

Assignment #4

Instructor: Karthik Narasimhan

Read all the instructions below carefully before you start working on the assignment, and before you make a submission. The course assignment policy is available at <http://nlp.cs.princeton.edu/cos484>.

- This assignment contains 2 theory problems and 1 programming problem.
- We *highly* recommended that you typeset your submissions in L^AT_EX. Use the template provided on the website for your answers. If you've never used L^AT_EX, you can refer to the short guide here: <http://bit.ly/WorkingWithLaTeX>. Include your name and NetIDs with your submission. If you wish to submit hand written answers, you can scan and upload the pdf.
- Assignments must be uploaded to Gradescope **before class (12pm Eastern)** on the due date mentioned above.
- As per the late day policy outlined on the course website, you have 96 allowed late hours (about 4 days) to use at your discretion throughout the semester. Once you run out of late hours, late submissions will incur a penalty of 10% for each day, up to a maximum of 3 days beyond which submissions will not be accepted.
- Each problem is on a different page for ease of readability.
- The solutions to this assignment will contain 2 files which need to be submitted on Gradescope:
 1. One file containing solutions to the theory problems. There is no need to include code in this file. This file should be submitted on the Gradescope assignment titled "Assignment 4 - Theory". Please be sure to tag each problem with the right pages if you submit a PDF file.
 2. Colab .ipynb for the two programming problems. **Please include your programming question answers in the colab notebook.** This file should be submitted the Gradescope assignment titled "Assignment 4 - Programming submission". Make sure the cells are properly executed and have the outputs printed.

Problem 1: LSTMs vs Transformers

(8 points)

Both LSTMs and Transformers can be thought of a transformation from a sequence of input vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$ to a sequence of outputs $\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n \in \mathbb{R}^d$ (for simplicity, we assume their dimensions are the same). Let's compare the running time of an LSTM layer and a (single-head) self-attention layer in a Transformer:

- LSTM:

$$\begin{aligned} \mathbf{i}_i &= \sigma(\mathbf{W}^{(i)}\mathbf{h}_{i-1} + \mathbf{U}^{(i)}\mathbf{x}_i) & \mathbf{f}_i &= \sigma(\mathbf{W}^{(f)}\mathbf{h}_{i-1} + \mathbf{U}^{(f)}\mathbf{x}_i) \\ \mathbf{o}_i &= \sigma(\mathbf{W}^{(o)}\mathbf{h}_{i-1} + \mathbf{U}^{(o)}\mathbf{x}_i) & \mathbf{g}_i &= \tanh(\mathbf{W}^{(g)}\mathbf{h}_{i-1} + \mathbf{U}^{(g)}\mathbf{x}_i) \\ \mathbf{c}_i &= \mathbf{f}_i \odot \mathbf{c}_{i-1} + \mathbf{i}_i \odot \mathbf{g}_i & \mathbf{h}_i &= \mathbf{o}_i \odot \tanh(\mathbf{c}_i) \end{aligned}$$

where $\mathbf{W}^{(i)}, \mathbf{W}^{(f)}, \mathbf{W}^{(o)}, \mathbf{W}^{(g)}, \mathbf{U}^{(i)}, \mathbf{U}^{(f)}, \mathbf{U}^{(o)}, \mathbf{U}^{(g)} \in \mathbb{R}^{d \times d}$ (the bias terms are omitted).

- Self-attention:

$$\begin{aligned} \mathbf{q}_i &= \mathbf{W}^{(q)}\mathbf{x}_i, \quad \mathbf{k}_i = \mathbf{W}^{(k)}\mathbf{x}_i, \quad \mathbf{v}_i = \mathbf{W}^{(v)}\mathbf{x}_i, \\ \mathbf{h}_i &= \mathbf{W}^{(o)} \sum_{j=1}^n \left(\frac{\exp(\mathbf{q}_i \cdot \mathbf{k}_j / \sqrt{d})}{\sum_{j'=1}^n \exp(\mathbf{q}_i \cdot \mathbf{k}_{j'} / \sqrt{d})} \mathbf{v}_j \right) \end{aligned}$$

where $\mathbf{W}^{(q)}, \mathbf{W}^{(k)}, \mathbf{W}^{(v)}, \mathbf{W}^{(o)} \in \mathbb{R}^{d \times d}$.

(a) (4 points) Compare the running time of the forward pass of these two layers by counting the total number of floating-point multiplication operations in terms of n and d . Use the following assumptions:

1. For attention, assume that $\mathbf{q}_i \cdot \mathbf{k}_j / \sqrt{d}$ is computed for all i and j and cached. This can then be re-used to compute the denominator of the softmax.
2. You can ignore the cost of dividing by \sqrt{d} and dividing by the normalization constant for a softmax.
3. You can ignore the cost of activation functions and exponentials.

(b) (4 points) Suppose that we want to run these two layers on long sequences $n \gg d$ (e.g., $d = 512, n = 4096$), which layer runs faster in theory? In practice, which layer is easier to parallelize? Which layer is better at capturing long-term dependencies? Provide a brief reasoning for each answer.

Problem 2: Attention for Time-Series Data

(16 points)

We are interested in modeling time-series data and would like to use a masked self-attention layer with a single head. The masking pattern in this self-attention layer is causal, i.e. queries cannot attend to keys and values at future time steps, and it is the type of self-attention layer used in a Transformer decoder. We make a small modification to the masked self-attention layer: The layer takes inputs $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$, and outputs a sequence of scalars $y_1, y_2, \dots, y_n \in \mathbb{R}$, via the following implementation:

$$\begin{aligned} \mathbf{q}_i &= \mathbf{W}^{(q)} \mathbf{x}_i, \quad \mathbf{k}_i = \mathbf{W}^{(k)} \mathbf{x}_i, \quad \mathbf{v}_i = \mathbf{W}^{(v)} \mathbf{x}_i, \\ A_{i,j} &= \begin{cases} \exp(\mathbf{q}_i \cdot \mathbf{k}_j / \tau) / \left(\sum_{j'=1}^i \exp(\mathbf{q}_i \cdot \mathbf{k}_{j'} / \tau) \right) & j \leq i \\ 0 & j > i \end{cases} \\ y_i &= \sum_{j=1}^i A_{i,j} \mathbf{v}_j \end{aligned}$$

where $A_{i,j}$ are the indices of the attention matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, τ is a hyper-parameter that controls the softmax temperature and $\mathbf{W}^{(q)} \in \mathbb{R}^{2 \times d}$ is the query projection matrix, $\mathbf{W}^{(k)} \in \mathbb{R}^{2 \times d}$ the key projection matrix and $\mathbf{W}^{(v)} \in \mathbb{R}^{1 \times d}$ is the value projection matrix. We will make the following assumptions throughout the question:

1. The input vector at time step i has the following structure:

$$\mathbf{x}_i = \begin{bmatrix} u_i \\ \mathbf{p}_i \end{bmatrix}, \quad \text{where} \quad \begin{aligned} \mathbf{p}_1 &= [1 \ 0 \ 0 \ \dots \ 0]^T \\ \mathbf{p}_2 &= [0 \ 1 \ 0 \ \dots \ 0]^T \\ \mathbf{p}_3 &= [0 \ 0 \ 1 \ \dots \ 0]^T \\ &\vdots \\ \mathbf{p}_{d-1} &= [0 \ 0 \ 0 \ \dots \ 1]^T \end{aligned}$$

Here, $u_i \in \mathbb{R}$ is the time-series value at time step i and $\mathbf{p}_i \in \mathbb{R}^{d-1}$ encodes the positional (or temporal) information as a one-hot vector. Therefore, we can only encode sequence up to a maximum length of $N = d - 1$, i.e., the sequence length must be $1 \leq n \leq N$.

2. The weights of the key projection matrix are defined by:

$$\mathbf{W}^{(k)} = \begin{bmatrix} 0 & \cos \frac{2\pi(1)}{N} & \cos \frac{2\pi(2)}{N} & \dots & \cos \frac{2\pi(N-1)}{N} & \cos \frac{2\pi(N)}{N} \\ 0 & \sin \frac{2\pi(1)}{N} & \sin \frac{2\pi(2)}{N} & \dots & \sin \frac{2\pi(N-1)}{N} & \sin \frac{2\pi(N)}{N} \end{bmatrix},$$

where $N = d - 1$ is again the maximum sequence length. This ignores the time-series values and maps the positional information to an \mathbb{R}^2 vector. The weights of the value projection matrix are defined by:

$$\mathbf{W}^{(v)} = [1 \ 0 \ 0 \ \dots \ 0 \ 0].$$

This selects the information of the time-series values and ignores the positional information.

(a) (2 points) Assume that the maximum sequence length is $N = 8$. Since the key vectors $\mathbf{k}_1, \mathbf{k}_2, \dots, \mathbf{k}_n$ only use the positional information of the input sequence, we can compute the key vectors for any valid input sequence which follows our assumptions. Sketch these key vectors on a 2-d plane and label the vectors by their time-step. (*You do not need to show numeric values for these vectors.*)

(b) (2 points) Consider a layer where the weights of the query matrix are all zero:

$$\mathbf{W}^{(q)} = \begin{bmatrix} 0 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 \end{bmatrix},$$

Describe the resulting attention matrix for an input sequence of length n . If the input contains a time-series with values u_1, u_2, \dots, u_n , what would the output of the layer be?

Hint: Remember that the attention mechanism is causal!

(c) (2 points)

Describe the behavior of the attention mechanism in the limit of $\tau \rightarrow 0$? Pay attention to the special case of “ties”, when multiple entries have the same value, i.e. $A_{i,j} = A_{i,j'}$ for $j \neq j'$.

(d) (10 points) We now assume that the layer operates in the limit of $\tau \rightarrow 0$. Consider the following time-series operations and, for each, either:

- A. Show how the operation can be implemented with the attention layer by finding the corresponding attention matrix \mathbf{A} and a query projection matrix $\mathbf{W}^{(q)}$, or
 - B. Give a reason why the operation cannot be implemented with the attention layer under our assumptions.
- (i) Always select the value at the m 'th time step, where m is a hyperparameter and greater than 1:

$$y_i = u_m$$

- (ii) The time-series values are shifted to the right by one time step and the value at the first position remains the same:

$$y_i = \begin{cases} u_1 & \text{if } i = 1 \\ u_{i-1} & \text{otherwise} \end{cases}$$

- (iii) A moving average of the current and the previous time-series value:

$$y_i = \begin{cases} u_1 & \text{if } i = 1 \\ \frac{u_i + u_{i-1}}{2} & \text{otherwise} \end{cases}$$

Hint: You need to use the special case of “ties” from your answer in (c).

- (iv) An exponentially-weighted moving average of the current and all previous time-series values, defined by:

$$y_i = \left(\sum_{j=1}^i (2^{j-i}) u_j \right) / \left(\sum_{j=1}^i (2^{j-i}) \right)$$

Programming 1: Neural Machine Translation

(26 points)

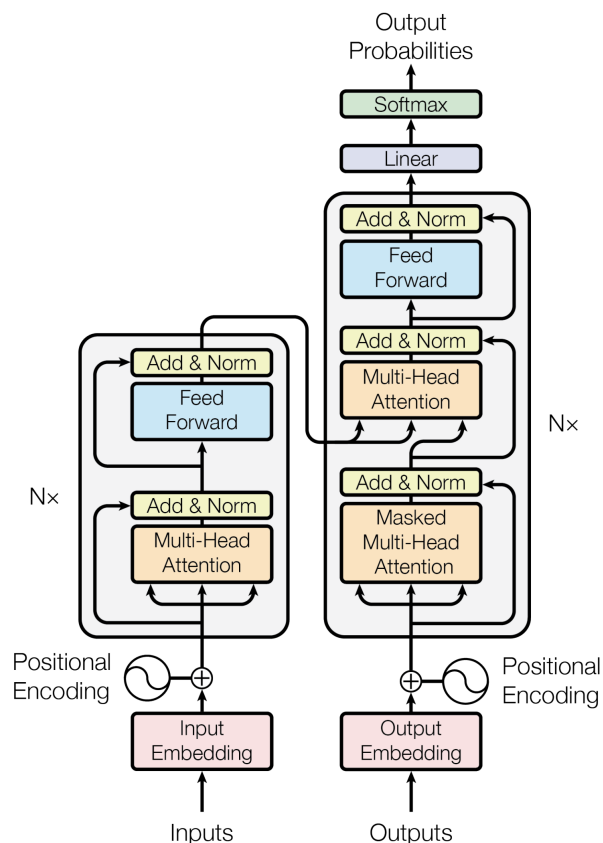
In this problem, you will implement a sequence-to-sequence (seq2seq) model based on the Transformer architecture to build a neural machine translation (NMT) system (translating from French to English).

Code setup Start from the code in this [Colab notebook](#). Please follow the instructions in the notebook to complete TODOs and train and evaluate the seq2seq model.

Data We will use a dataset consisting of parallel French and English sentences and partitioned into training, validation and test splits. You will need to tokenize the data using French and English tokenizers which are provided with the assignment resources. We will use BPE (byte pair encoding) tokenization, which can split a less common word into multiple subword tokens. If you are interested in learning more about BPE, see the paper [Neural Machine Translation of Rare Words with Subword Units](#), or this [blog post](#).

Model The model we will be implementing is an encoder-decoder Transformer model. For simplicity, we make several modifications to the architecture described in this paper:

- We use learned positional embeddings instead of sinusoidal positional embeddings.¹
- No dropout in the embedding layer and the attention weights.
- We use weight tying between the decoder's input embeddings and the output projection matrix (we have learned this in the class!).



We will describe each block of the model. Let d be the embedding dimension of input embeddings and hidden states, N the maximum sequence length and $V^{(e)}$ and $V^{(d)}$ the vocab sizes for encoder and decoder vocabulary, respectively. Let n be the length of a particular input sequence. In practice, the model will process B sequences in a batch in parallel and the sequence might contain pad tokens, which should be excluded from the attention mechanism and the loss function.

¹The authors find that this does not impact performance in their experiments. See Table 3 in the [Transformers](#) paper.

- **Embedding:** Let $\mathbf{E}^{(e)} \in \mathbb{R}^{V^{(e)} \times d}$ be the token embedding tables and $\mathbf{P}^{(e)} \in \mathbb{R}^{N \times d}$ be the positional embedding tables for the encoder. To embed a token with token id t at position i , we look up the token embedding at index t and position embedding at index i and add the two embeddings. The same procedure is done in the decoder embedding layer with separate matrices $\mathbf{E}^{(d)} \in \mathbb{R}^{V^{(d)} \times d}$ and $\mathbf{P}^{(d)} \in \mathbb{R}^{N \times d}$. The output of each embedding layer is a sequence of vectors $\mathbf{h}_1^{(0)}, \mathbf{h}_2^{(0)}, \dots, \mathbf{h}_n^{(0)} \in \mathbb{R}^d$.
- **Multi-head attention:** The input to this layer are a sequence of hidden state vectors from the previous layer $\mathbf{h}_1^{(l-1)}, \mathbf{h}_2^{(l-1)}, \dots, \mathbf{h}_n^{(l-1)} \in \mathbb{R}^d$. We project each of these vectors to query, key and value vectors:

$$\mathbf{q}_i^{(l)} = \mathbf{W}^{(q)} \mathbf{h}_i^{(l-1)}, \quad \mathbf{k}_i^{(l)} = \mathbf{W}^{(k)} \mathbf{h}_i^{(l-1)}, \quad \mathbf{v}_i^{(l)} = \mathbf{W}^{(v)} \mathbf{h}_i^{(l-1)},$$

where $\mathbf{q}_i^{(l)}, \mathbf{k}_i^{(l)}, \mathbf{v}_i^{(l)} \in \mathbb{R}^d$. For the sake of clarity, we are omitting the layer index (l) from the weight matrices and intermediate variables. In cross-attention, the input to key and value matrices would be the output states of the encoder, whereas the queries would be computed from the hidden states of the decoder.

We then split these vectors into H separate vectors, where H is the number of attention heads:

$$\mathbf{q}_i = \begin{bmatrix} \mathbf{q}_{i,1} \\ \mathbf{q}_{i,2} \\ \vdots \\ \mathbf{q}_{i,H} \end{bmatrix}, \quad \mathbf{k}_i = \begin{bmatrix} \mathbf{k}_{i,1} \\ \mathbf{k}_{i,2} \\ \vdots \\ \mathbf{k}_{i,H} \end{bmatrix}, \quad \mathbf{v}_i = \begin{bmatrix} \mathbf{v}_{i,1} \\ \mathbf{v}_{i,2} \\ \vdots \\ \mathbf{v}_{i,H} \end{bmatrix},$$

where $\mathbf{q}_{i,h}, \mathbf{k}_{i,h}, \mathbf{v}_{i,h} \in \mathbb{R}^{d_H}$ and $d_H = \frac{d}{H}$ is the head dimension. For each head h , we compute a scaled dot product attention:

$$\mathbf{y}_{i,h} = \sum_{j=1}^n \left(\frac{\exp(\mathbf{q}_{i,h} \cdot \mathbf{k}_{j,h} / \sqrt{d_H})}{\sum_{j'=1}^n \exp(\mathbf{q}_{i,h} \cdot \mathbf{k}_{j',h} / \sqrt{d_H})} \mathbf{v}_{j,h} \right),$$

where each $\mathbf{y}_{i,h} \in \mathbb{R}^{d_H}$.

We need to make sure to avoid attending to pad tokens, which should not affect the model's output. In practice, this is achieved by setting the attention score $\mathbf{q}_i^{(h)} \cdot \mathbf{k}_j^{(h)} / \sqrt{d_H}$ to a very large negative value if there is a pad token at position j . When computing the self-attention in the decoder, we use causal masking to avoid attending to future values:

$$\mathbf{y}_{i,h} = \sum_{j=1}^i \left(\frac{\exp(\mathbf{q}_{i,h} \cdot \mathbf{k}_{j,h} / \sqrt{d_H})}{\sum_{j'=1}^i \exp(\mathbf{q}_{i,h} \cdot \mathbf{k}_{j',h} / \sqrt{d_H})} \mathbf{v}_{j,h} \right),$$

Finally we stack the attention outputs and project each token to obtain a sequence of output vectors $\mathbf{h}_1^{(l)}, \mathbf{h}_2^{(l)}, \dots, \mathbf{h}_n^{(l)} \in \mathbb{R}^d$:

$$\mathbf{h}_i^{(l)} = \mathbf{W}^{(o)} \begin{bmatrix} \mathbf{y}_{i,1} \\ \mathbf{y}_{i,2} \\ \vdots \\ \mathbf{y}_{i,H} \end{bmatrix}$$

- **Feedforward layers:** *These are already implemented for you in the notebook!* The input to this layer are a sequence of hidden state vectors from the previous layer $\mathbf{h}_1^{(l-1)}, \mathbf{h}_2^{(l-1)}, \dots, \mathbf{h}_n^{(l-1)} \in \mathbb{R}^d$. Each feedforward layer learns two feedforward matrices: $\mathbf{W}^{(1)} \in \mathbb{R}^{d_I \times d}$ and $\mathbf{W}^{(2)} \in \mathbb{R}^{d \times d_I}$ (we are omitting the layer index again). These matrices are used to project each input vector to a larger intermediate dimension d_I , apply a ReLU activation and then project back to the original embedding space. Finally, dropout is applied.

$$\mathbf{h}_i^{(l)} = \text{Dropout}(\mathbf{W}^{(2)} \text{ReLU}(\mathbf{W}^{(1)} \mathbf{h}_i^{(l-1)}))$$

- **Add & Norm:** *These are already implemented for you in the notebook!* After each attention and feedforward block, we add a residual network connection and perform layer normalization. This helps with optimization issues of training deep Transformer networks.
- **Output layer:** We will re-use the token input embeddings and perform next token prediction at each position i in the decoder:

$$\text{logits}_i = \mathbf{E}^{(d)} \mathbf{h}_i^{(L)} \in \mathbb{R}^{V^{(d)}}$$

Note that in Pytorch, we do not have to compute the softmax when using the cross entropy loss function.

Tips

- Carefully read the function signatures and docstrings for the functions that you need to implement, particularly the type hints and tensor shapes.
- Before you get started on implementing the missing parts of the model implementation, make sure you read the rest of the provided code carefully. This can be useful to understand how each module is going to be used.
- Before starting to train models, make sure you visualize the attention weights to convince yourself that attention masking is working correctly. You are encouraged to implement additional tests and sanity checks for the rest of the code.
- Although this is an extensive coding assignment, we provide detailed documentation for each part you need to implement. The total number of lines that you need to implement can be below 100 lines. You should make use of PyTorch's documentation when needed <https://pytorch.org/docs/stable/>.

(a) (16 points) Complete the TODO items in the code in order:

1. Complete the tokenization code to produce tokenized datasets.
2. Implement the multi-head attention module. You are not allowed to use `nn.MultiheadAttention` or `nn.functional.scaled_dot_product_attention`. For full credit, avoid using python loops.
3. Before moving on to completing the other sub-modules, complete the sanity check for the multi-head attention. The generated plots should show the correct attention masking patterns. Make sure that they are included in your submitted notebook.
4. Implement the embedding layer, which computes and sums the token embeddings and the positional embeddings.
5. Put all the sub-modules together by implementing the main `EncoderDecoderModel`.

(b) (4 points) Now you should be ready to train an NMT system on the real data. Start the training process using the model that you just completed. Take a look at the hyperparameters defined in colab (don't change them!) and observe the training progress in the training log. Make sure to include the training log in your answer. Provide answers at the end of the notebook to the following questions:

- (i) What vocabulary size are we using for the source and target language including special tokens?
- (ii) Approximately how many source and target tokens are on average contained in a training batch? What proportion of these tokens are `<pad>` tokens on average?
- (iii) What is the specific purpose of saving the model parameters in a file `model.pt` throughout training in the code we provide?

You can alter the provided code to obtain the answers, but be careful not to break anything!

(c) (2 points) Load the trained model and evaluate the model on the test set. Report the BLEU score you've obtained on the test set. Manually look at some results and compare them with the gold answers. What do you think of the quality of the translations? Are these grammatical English sentences? Can you identify any common mistakes?

(d) (4 points) Consider the provided beam search method. This implementation is not efficient and performs a lot of repeated computation. Identify the issue and propose how you can fix it. In particular, describe how would you have to change the arguments and return values of your `EncoderDecoderModel` and its sub-modules?