

HOMework 3

16824 VISUAL LEARNING AND RECOGNITION (FALL 2023)

<https://piazza.com/class/llo9w21ejlp2f>

RELEASED: Wed, 25th Oct 2023

DUE: Wed, 20th Nov 2023

Instructor: Jun-Yan Zhu

TAs: Zhipeng Bao, Wen-Hsuan Chu, Yeojin Jung, Rutika Moharir

START HERE: Instructions

- **Collaboration policy:** All are encouraged to work together BUT you must do your own work (code and write up). If you work with someone, please include their name in your write-up and cite any code that has been discussed. If we find highly identical write-ups or code or lack of proper accreditation of collaborators, we will take action according to strict university policies. See the [Academic Integrity Section](#) detailed in the initial lecture for more information.
- **Late Submission Policy:** There are a **total of 10** late days across all homework submissions. Submissions that use additional late days will incur a 10% penalty per late day.
- **Submitting your work:**
 - We will be using Gradescope (<https://gradescope.com/>) to submit the Problem Sets. Please use the provided template only. Submissions must be written in LaTeX. All submissions not adhering to the template will not be graded and receive a zero.
 - **Deliverables:** Please submit all the `.py` files. Add all relevant plots and text answers in the boxes provided in this file. TO include plots you can simply modify the already provided latex code. Submit the compiled `.pdf` report as well.

NOTE: Partial points will be given for implementing parts of the homework even if you don't get the mentioned numbers as long as you include partial results in this pdf.

1 Image Captioning with Transformers (70 points)

We will be implementing the different pieces of a Transformer decoder ([Transformers](#)), and train it for image captioning on a subset of the [COCO dataset](#).

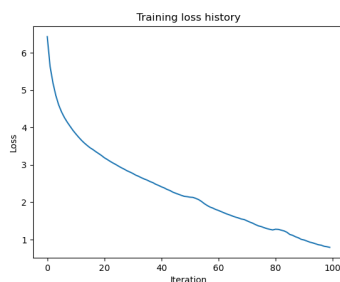
- **Setup:** Run the following command to extract COCO data, in the `transformer_captioning/datasets` folder: `./get_coco_captioning.sh`
- **Question:** Follow the instructions in the `README.md` file in the `transformer_captioning` folder to complete the implementation of the transformer decoder.
- **Deliverables:** After implementing all parts, use `run.py` for training the full model. The code will log plots to `plots`. Extract plots and paste them into the appropriate section below.
- **Expected results:** These are expected training losses after 100 epochs. Do not change the seed in `run.py`.
 - 2-heads, 2-layers, lr 1e-4: Final loss ≤ 1
 - 4-heads, 6-layers, lr 1e-4: Final loss ≤ 0.3
 - 4-heads, 6-layers, lr 1e-3: Final loss ≤ 0.05

1. Paste training loss plots for each of the three hyper-param configs

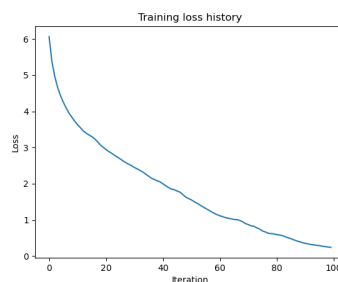
2-heads-2-layers-lr-1e-4: **0.793922**

4-heads-6-layers-lr-1e-4: **0.241308**

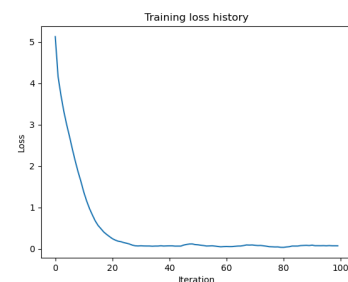
4-heads-6-layers-lr-1e-3: **0.074889** (**Best Value obtained 0.038355**)



(a) 2-heads-2-layers-lr-1e-4



(b) 4-heads-6-layers-lr-1e-4



(c) 4-heads-6-layers-lr-1e-3

2. Paste any three generated captioning samples from the training set with the three different settings. The provided code creates these plots at the end of training.

train
the woman in the <UNK> of the <UNK> the woman in the woman <END>
GT:<START> the woman in the <UNK> hold the box of donuts <END>



train
the woman in the <UNK> hold the box of donuts <END>
GT:<START> the woman in the <UNK> hold the box of donuts <END>



train
the woman in the <UNK> hold the box of donuts <END>
GT:<START> the woman in the <UNK> hold the box of donuts <END>



(a) Sample1: 2-heads-2-layers-lr-1e-4

(b) Sample2: 4-heads-6-layers-lr-1e-4

(c) Sample3: 4-heads-6-layers-lr-1e-3

3. Based on the observations of the three different settings, What would you change in the training procedure to get better validation performance? Why tweaking these hyper-parameters will lead to better performances?

Solution:

In the training, I can think of multiple ways to improve the validation performance. A few of them are:

1. Experiment with Learning Rate: I feel experimenting with an even smaller learning rate might help.
2. Experiment with Learning Rate Scheduler: Experimenting with an LR scheduler would possibly help in getting better results and using cyclic LR scheduler might even help in avoiding local optima. This can possibly be done in conjunction with the learning rate as well.
3. Experiment with number of heads and number of layers: Possibly increasing the number of heads and layers might help in capturing intricacies that might be missed due to the less number and it helps in noticing intricate patterns. Although increasing it way too much could also lead to some problems.
4. Experiment with Dropout: I would probably increase the dropout values and may even use different dropout values based on experimentation. Generally, based on my experience, increasing the dropout value to 0.3 or more although reduces the training accuracy, improves the validation accuracy, in turn helping the model to generalize slightly better.
5. Experiment with embedding dimensions: Maybe increasing the embedding dimension might help but that comes at the cost of computation, so if I have a lot of compute available, this might be a good hyperparameter to experiment with. It helps in capturing more.
6. Experiment with number of epochs: This is also related to the learning rate and might even scale with it.
7. Experiment with optimizer: In my experience, sometimes changing the optimizer with the appropriate learning rate also affects both the training and the validation datasets.

Tweaking these hyperparameters will lead to better performance because all of them help in capturing details that a shallow or the original hyperparameters might have missed.

2 Classification with Vision Transformers (30 points)

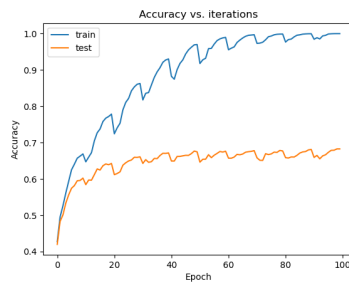
We will use the transformer you implemented in the previous part to implement a Vision Transformer (ViT), for classification on CIFAR10.

- **Question:** Follow the instructions in the `README.md` file in the `vit_classification` folder. You are encouraged to reuse code from the previous question.
- **Deliverables:** Run training using `run.py` for training the full model. The code will log plots `acc_out.png` (train and test accuracy) and `loss_out.png` (train loss).
- **Expected Results:** After 100 epochs, test accuracy should be $\geq 65\%$, train accuracy should be $\approx 100\%$, and training loss ≤ 0.3 .

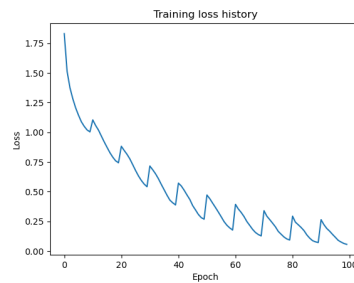
Final loss: **0.055999**

Final train accuracy: **0.999900** or **99.99 %**

Final test accuracy: **0.682700** or **68.27 %**



(a) Train/test accuracy



(b) Training loss

Collaboration Survey Please answer the following:

1. Did you receive any help whatsoever from anyone in solving this assignment?

☐ Yes

☒ No

- If you answered 'Yes', give full details:
- (e.g. "Jane Doe explained to me what is asked in Question 3.4")

2. Did you give any help whatsoever to anyone in solving this assignment?

☐ Yes

☒ No

- If you answered 'Yes', give full details:
- (e.g. "I pointed Joe Smith to section 2.3 since he didn't know how to proceed with Question 2")

3. Note that copying code or writeup even from a collaborator or anywhere on the internet violates the [Academic Integrity Code of Conduct](#).