

Due Date: March 24th 23:00, 2022

Instructions

- *This assignment is involved – please start well ahead of time.*
- *For all questions, show your work!*
- *Submit your report (PDF) and your code electronically via the course Gradescope page. Your report must contain answers to Problem 1 (Question 3), Problem 2 (Question 5), and to Problem 3 (all questions).*
- *For open-ended experiments (i.e., experiments that do not have associated test-cases), you do not need to submit code – a report will suffice.*
- *TAs for this assignment are **Ankit Vani** and **Sai Aravind Sreeramadas**.*

Summary: In this assignment, you will implement a sequential language model (an **LSTM**) and an image classifier (a **Vision Transformer**).

In problem 1, you will use built-in PyTorch modules to implement an LSTM and perform language modelling on wikitext and run some LSTM configurations. Download the LSTM embedding file from [here](#) and place it in `./data` folder.

In problem 2, you will implement various building blocks of a transformer, including **LayerNorm** (layer normalization) and the **Attention** mechanism for vision Transformer and build an image classifier on CIFAR10.

In problem 3, you will run the different transformer architectures and compare their performance to a simple CNN network.

The Wikitext-2 dataset comprises 2 million words extracted from the set of verified “Good” and “Featured” articles on Wikipedia. See this [blog post](#) for details about the Wikitext dataset and sample data. The dataset you get with the assignment has already been preprocessed using OpenAI’s GPT vocabulary, and each file is a compressed numpy array containing two arrays: **tokens** containing a flattened list of (integer) tokens, and **sizes** containing the size of each document.

You are provided a PyTorch dataset class (`torch.utils.data.Dataset`) named `Wikitext2` in the `utils` folder. This class loads the Wikitext-2 dataset and generates fixed-length sequences from it. Throughout this assignment, **all sequences will have length 256**, and we will use zero-padding to pad shorter sequences.

In practice though, you will work with mini-batches of data, each with batchsize `B` elements. You can wrap this dataset object into a `torch.utils.data.DataLoader`, which will return a dictionary with keys `source`, `target`, and `mask`, each of shape `(B, 256)`.

CIFAR10 dataset The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class. There are 50000 training images and 10000 test images. You are provided a PyTorch dataset class (`torch.utils.data.Dataset`) named `CIFAR10` from the `torchvision` and have the train, valid and test splits given. Throughout this assignment, the shape of CIFAR10 data is (Batch, Channels, Height, Width)

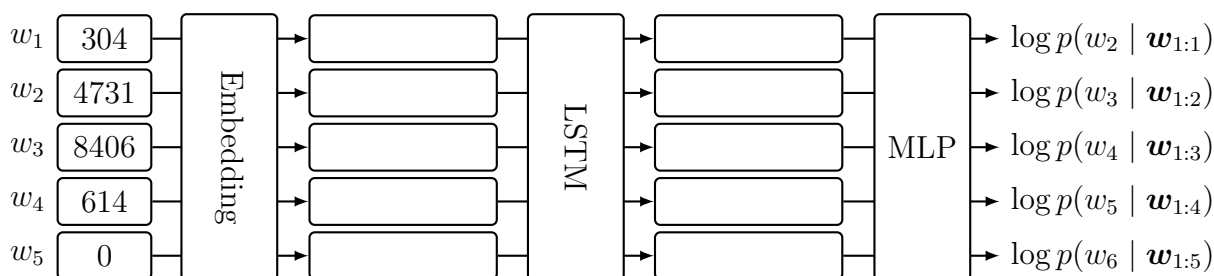
For students using Google Colab to complete their assignments, a cell with this command is available in the `main.ipynb` notebook.

If the tests on Gradescope fail: as a rule of thumb, x corresponds to the value in your assignment (e.g. the value returned by your function), and y is the expected value.

Coding instructions You will be required to use PyTorch to complete all questions. Moreover, this assignment **requires running the models on GPU** (otherwise it will take an incredibly long time); if you don't have access to your own resources (e.g. your own machine, a cluster), please use Google Colab (the notebook `main.ipynb` is here to help you). For some questions, you will be asked to not use certain functions in PyTorch and implement these yourself using primitive functions from `torch`; in that case, the functions in question are explicitly disabled in the tests on Gradescope.

Problem 1

Implementing an LSTM (13pts) In this problem, you will be using PyTorch's built-in modules in order to implement an LSTM. The architecture you will be asked to implement is the following:



In the file `lstm_solution.py`, you are given an `LSTM` class containing all the blocks necessary to create this model. In particular, `self.embedding` is a `nn.Embedding` module that converts sequences of token indices into embeddings, `self.lstm` is a `nn.LSTM` module that runs an LSTM over a sequence of vectors, and `self.classifier` is a 2-layer MLP responsible for classification. To do : Points per question

1. (5 pts) Using the different modules described above, complete the `forward()` function. This function must return the log-probabilities (not the logits) of the next words in the sequence, as well as the final hidden state of the LSTM.

2. (4 pts) Complete the `loss()` function, that returns the mean negative log-likelihood of the entire sequences in the minibatch (and also averaged over the mini-batch dimension). More precisely, for a single sequence in the mini-batch

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{w}_{1:T+1}) = -\frac{1}{T} \sum_{t=1}^T \sum_{i=0}^N \log p(\tilde{w}_{t+1} = i \mid \mathbf{w}_{1:t}; \boldsymbol{\theta}) \mathbf{1}(i = w_{t+1}),$$

where \tilde{w} are the predictions made by the model, and $\mathbf{1}(i = w_{t+1})$ is the indicator function which equals 1 if $i = w_{t+1}$, and 0 otherwise. Note that here T might be smaller than 256 (called `sequence_length` in the code), because the sequence might be zero-padded; you may use `mask` for this. The `loss` function directly takes the log-probabilities as input (e.g. returned by the `forward` function).

Training language models Unlike in classification problems, where the performance metric is typically accuracy, in language modelling, the performance metric is typically based directly on the cross-entropy loss, i.e. the negative log-likelihood (*NLL*) the model assigns to the tokens. For word-level language modelling it is standard to report **perplexity (PPL)**, which is the exponentiated average per-token NLL (over all tokens):

$$\exp \left(\frac{1}{TM} \sum_{t=1}^T \sum_{j=1}^M -\log p(\mathbf{w}_t^{(j)} \mid \mathbf{w}_1^{(j)}, \dots, \mathbf{w}_{t-1}^{(j)}; \boldsymbol{\theta}) \right),$$

where t is the index with the sequence, and j indexes different sequences. The purpose of this question is to perform model exploration, which is done using a validation set. As such, we do not require you to run your models on the test set.

3. (2 pts) Run the 6 configurations listed in `run_lstm.py`. For each of these experiments, plot learning curves (train and validation) of perplexity and loss over epochs. Figures should have labeled axes and a legend and an explanatory caption
4. (2 pts) Among the 6 configurations, which hyperparameters + optimizer would you use if you were most concerned with wall-clock time? With generalization performance?

Problem 2

Implementing a Vision Transformer (29pts) While typical RNNs “remember” past information by taking their previous hidden state as input at each step, recent years have seen a profusion of methodologies for making use of past information in different ways. The transformer¹ is one such fairly new architecture which uses several self-attention networks (“heads”) in parallel, among other architectural specifics. Implementing a transformer is a fairly involved process – so we provide

¹See <https://arxiv.org/abs/1706.03762> for more details.

most of the boilerplate code and your task is only to implement the multi-head scaled dot-product attention mechanism, as well as the layernorm operation.

Implementing Layer Normalization (5pts): You will first implement the layer normalization (LayerNorm) technique that we have seen in class. For this assignment, **you are not allowed** to use the PyTorch `nn.LayerNorm` module (nor any function calling `torch.layer_norm`).

As defined in the [layer normalization paper](#), the layernorm operation over a minibatch of inputs x is defined as

$$\text{layernorm}(x) = \frac{x - \mathbb{E}[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \text{weight} + \text{bias}$$

where $\mathbb{E}[x]$ denotes the expectation over x , $\text{Var}[x]$ denotes the variance of x , both of which are only taken over the last dimension of the tensor x here. `weight` and `bias` are learnable affine parameters.

1. (5pts) In the file `vit_solution_template.py`, implement the `forward()` function of the `LayerNorm` class. Pay extra attention to the lecture slides on the exact details of how $\mathbb{E}[x]$ and $\text{Var}[x]$ are computed. In particular, PyTorch's function `torch.var` uses an unbiased estimate of the variance by default, defined as the formula on the left-hand side

$$\overline{\text{Var}}(X)_{\text{unbiased}} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \qquad \overline{\text{Var}}(X)_{\text{biased}} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

whereas LayerNorm uses the biased estimate on the right-hand side (where \bar{X} here is the mean estimate). Please refer to the docstrings of this function for more information on input/output signatures.

Implementing the attention mechanism (17pts): You will now implement the core module of the transformer architecture – the multi-head attention mechanism. Assuming there are m attention heads, the attention vector for the head at index i is given by:

$$\begin{aligned} [q_1, \dots, q_m] &= QW_Q + b_Q & [k_1, \dots, k_m] &= KW_K + b_K & [v_1, \dots, v_m] &= VW_V + b_V \\ A_i &= \text{softmax} \left(\frac{q_i k_i^\top}{\sqrt{d}} \right) \\ h_i &= A_i v_i \\ A(Q, K, V) &= \text{concat}(h_1, \dots, h_m) W_O + b_O \end{aligned}$$

Here Q, K, V are queries, keys, and values respectively, where all the heads have been concatenated into a single vector (e.g. here $K \in \mathbb{R}^{T \times md}$, where d is the dimension of a single key vector, and T the length of the sequence). W_Q, W_K, W_V are the corresponding projection matrices (with biases b), and W_O is the output projection (with bias b_O). Q, K , and V are determined by the output of the previous layer in the main network. A_i are the attention values, which specify which elements of the input sequence each attention head attends to. In this question, **you are not allowed** to use the module `nn.MultiheadAttention` (or any function calling

`torch.nn.functional.multi_head_attention_forward`). Please refer to the docstrings of each function for a precise description of what each function is expected to do, and the expected input/output tensors and their shapes.

2. (4pts) The equations above require many vector manipulations in order to split and combine head vectors together. For example, the concatenated queries \mathbf{Q} are split into m vectors $[\mathbf{q}_1, \dots, \mathbf{q}_m]$ (one for each head) after an affine projection by \mathbf{W}_Q , and the \mathbf{h}_i 's are then concatenated back for the affine projection with \mathbf{W}_O . In the class `MultiHeadedAttention`, implement the utility functions `split_heads()` and `merge_heads()` to do both of these operations, as well as a transposition for convenience later. For example, for the 1st sequence in the mini-batch:

$$\begin{aligned} \mathbf{y} = \text{split_heads}(\mathbf{x}) &\rightarrow \mathbf{y}[0, 1, 2, 3] = \mathbf{x}[0, 2, \text{num_heads} * 1 + 3] \\ \mathbf{x} = \text{merge_heads}(\mathbf{y}) &\rightarrow \mathbf{x}[0, 1, \text{num_heads} * 2 + 3] = \mathbf{y}[0, 2, 1, 3] \end{aligned}$$

These two functions are exactly inverse from one another. Note that in the code, the number of heads m is called `self.num_heads`, and the head dimension d is `self.head_size`. Your functions must handle mini-batches of sequences of vectors, see the docstring for details about the input/output signatures.

3. (8pts) In the class `MultiHeadedAttention`, implement the function `get_attention_weights()`, which is responsible for returning \mathbf{A}_i 's (for all the heads at the same time) from \mathbf{q}_i 's and \mathbf{k}_i 's. Concretely, this means taking the softmax over the whole sequence. The softmax is then

$$[\text{softmax}(\mathbf{x})]_\tau = \frac{\exp(x_\tau)}{\sum_i \exp(x_i)}$$

4. (2pts) Using the functions you have implemented, complete the function `apply_attention()` in the class `MultiHeadedAttention`, which computes the vectors \mathbf{h}_i 's as a function of \mathbf{q}_i 's, \mathbf{k}_i 's and \mathbf{v}_i 's, and concatenates the head vectors.

$$\text{apply_attention}(\{\mathbf{q}_i\}_{i=1}^m, \{\mathbf{k}_i\}_{i=1}^m, \{\mathbf{v}_i\}_{i=1}^m) = \text{concat}(\mathbf{h}_1, \dots, \mathbf{h}_m).$$

5. (3pts) Using the functions you have implemented, complete the function `forward()` in the class `MultiHeadedAttention`. You may implement the different affine projections however you want (do not forget the biases), and you can add modules to the `__init__()` function. How many learnable parameters does your module have, as a function of `num_heads` and `head_size`?

The ViT forward pass (6pts): You now have all building blocks to implement the forward pass of a miniature ViT model. You are provided a module `PostNormAttentionBlock` which corresponds to a full block with self-attention and a feed-forward neural network, with skip-connections, using the modules `LayerNorm` and `MultiHeadedAttention` you implemented before.

In this part of the exercise, you will fill in the `VisionTransformer` class in `vit_solution_template.py`. This module contains all the blocks necessary to create this model. In particular, `get_patches()` is a module responsible for converting images in to token which are then converted into embeddings

(using input and positional embeddings), `self.layers` is a `nn.ModuleList` containing the different Attention Block layers, and `self.classifier` is a linear layer responsible for classification.

6. (2pts) Implement the function `get_patches()` which converts the image in to a sequence of patches based on the given patch size.
7. (1pts) By taken inspiration from the `PostNormAttentionBlock`, implement the `PreNormAttentionBlock`. You can look at the implementation of the forward function of the `PostNorm` block to complete the forward function in `PreNormAttentionBlock`. See the figure below for a comparison of the post-norm and pre-norm.

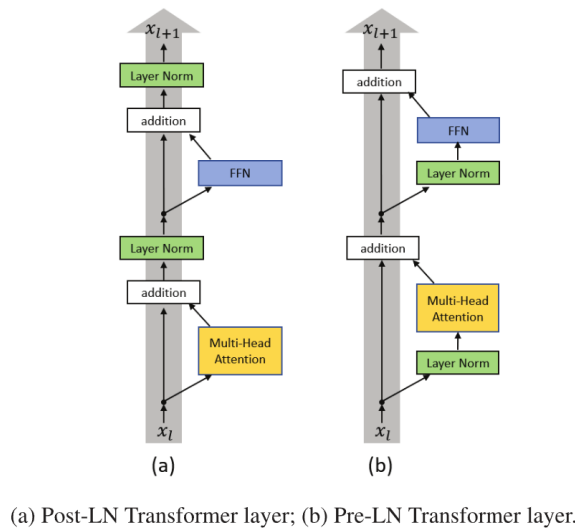


Figure 1: Image from Ruibin Xiong et al [On Layer Normalization in the Transformer Architecture](#)

8. (2pts) In the class `VisionTransformer`, complete the function `forward()` using the different modules described above.
9. (1pts) Complete the `loss()` function, that returns the cross entropy of the mini-batch.

Problem 3

Training ViT models (22pts) You will train each of the following architectures using an optimization technique and scheduler of your choice. For reference, we have provided a *feature-complete* training script (`run_exp_vit.py`) that uses the `ADAMW` optimizer. You are free to modify this script as you deem fit. You do not need to submit code for this part of the assignment. However, you are required to create a report that presents the accuracy and training curve comparisons as specified in the following questions.

Note: For each experiment, closely observe the training curves, and report the best validation accuracy score across epochs (not necessarily the validation score for the last epoch)

Configurations to run: At the top of the runner file (`run_exp_vit.py`), we have provided 6 experiment configurations for you to run. Together, these configurations span several choices of neural network architecture, optimizer, and weight-decay parameters. Perform the following analysis on the logs.

1. (4pts) You are asked to run 7 experiments with different optimizers, and hyperparameters settings. These parameter settings are given to you at the top of the runner file (`run_exp_vit.py`). For each of these experiments, plot learning curves (train and validation) of accuracy over both **epochs** and wall-clock time. Figures should have labeled axes and a legend and an explanatory caption.
2. (3pts) Make a table of results summarizing the train and validation performance for each experiment, indicating the architecture and optimizer.² Sort by architecture, then number of layers, then optimizer, and use the same experiment numbers as in the runner script for easy reference. Bold the best result for each architecture. The table should have an explanatory caption, and appropriate column and/or row headers. Any shorthand or symbols in the table should be explained in the caption.
3. (2pts) Among the first 6 configurations, which hyperparameters + optimizer would you use if you were most concerned with wall-clock time? With generalization performance?
4. (3pts) Between the experiment configurations 1-4 at the top of `run_exp_vit.py`, only the optimizer changed. What difference did you notice about the four optimizers used? What was the impact of weight decay, momentum, and ADAM?
5. (3pts) Compare experiments 6 and 7. Which model did you think performed better (PreNorm or PostNorm)? Why?
6. (3pts) In configurations 1- 7, you trained a transformer with various hyper-parameter settings. Given the recent high profile transformer based models, are the results as you expected? Speculate as to why or why not. How do they compare with CNN architectures given below.

Model	Val Accuracy	Test Accuracy	Num Parameters
GoogleNet	90.40%	89.70%	260,650
ResNet	91.84%	91.06%	272,378
ResNetPreAct	91.80%	91.07%	272,250
DenseNet	90.72%	90.23%	239,146

7. (2pts) For each of the experiment configurations above, measure the average steady-state GPU memory usage (`nvidia-smi` is your friend!). Comment about the GPU memory footprints

²You can also make the table in LaTeX; for convenience you can use tools like [LaTeX table generator](#) to generate tables online and get the corresponding LaTeX code.

of each model, discussing reasons behind increased or decreased memory consumption where applicable.

8. (2pts) Comment on the overfitting behavior of the various models you trained, under different hyperparameter settings. Did a particular class of models overfit more easily than the others? Can you make an informed guess of the various steps a practitioner can take to prevent overfitting in this case?