

Homework 3: Transformer

Deep Learning

Spring 2025

The goal of homework 3 is to test your understanding of attention and Transformers.

In part 1, you should submit all your answers in a pdf file. As before, we recommend using L^AT_EX.

For part 2, you will implement some neural networks by adding your code to the provided ipynb file.

The due date of homework 3 is 11:59pm 03/09. Submit the following files in a zip file `your_net_id.zip` through NYU classes:

- `hw3_theory.pdf`
- `hw3_impl.ipynb`

The following behaviors will result in penalty of your final score:

1. 10% penalty for submitting your file without using the correct naming format (including naming the zip file, PDF file or python file wrong, adding extra files in the zip folder, like the testing scripts in your zip file).
2. 10% penalty for every extra day of lateness. Up to 4 days max (after that we won't accept submission)
3. 20% penalty for code submission that cannot be executed following the steps we mentioned.

1 Theory (50pt)

1.1 Attention (15pts)

This question tests your intuitive understanding of attention and its property.

- (a) (2pts) Given queries $\mathbf{Q} \in \mathbb{R}^{d \times n}$, keys $\mathbf{K} \in \mathbb{R}^{d \times m}$ and values $\mathbf{V} \in \mathbb{R}^{t \times m}$, describe the operations needed to calculate the output \mathbf{H} of the standard dot-product

attention. What is the output dimension? (You can use the softargmax_β function directly. It is applied to the column of each matrix).

- (b) (2pts) The scaling value of softargmax_β controls the sharpness of the output distribution. It is also what makes the usual attention operation /textitscaled dot product attention. Why is this scaling value necessary, and what is the usual choice for it?
- (c) (2pts) One advantage of the attention operation is that it can preserve a particular value vector \mathbf{v} to the output \mathbf{h} . When does this occur and, what should the scale β be in this case? Which of the four types of attention we are referring to? How can this be done when using fully connected architectures?
- (d) (2pts) On the other hand, the attention operation can incorporate many value vectors \mathbf{v} to generate a new output \mathbf{h} . Elaborate on when this occurs, and β should be for this effect. Which of the four types of attention we are referring to? How can this be done when using fully connected architectures?
- (e) (2pts) If we have a small perturbation to one of the \mathbf{k}_i (you could assume the perturbation is a zero-mean Gaussian with small variance, so the new $\hat{\mathbf{k}}_i = \mathbf{k}_i + \boldsymbol{\epsilon}$), how will the output of the \mathbf{H} change?
- (f) (2pts) If we have a small perturbation to one of the queries \mathbf{q}_i , how will the output of the \mathbf{H} change? How would this differ from the previous case?
- (g) (3pts) What should we do if some of our query outputs should NOT be a function of certain key-query inputs? What technique can we use to achieve this, and where does this operation take place?

1.2 Multi-headed Attention (3pts)

This question tests your intuitive understanding of Multi-headed Attention and its property.

- (a) (1pts) Given queries $\mathbf{Q} \in \mathbb{R}^{d \times n}$, $\mathbf{K} \in \mathbb{R}^{d \times m}$ and $\mathbf{V} \in \mathbb{R}^{t \times m}$, describe the operations for calculating the output \mathbf{H} of the standard multi-headed scaled dot-product attention? Assume we have h heads.
- (b) (2pts) Is there anything similar to multi-headed attention for convolutional networks? Explain why do you think they are similar.

1.3 Self Attention (10pts)

This question tests your intuitive understanding of Self Attention and its property.

- (a) (3pts) Attention aggregates information from a collection of key-value pairs based on a query. Describe how the attention operation used to implement self-attention. Furthermore, what are the conceptual benefits to self-attention and are there any drawbacks in comparison to sequential methods like recurrent networks?
- (b) (3pts) Explain what is positional encoding and why is it essential to self-attention in particular. What is the difference between absolute and relative positional encoding and when might one be preferred over the other?
- (c) (2pts) Show us one situation that the self attention layer behaves like an identity layer or permutation layer.
- (d) (2pts) Show us one situation that the self attention layer behaves like a convolution layer with a kernel larger than 1. You can assume we use positional encoding.

1.4 Transformer (11pts)

Read the original paper on the Transformer model: "Attention is All You Need" by Vaswani et al. (2017).

- (a) (3pts) Explain the primary differences between the Transformer architecture and previous sequence-to-sequence models (such as RNNs and LSTMs). Your answer should pertain to: (1) training-time computation (2) information flow and (3) how sequential information is managed.
- (b) (3pts) Explain how self-attention in particular is beneficial for the Transformer model and the task of machine translation.
- (c) (3pts) Explain the feed-forward neural networks used in the model and their purpose.
- (d) (2pts) Name two techniques used in the paper to improve training stability of the transformer model, in particular regards to the issue of exploding / vanishing gradients. And briefly explain how they do so.

1.5 Vision Transformer (11pts)

Read the paper on the Transformer model: "An Image is Worth 16×16 Words: Transformers for Image Recognition at Scale".

- (a) (3pts) What is the key difference between the Vision Transformer (ViT) and traditional convolutional neural networks (CNNs) in terms of handling input images? Can you spot a convolution layer in the ViT architecture?

- (b) (3pts) What is the role of positional embeddings in the Vision Transformer model, and how do they differ from positional encodings used in the original Transformer architecture?
- (c) (2pts) How does the Vision Transformer model generate the final classification output? Describe the process and components involved in this step.
- (d) (3pts) How does ViT compare with CNN in terms of performance across different data regimes? What explains this trend?

2 Implementation (50pt)

Please add your solutions to this notebook [hw3-impl.ipynb](#). **Please use your NYU account to access the notebook.** The notebook contains parts marked as TODO, where you should put your code or explanations. The notebook is a Google Colab notebook, you should copy it to your drive, add your solutions, and then download and submit it to NYU Classes. You're also free to run it on any other machine, as long as the version you send us can be run on Google Colab.