

UNIVERSITY OF PENNSYLVANIA

ESE 546: PRINCIPLES OF DEEP LEARNING

FALL 2020

[11/03] HOMEWORK 4

DUE: 11/18 WED 1.30P ET

---

## Changelog

---

### Instructions

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

- You will submit solutions typeset in L<sup>A</sup>T<sub>E</sub>X 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.
- Please 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.
- Clearly indicate the name and Penn email ID of all your collaborators on your submitted solutions.
- For each problem in the homework, you should mention the total amount of time you spent on it.
- 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 1 Code” and “HW 4 Problem 2 Code” where you will upload your solution for the respective problems. **For each programming problem, you should create a fresh Google Colab notebook.** This notebook should contain all the code to reproduce the results of the problem and you will upload the .ipynb file obtained from Colab. Name your notebook to be “pennkey\_hw3\_problem4.ipynb”, e.g., I will name my code for Problem 4 as “pratikac\_hw3\_problem4.ipynb”.
- **Remember that you should include all the relevant plots in the PDF, without doing so you will not get full credit. Your PDF solutions should be completely self-contained, we should not have have to look at/run the Colab notebook to check your solutions.**

**Credit** The points for the problems add up to 120. You only 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 (60 points). (Can be done on your laptop)** In this problem, we will implement logistic  
2   regression for classifying two classes (zero and one) from MNIST. You may not use any routines  
3   from PyTorch other than the ones that help download the data. Use Numpy to code up the optimizers.

4   (a) **(0 points)** First, prepare the dataset. Select all the samples belonging to the first two classes from  
5   MNIST's training dataset, this will be your training dataset. Similarly, create a validation dataset for  
6   the two classes by picking samples corresponding to the first two classes from the validation set of  
7   the MNIST dataset. You can subsample input images from  $28 \times 28$  to  $14 \times 14$  if you need.

8   (b) **(15 points)** Logistic regression solves for

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

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

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

17   (c) **(10 points)** Write down the Hessian of the loss function in (1). Without assuming any special  
18   conditions about the dataset  $\{(x_i, y_i)\}_{i=1, \dots, n}$ , is this problem strongly convex? What is the best  
19   strong convexity parameter for the loss function in (1)?

20   (d) **(15 points)** Optimize the objective in (1) again, this time using gradient descent with Nesterov's  
21   acceleration. If we knew the condition number  $\kappa$ , what momentum coefficient would we use for this  
22   problem? It is difficult, although possible, to evaluate  $\kappa$ , so let's use

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

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

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

30   (e) **(20 points)** We will now optimize the same problem with stochastic gradient descent (SGD). Use  
31   a batch-size  $b = 128$  and optimize (1) using SGD with and without Nesterov's acceleration. Plot the  
32   training loss against the number of parameter updates on a semi-log scale (log scale on the Y-axis) for  
33   (i) gradient descent with Nesterov's updates (can be the same plot from your previous solutions), (ii)  
34   SGD without Nesterov's acceleration and, (iii) SGD with Nesterov's acceleration. Is the convergence  
35   for (iii) faster than that of (ii)? Comment on the differences for the convergence curves of (i) and (ii).

36   **Problem 2 (60 points). (You need Colab/AWS for this problem)** In this problem, we will see a  
37   heuristic that helps find out a good schedule for the learning rate. You will use the All-CNN network  
38   from Homework 2 in

39   <https://gist.github.com/pratikac/68d6d94e4739786798e90691fb1a581b> and train it on the CIFAR-10

40 dataset. Here is one learning rate schedule that has been noticed to work in practice

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

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

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

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

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

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

49 The only parameter left to choose is  $\eta_{\max}$ .

50 (a) **(5 points)** Why does the cross-entropy loss for CIFAR-10 start at around 2.3 for the All-CNN  
51 network with the default PyTorch initialization?

52 (b) **(15 points)** We saw in the lecture that for gradient descent to converge, in particular to make  
53 monotonic progress, the learning rate is limited by the smoothness coefficient of the loss function. It  
54 is difficult to estimate the Lipschitz constant of the gradients of a typical deep network and therefore  
55 to pick a good learning rate: pick it too small and you train too slowly, pick it to be too large and the  
56 loss does not go down quickly enough. We will use exactly this property to pick  $\eta_{\max}$ . Here is the  
57 algorithm, often called the “learning-rate finder”.

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

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

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

71 Run the above algorithm and draw the plot in part (ii). Identify the global minimum of this plot  
72 (doing it manually is enough) and call the corresponding learning rate  $\eta^*$ . You do not need to be  
73 very precise in identifying the global minimum, so long as you eyeball something close to the global  
74 minimum you can proceed to the next step.

75 (c) (20 points) We cannot use the learning rate  $\eta^*$  directly to choose a value of  $\eta_{\max}$  in (2). This is  
 76 because the network had already trained for a few time-steps before we tried out  $\eta^*$  in part (a.ii). We  
 77 should be more conservative in picking the maximum learning rate and will set

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

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

```
82 for g in optimizer.param_groups:  
83     g['lr'] = lr
```

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

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

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

92 is kept unchanged, the validation error more or less remains the same. You will train the network for  
 93 50 epochs and measure the validation error at the end of the 50 epochs for three settings:

- 94 (i)  $\eta_{\max}$  and  $\rho = 0.9$
- 95 (ii)  $\eta_{\max} \leftarrow 5\eta_{\max}$  and  $\rho = 0.5$
- 96 (iii)  $\eta_{\max} \leftarrow \eta_{\max}$  and  $\rho = 0.5$

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

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

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

100 is called the effective learning rate in the hyper-parameter optimization literature and often determines  
 101 the validation error if everything else remaining the same. Here  $\lambda$  is the coefficient of weight-decay,  
 102  $\rho$  is the momentum parameter as above and  $\theta$  is the batch-size.

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