shallow vs. narrow, deeper) better than the other? Give an intuition for why or why not.

Using PyTorch: For your solution, you may not use any functionality from the `torch.nn` module except for `torch.nn.functional.relu`, `torch.nn.functional.cross_entropy`, `torch.nn.parameter.Parameter` and `torch.nn.Module`. You must implement the networks \mathcal{F}_1 and \mathcal{F}_2 from scratch.

What to Submit:

- **Parts a-b:** Provide a plot of the training loss versus epoch. In addition evaluate the model trained on the test data and report the accuracy and loss.
- **Part c:** Report the number of parameters for the network trained in part (a) and for the network trained in part (b). Provide a comparison of the two networks as described in part in 1-2 sentences.
- **Code** on Gradescope through coding submission.

Administrative

[2 points] About how many hours did you spend on this homework? There is no right or wrong answer :)