

UNIVERSITY OF PENNSYLVANIA  
ESE 546: PRINCIPLES OF DEEP LEARNING  
HOMEWORK 3

- You will submit solutions typeset in  $\text{\LaTeX}$  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 1 PDF” where you will upload the PDF of your solutions. You will also see entries like “HW 1 Problem 1 Code” and “HW 1 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 1 Problem 3 Code” or “HW 1 Problem 3 Code”. Name your file as `pennkey_hw1_problem3.py`, e.g., I will name my code as `pratikac_hw1_problem3.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.
- **In addition to submitting the code, you should append the entire Python code for the particular problem to the solution in the PDF. If you are using Latex, you can do something like the screenshot below. The instructors will execute the code to check it. Your code should run without any errors and should create all output/plots required in the problem.**

```
\includepackage{pythonhighlight}

\begin{python}

a = np.array(10)

...

\end{python}
```

**Credit** The points for the problems add up to 110. 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 (35 points, do this problem on your laptop when you are debugging but Colab with**  
2 **GPU will train much faster).** The pre-training procedure for a large language model (LLM) involves  
3 minimizing the cross-entropy loss for predicting the next token on a large corpus of data. In  
4 this problem, you will fit a small Transformer on the dataset that you used for problem 4 of HW 2 (a  
5 concatenation of the complete works of Shakespeare, Leo Tolstoy’s War and Peace and any English  
6 book of your choice from Project Gutenberg). We wish to build a network that can complete a given  
7 sentence.

8 (a) **(3 points)** First, you should tokenize the dataset using a more sophisticated tokenizer than  
9 the characters that you used last time. You will use what is called a byte pair encoding (BPE)  
10 [https://en.wikipedia.org/wiki/Byte-pair\\_encoding](https://en.wikipedia.org/wiki/Byte-pair_encoding).  
11 See the example shown at [https://en.wikipedia.org/wiki/Byte-pair\\_encoding#Example](https://en.wikipedia.org/wiki/Byte-pair_encoding#Example). BPE works  
12 by scanning the text to identify frequently occurring byte pairs and collapsing them into a  
13 new symbol into the vocabulary. This is straight-forward to implement but it is quite slow in  
14 Python (it requires multiple passes over the dataset to compute the frequencies and compress the  
15 text). We will therefore use an existing library to build the vocabulary. Read the webpage at  
16 <https://huggingface.co/docs/tokenizers/quicktour> carefully to understand how to use Huggingface’s  
17 tokenizers library (install it using “pip install tokenizers”) to train and run a tokenizer.

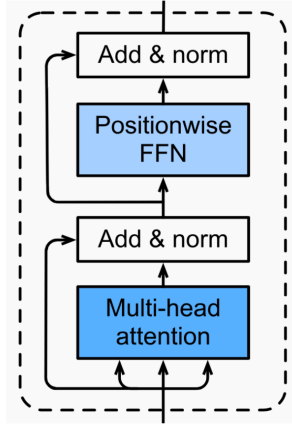
18

```
1 from tokenizers import Tokenizer
2 from tokenizers.trainers import BpeTrainer
3 from tokenizers.models import BPE
4 from tokenizers.pre_tokenizers import BertPreTokenizer
5
6 tokenizer = Tokenizer(BPE(unk_token="[UNK]"))
7 tokenizer.pre_tokenizer=BertPreTokenizer()
8
9 trainer = BpeTrainer(special_tokens=["[UNK]", "[CLS]", "[SEP]", "[PAD]", "[MASK]"])
10 tokenizer.train(['enwiki8.txt'],trainer)
```

19 You can set up your code to do something like the picture above. Two things to note here. First,  
20 we have showed an example that uses a few special tokens while tokenizing the text, e.g., [UNK]  
21 stands for “unknown token” for when the data might contain new characters that were not a part of the  
22 vocabulary, the [CLS] token is sometimes used for downstream tasks using pre-trained BERT/GPT  
23 representations, [PAD] stands for padding which allows the user to create a mini-batch of short and  
24 long sentences, etc. The vocabulary of chatbots contains many such special tokens for start of turn,  
25 end of turn, beginning of “thinking”, end of “thinking” etc. Second, we are using a pre-tokenizer  
26 from Bert, this procedure simply sanitizes the text, e.g., splitting apostrophes with a space on both  
27 sides (so that the tokenizer will have the apostrophe symbol as an element of the vocabulary), etc.  
28 Such procedures are broadly called stemming and lemmatization in the NLP literature and they are a  
29 very important step before any language modeling can be done (Spacy <https://spacy.io/api/tokenizer>  
30 is an excellent library for this purpose).

31 You can give options to BpeTrainer to use a particular vocabulary size, use something in the region of  
32 1024–2048.

33 (b) **(17 points)** Next you will write code for a (decoder-only) Transformer-based network. You  
34 may not use pre-coded Transformer blocks from any library (this includes PyTorch layers such as  
35 nn.Transformer, nn.TransformerEncoder etc.). You should write the following code by yourself using  
36 basic PyTorch tensors. Here is a rough guideline for the various parameters is as follows.



- A context length of  $T = 128$ .
- Each token is embedded as a vector into  $d = 256$  dimensions with sine-cosine-based position encoding.
- The feature dimension of keys, queries and values is  $p = 32$ , the number of heads is 2.
- Each Transformer block has one self-attention-based layer and one MLP layer, with layer normalization and residual connections as shown in the picture above. We will code up all these blocks ourselves. You can set the bias to False for all operations (linear operations, key/value computation, as well as normalization).
- For multi-head self-attention if input features are  $h \in \mathbb{R}^{T \times p}$ , the computations of the multi-head attention block for each head denoted by the index  $m = 1, 2$  are:

$$\mathbb{R}^{T \times p} \ni k^m = \sigma(h w_k^m) \text{ with } w_k^m \in \mathbb{R}^{p \times p}$$

$$\mathbb{R}^{T \times p} \ni q^m = \sigma(h w_q^m)$$

$$\mathbb{R}^{T \times p} \ni v^m = \sigma(h w_v^m)$$

$$\forall t : \mathbb{R}^p \ni h_t'^m = \text{dropout} \left( \sum_{s=1}^{t-1} \left( \frac{\exp(k_s^{m\top} q_t^m / \sqrt{p})}{\sum_{u=1}^{t-1} \exp(k_u^{m\top} q_t^m / \sqrt{p})} \right) v_t^m \right).$$

We then concatenate the output of each head to get features  $h' \in \mathbb{R}^{T \times 2p}$  before an output projection layer

$$\mathbb{R}^{T \times p} \ni h_{\text{attn}} = h' w_o \text{ with } w_o \in \mathbb{R}^{2p \times p}.$$

- Residual connection and normalization are easy to write

$$\mathbb{R}^{T \times p} \ni h_{\text{attn}} = \text{layer norm}(h + h_{\text{attn}})$$

You can use PyTorch's `nn.LayerNorm` to implement this operation.

- The next thing we should is to implement a position-wise MLP to mix the features computed by self-attention. This looks like

$$\mathbb{R}^{T \times p} \ni h = \text{layer norm}(h_{\text{attn}}(t) + \text{MLP}(h_{\text{attn}}(t)))$$

54 with the MLP being a two-layer network with ReLU nonlinearities,  $p$  input neurons,  $2p$  hidden  
 55 neurons and  $p$  output neurons. Notice that the MLP acts on the input data position-wise.  
 56 • All this together forms one “block” of a Transformer. Your network should have 2–4 such  
 57 blocks.

58 (c) **(10 points)** Train this network on the tokenized data for predicting the next token using the past  
 59 tokens in the context (you will see that the equations above have causal attention). In this problem, we  
 60 are only interested in completing a given sentence, so we will create a mini-batch where each datum  
 61 is a sentence. Remember that you will need to use padding while creating mini-batches that contain  
 62 sentences of unequal lengths.

63 You should plot the (i) training loss of successive mini-batches as a function of the weight updates  
 64 (feel free to do a moving average), (ii) periodically, say after every 1000 weight updates, you can also  
 65 compute the validation loss on the validation dataset.

66 (d) **(5 points)** Write code for auto-regressive sampling using this network. Given a prompt, say “The  
 67 Eiffel Tower is in Paris”, your code will tokenize the prompt (let’s say this prompt contains 6 tokens  
 68  $x_1, x_2, \dots, x_6$ ), and run the sampler step-by-step as

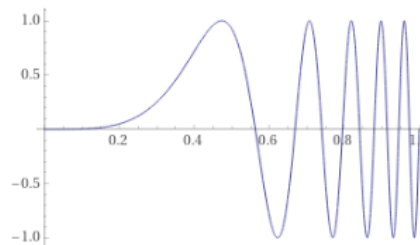
$$\begin{aligned}\hat{x}_7 &\sim p(\cdot \mid x_1, \dots, x_6) \\ \hat{x}_8 &\sim p(\cdot \mid x_1, \dots, x_6, \hat{x}_7) \\ \hat{x}_9 &\sim p(\cdot \mid x_1, \dots, x_6, \hat{x}_7, \hat{x}_8),\end{aligned}$$

69 each sample involves a call to the network and a random draw from the probability distribution of the  
 70 logits. Remember that you need to use position-encoding when you feed-in the context. You should  
 71 report a few examples of the network completing your prompts (use prompts of length 5–6 words,  
 72 completions of length 7–8 words).

73 **Problem 2 (20 points. Do this problem on your laptop; you do not really need a GPU for this**  
 74 **problem.).** In this problem, we will use nonlinear regression to understand one possible reason why  
 75 deep networks can avoid overfitting even when the number of samples in the dataset is smaller than  
 76 the number of weights. We will effectively show that generalization occurs when the test data lies  
 77 within the “convex hull” of the training data. If the test samples lie far away from the train data, then  
 78 one should not expect any generalization with networks with lots of parameters such as deep networks.  
 79 A second purpose of working on this problem is to see that we can do interesting experiments and  
 80 train 1000s of networks on a laptop (even without a GPU).

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

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



82

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

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

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

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

96 Write a function to evaluate the quantity

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

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

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

102 (c) **(10 points)** For each of 20 different values of  $n$ , e.g.,  $ns = \text{np.logspace}(1, 3, 20).astype(int)$ , create  
 103 5 training datasets (i.e., you will create 100 different datasets in total) like it is described in part  
 104 (a), fit an MLP to this dataset and evaluate the mean and standard deviation (across the 5 datasets  
 105 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  
 106 deviation) versus  $n$ . On the same plot also draw  $\log \delta f_{\text{out}}(n)$  (mean and standard deviation) versus  $n$ .  
 107 Interpret your plot.

108 **Problem 3 (20 points).** In this problem you will show that co-coercivity of the gradient and its  
 109 Lipschitz continuity are equivalent. Assume that the function  $f(x)$  is convex.

110 (a) **(5 points)** For a differentiable function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , show that co-coercivity of  $\nabla f$  implies  
 111 Lipschitz continuity of  $\nabla f$ . In other words, show that for all  $x, y$

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

112 Hint: use the Cauchy-Schwartz inequality.

113 (b) **(8 points)** Show the converse, i.e., show that Lipschitz continuity of  $\nabla f$  implies co-coercivity.

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

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

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

117 **Problem 4 (35 points). (Do this on your laptop)** In this problem, we will implement logistic  
 118 regression for classifying two classes (zero and one) from MNIST. You may not use any routines from  
 119 PyTorch other than the ones that help download the data. Use Numpy to code up the optimization  
 120 algorithm.

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

125 (b) **(10 points)** Logistic regression solves for

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

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

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

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

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

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

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

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

147 (e) **(10 points)** We will now optimize the same problem with stochastic gradient descent (SGD). Use  
 148 a batch-size  $\ell = 128$  (feel free to try different batch-size such as 64 and 8) and optimize Eq. (1) using  
 149 SGD with and without Nesterov's acceleration. Plot the training loss against the number of parameter  
 150 updates on a semi-log scale (log scale on the Y-axis) for (i) gradient descent with Nesterov's updates  
 151 (can be the same plot from your previous solutions), (ii) SGD without Nesterov's acceleration and,

152 (iii) SGD with Nesterov's acceleration. Is the convergence for (iii) faster than that of (ii)? Comment  
153 on the differences for the convergence curves of (i) and (ii).