

UNIVERSITY OF PENNSYLVANIA
ESE 546: PRINCIPLES OF DEEP LEARNING
FALL 2024
[11/04] HOMEWORK 4
DUE: 11/25 MON 11:59P ET

Changelog. **None**

Read the following instructions carefully before beginning to work on the homework.

- You will submit solutions typeset in \LaTeX on Gradescope (strongly encouraged). You can use `hw_template.tex` on Canvas in the “Homeworks” folder to do so. If your handwriting is unambiguously legible, you can submit PDF scans/tablet-created PDFs.
- Clearly indicate the name and Penn email ID of all your collaborators on your submitted solutions.
- Start a new problem on a fresh page and mark all the pages corresponding to each problem. Failure to do so may result in your work not graded completely.
- For each problem in the homework, you should mention the total amount of time you spent on it. This helps us keep track of which problems most students are finding difficult.
- You can be informal while typesetting the solutions, e.g., if you want to draw a picture feel free to draw it on paper clearly, click a picture and include it in your solution. Do not spend undue time on typesetting solutions.
- You will see an entry of the form “HW 4 PDF” where you will upload the PDF of your solutions. You will also see entries like “HW 4 Problem 2 Code” and “HW 4 Problem 3 Code” where you will upload your solution for the respective problems.
- **For each programming problem/sub-problem, you should create a fresh .py file.** This file should contain **all** the code to reproduce the results of the problem/sub-problem, e.g., it should save the plot that is required (correctly with all the axes, title and legend) as a PDF in the same directory. You will upload the .py file as your solution for “HW 4 Problem 2 Code”. Name your file as `pennkey_hw4_problem2.py`, e.g., I will name my code as `pratikac_hw4_problem2.py`. Note, we will not accept .ipynb files (i.e., Jupyter notebooks), you should only upload .py files. If you are using Google Colab to do your homework (and I suggest that you don’t...), you can export the notebook to a .py file.
- **This is very important.** Note that the instructors will download your code and execute it themselves, so your code should be such that it can be executed independently without any errors to create all output/plots required in the problem.

- **This is is even more important.** If you do not include all the plots in the PDF, you will not get credit for the corresponding questions.

Credit The points for the problems add up to 130. You need to solve for 100 points to get full credit, i.e., your final score will be $\min(\text{your total points}, 100)$.

1 **Problem 1 (30 points).** In this problem you will show that co-coercivity of the gradient and its
 2 Lipschitz continuity are equivalent. Assume that the function $f(x)$ is convex.

3 (a) (10 points) For a differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, show that co-coercivity of ∇f implies
 4 Lipschitz continuity of ∇f . In other words, show that for all x, y

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2 \implies \|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|.$$

5 Hint: use the Cauchy-Schwartz inequality.

6 (b) (10 points) Show the converse, i.e., show that Lipschitz continuity of ∇f implies co-coercivity.

7 (c) (10 points) For a twice-differentiable function $f(x)$ with L -Lipschitz gradients and strong-convexity
 8 parameter m , show that

$$m \leq \|\nabla^2 f(x)\|_2 \leq L$$

9 for all x . Hint: use the mean-value theorem on a line that joins two points x, y .

10 **Problem 2 (40 points). (Do this on your laptop)** In this problem, we will implement logistic
 11 regression for classifying two classes (zero and one) from MNIST. You may not use any routines from
 12 PyTorch other than the ones that help download the data. Use Numpy to code up the optimization
 13 algorithm.

14 (a) (0 points) First, prepare the dataset. Select all the samples belonging to the first two classes from
 15 MNIST's training dataset, this will be your training dataset. Similarly, create a validation dataset for
 16 the two classes by picking samples corresponding to the first two classes from the validation set of the
 17 MNIST dataset. You can subsample input images from 28×28 to 14×14 if you need.

18 (b) (10 points) Logistic regression solves for

$$\underset{w \in \mathbb{R}^d, w_0 \in \mathbb{R}}{\operatorname{argmin}} \quad \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y_i (w^\top x_i + w_0)} \right) + \frac{\lambda}{2} \left(\|w\|_2^2 + w_0^2 \right) \quad (1)$$

19 where $x \in \mathbb{R}^{196}$. Set $y_i = 1$ for MNIST images labeled zero and $y_i = -1$ for MNIST images labeled
 20 one. Initialize the weights randomly for both the following parts but make sure that they are the same
 21 for both gradient descent and gradient descent with Nesterov's acceleration in part (c). You can try a
 22 few different values of λ and pick the one that gives the best validation error for the following parts.

23 Optimize the objective in Eq. (1) using gradient descent (note, not stochastic gradient descent) and
 24 plot the training *loss* as a function of the number of parameter updates on a semi-log scale (log scale
 25 on the Y-axis). This plot should be a straight line. As we saw in the class, the slope of this line should
 26 be about $-\kappa^{-1}$ for gradient descent. Compute the slope of the line in your plot and mention it clearly.

27 (c) (5 points) Write down the Hessian of the loss function in Eq. (1). Without assuming any special
 28 conditions about the dataset $\{(x_i, y_i)\}_{i=1, \dots, n}$, is this problem strongly convex? What is the best
 29 strong convexity parameter for the loss function in Eq. (1)?

30 (d) (10 points) Optimize the objective in Eq. (1) again, this time using gradient descent with Nesterov's
 31 acceleration. If we knew the condition number κ , what momentum coefficient would we use for this
 32 problem? It is difficult, although possible, to evaluate κ , so let's use

$$\kappa = \frac{L}{m}$$

where we treat L as a hyper-parameter, i.e., we *choose* a value for L and set m to be the solution of part (b). Again, choose the value of L by trying out a few values and looking at the slope of the training loss for each setting. (Hint: the momentum parameter is typically between 0.75-0.95).

The slope of the semi-log plot of training loss versus the number of parameter updates in this case will be about $-\kappa^{-1/2}$ for Nesterov's updates. Note that we do not know the correct κ in this problem, so your slope may not match the value you chose for κ above. You should however see that the slope of the plot for Nesterov's acceleration is better than the slope of the plot you obtained in part (b).

(e) **(15 points)** We will now optimize the same problem with stochastic gradient descent (SGD). Use a batch-size $b = 128$ (feel free to try different batch-size such as 64 and 8) and optimize Eq. (1) using SGD with and without Nesterov's acceleration. Plot the training loss against the number of parameter updates on a semi-log scale (log scale on the Y-axis) for (i) gradient descent with Nesterov's updates (can be the same plot from your previous solutions), (ii) SGD without Nesterov's acceleration and, (iii) SGD with Nesterov's acceleration. Is the convergence for (iii) faster than that of (ii)? Comment on the differences for the convergence curves of (i) and (ii).

Problem 3 (40 points). (You need Colab/AWS for this problem) In this problem, we will see a heuristic that helps find out a good schedule for the learning rate. You will use the All-CNN network from Homework 2 in

<https://gist.github.com/pratikac/68d6d94e4739786798e90691fb1a581b> and train it on the CIFAR-10 dataset. Here is one learning rate schedule that has been noticed to work in practice

$$\eta(t) = \begin{cases} 10^{-4} + \frac{t}{T_0} \eta_{\max} & \text{if } t \leq T_0 \\ \eta_{\max} \cos\left(\frac{\pi}{2} \frac{t-T_0}{T-T_0}\right) + 10^{-6} & \text{if } T_0 \leq t \leq T. \end{cases} \quad (2)$$

Here t is the number of weight updates of the network. This is called the cosine learning rate schedule with a warmup. The first phase until $t \leq T_0$ is called the warmup phase and the second phase where the learning rate decreases along a cosine is called the annealing phase. The constants 10^{-4} and 10^{-6} at the beginning and the end are your choices; they do not typically matter much in practice.

There are three parameters in Eq. (2): the length of warmup T_0 , the maximum learning rate η_{\max} and the total number of weight updates T . Let us eliminate one parameter and set

$$T_0 = \frac{T}{5}.$$

For a mini-batch size of 128 in CIFAR-10, there are 391 weight updates in every training epoch. We will train for 50 epochs for this problem which gives

$$T = 50 \times 391 = 19,550.$$

The only parameter left to choose is η_{\max} .

(a) **(0 points)** Why does the cross-entropy loss for CIFAR-10 start at around 2.3 for the All-CNN network with the default PyTorch initialization?

(b) **(10 points)** We saw in the lecture that for gradient descent to converge, in particular to make monotonic progress, the learning rate is limited by the smoothness coefficient of the loss function. It is difficult to estimate the Lipschitz constant of the gradients of a typical deep network and therefore to pick a good learning rate: pick it too small and you train too slowly, pick it to be too large and the

67 loss does not go down quickly enough. We will use exactly this property to pick η_{\max} . Here is the
68 algorithm, often called the “learning-rate finder”.

- 69 (1) Fix all hyper-parameters, e.g., dropout probability, weight-decay and momentum to some
70 reasonable values. Say, pick weight-decay to be 10^{-3} , momentum coefficient to be 0.9 and
71 dropout to 0.5.
72 (2) Start from $\eta(0) = 10^{-5}$. For the mini-batch of the training set at time t , set the learning rate
73 to

$$\eta(t+1) = 1.1 \eta(t);$$

74 you can also pick some other constant greater than 1 (typically, something around 1.1 works
75 well) to increase the learning rate. Record the average training loss of each mini-batch
76 separately and the learning rate $\eta(t)$ that was used for it for about 100 iterations. Plot the
77 training loss (Y-axis) as a function of the learning rate (X-axis); use a log-scale for the X-axis.
78 You should see $\ell(w^{(t)})$ start off at high value, decrease rapidly with the increasing learning
79 rate and then start increasing again after a certain value. The local minimum of this plot is
80 the maximum learning rate one can use for this particular network, with this dataset and
81 hyper-parameters.

82 Run the above algorithm and draw the plot in part (ii). Identify the minimum of this curve and denote
83 the corresponding learning rate η^* . You do not need to be very precise in identifying the minimum, so
84 long as you can find it approximately simply by looking at the curve, you can proceed to the next step.

85 (c) **(10 points)** We cannot use the learning rate η^* directly to choose a value of η_{\max} in Eq. (2). This
86 is because the network had already trained for a few time-steps before we tried out η^* in part (a.ii).
87 We should be more conservative in picking the maximum learning rate and will set

$$\eta_{\max} = \frac{\eta^*}{10}.$$

88 Train the network with this value of η_{\max} for 100 epochs as we had calculated above. Take care to
89 ensure that you change the learning rate according to the schedule in Eq. (2) *before each weight update*.
90 Note that you *cannot* simply create a new PyTorch optimizer because it will reset the momentum
91 buffer. You will have to do something like

```
92 for g in optimizer.param_groups:  
93     g['lr'] = lr  
94
```

96 before feeding each mini-batch. Plot the learning rate, training and validation loss and error as a
97 function of the number of weight updates. You should compute the validation loss and error after
98 every epoch.

99 (d) **(20 points) (This part will require some time to train)** There are heuristics in the optimization
100 process that help pick hyper-parameters in practice. We picked the momentum parameter to be
101 $\rho = 0.9$ in the previous two parts. We will now show experimentally that if

$$\frac{\eta_{\max}}{1 - \rho}$$

102 is kept unchanged, the validation error more or less remains the same. You will train the network for
103 50 epochs and measure the validation error at the end of the 50 epochs for three settings:

- 104 (i) η_{\max} and $\rho = 0.9$
 105 (ii) $\eta_{\max} \leftarrow 5\eta_{\max}$ and $\rho = 0.5$
 106 (iii) $\eta_{\max} \leftarrow \eta_{\max}$ and $\rho = 0.5$

107 You will notice that the validation error is about the same for the first two but increases for the third
 108 setting.

109 **Note:** There are more such heuristics, e.g.,

$$\frac{\eta \lambda}{\ell(1 - \rho)};$$

110 is called the effective learning rate in the hyper-parameter optimization literature and often determines
 111 the validation error if everything else remaining the same. Here λ is the coefficient of weight-decay, ρ
 112 is the momentum parameter as above and ℓ is the batch-size.

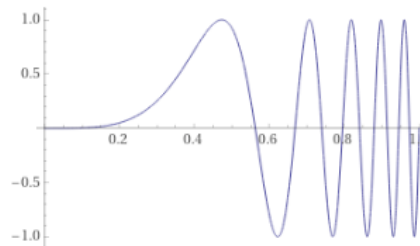
113 Such heuristics are incredibly useful in practice, for instance if one were to search for the best value
 114 for these 4 hyper-parameters using a grid-search, even with only 3 candidate values for each of
 115 them, one would need to train the network 81 times. You can perform a bisection search with the
 116 knowledge of this invariant to significantly reduce this cost. Note that these heuristics work well for
 117 image-classification problems, they do not work so well for recurrent networks.

118 **Problem 4 (20 points. Do this problem on your laptop; you do not need a GPU for this problem.).**

119 In this problem, we will use nonlinear regression to understand one possible reason why deep networks
 120 can avoid overfitting even when the number of samples in the dataset is smaller than the number of
 121 weights.

122 (a) **(0 points)** You will create the dataset yourself. Let the true data come from

$$\mathbb{R} \ni y = f^*(x) \equiv \sin(10\pi x^4); \text{ for } x \in [0, 1].$$



123

124 You will write a function to sample n inputs $x_i \in [0, 1]$ and their corresponding outputs $y_i =$
 125 $\sin(10\pi x_i^4)$ to create a dataset $D_n = \{(x_i, y_i)\}_{i=1}^n$.

126 (b) **(10 points)** Build a multi-layer perceptron (MLP) in PyTorch to regress the dataset. You can use
 127 an MLP with 2-3 layers, ReLU nonlinearities with batch-normalization after each layer and train it
 128 with SGD. You can use any PyTorch function for this that you would like; but it is actually easier to
 129 code the mini-batch sampling etc. yourself simply using a for loop for the iterations instead of using a
 130 dataloader. You should choose the hyper-parameters (weight-decay, learning rate, number of epochs)
 131 to ensure that you are getting zero or very close to zero training error. Denote the learned function as

$$\hat{y} = f_w(x; n)$$

132 where w are the weights of the MLP, x is the test input and we have denoted by n the fact that this
 133 model was trained with n data points. Make sure your code for training the MLP is efficient, we will
 134 be fitting about 100 different models. Solving this problem will also convince you that your laptops
 135 are plenty powerful to ask serious questions about why deep networks work; you do not need GPUs
 136 for everything in deep learning.

137 Write a function to evaluate the quantity

$$\delta f_{\text{in}}(n) = \max_{x \in [0, 1]} |f_w(x; n) - f^*(x)|$$

$$\delta f_{\text{out}}(n) = \max_{x \in [0, 1.5]} |f_w(x; n) - f^*(x)|.$$

138 You will evaluate the maximum in the above expression by sampling many (say 1000) test inputs
 139 $x \in [0, 1]$ or $x \in [0, 1.5]$ and taking the maximum of their predictions. The quantity $\delta f_{\text{in}}(n)$ is the
 140 largest discrepancy within the support of the dataset, i.e., within $[0, 1]$, between the true function
 141 $f^*(x)$ and our fitted function $f_w(x; n)$ using n training data points. Similarly the quantity $\delta f_{\text{out}}(n)$ is
 142 the largest discrepancy outside the support, for example in $x \in [0, 1.5]$.

143 (c) **(10 points)** For each of 20 different values of n , e.g., $\text{ns} = \text{np.logspace}(1, 3, 20).astype(\text{int})$, create
 144 5 training datasets (i.e., you will create 100 different datasets in total) like it is described in part
 145 (a), fit an MLP to this dataset and evaluate the mean and standard deviation (across the 5 datasets
 146 corresponding to each n) of both $\delta f_{\text{in}}(n)$ and $\delta f_{\text{out}}(n)$. Draw a plot of $\log \delta f_{\text{in}}(n)$ (mean and standard
 147 deviation) versus n . On the same plot also plot $\log \delta f_{\text{out}}(n)$ (mean and standard deviation) versus n .
 148 Interpret your plot.