

Natural Language Processing

Tel Aviv University

Assignment 2: Language ModelsDue Date: *Monday, July 22nd, 2024*

Lecturer: Maor Ivgi, TA: Noa Mark

1 Word-Level Neural Bi-gram Language Model

In this question, you will implement and train neural language model, and evaluate it on using perplexity.

- (a) Derive the gradient with respect to the input of a softmax function when cross entropy loss is used for evaluation, i.e., find the gradients with respect to the softmax input vector θ , when the prediction is made by $\hat{y} = \text{softmax}(\theta)$. Cross entropy and softmax are defined as:

$$\text{CE}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \cdot \log(\hat{y}_i)$$

$$\text{softmax}(\theta)_i = \frac{\exp(\theta_i)}{\sum_j \exp(\theta_j)}$$

The gold vector \mathbf{y} is a one-hot vector, and the predicted vector $\hat{\mathbf{y}}$ is a probability distribution over the output space.

- (b) Derive the gradients with respect to the input \mathbf{x} in a one-hidden-layer neural network (i.e., find $\frac{\partial J}{\partial \mathbf{x}}$, where J is the cross entropy loss $\text{CE}(\mathbf{y}, \hat{\mathbf{y}})$). The neural network employs a sigmoid activation function for the hidden layer, and a softmax for the output layer. Assume a one-hot label vector \mathbf{y} is used. The network is defined as:

$$\mathbf{h} = \sigma(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1),$$

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{h}\mathbf{W}_2 + \mathbf{b}_2).$$

The dimensions of the vectors and matrices are $\mathbf{x} \in \mathbb{R}^{1 \times D_x}$, $\mathbf{h} \in \mathbb{R}^{1 \times D_h}$, $\hat{\mathbf{y}} \in \mathbb{R}^{1 \times D_y}$, $\mathbf{y} \in \mathbb{R}^{1 \times D_y}$. The dimensions of the parameters are $\mathbf{W}_1 \in \mathbb{R}^{D_x \times D_h}$, $\mathbf{W}_2 \in \mathbb{R}^{D_h \times D_y}$, $\mathbf{b}_1 \in \mathbb{R}^{1 \times D_h}$, $\mathbf{b}_2 \in \mathbb{R}^{1 \times D_y}$.

- (c) Implement the forward and backward passes for a neural network with one sigmoid hidden layer. Fill in your implementation in `q1c_neural.py`. Sanity check your implementation with `python q1c_neural.py`.
- (d) GloVe ([Global Vectors](#)) embeddings are a type of word embeddings that represent words as vectors in a high-dimensional space, based on the co-occurrence statistics of words in a corpus. They are related to the skip-gram embeddings you saw in class in that they both aim to capture the semantic and syntactic relationships between words, but GloVe embeddings incorporate global corpus-level information in addition to local context information. In this section you will be using GloVe embeddings to represent the vocabulary. Use the neural network to implement a bigram language model in `q1d_neural_lm.py`. Use GloVe embeddings to represent the vocabulary (`data/lm/vocab.embeddings.glove.txt`). Implement the `lm_wrapper` function, that is used by `sgd` to sample the gradient, and the `eval_neural_lm` function that is used for model evaluation. Report the dev perplexity in your written solution. Don't forget to include `saved_params_40000.npy` in your submission zip!

2 Generating Shakespeare Using a Character-level Language Model

In this section we will train a language model and use it to generate text.

Follow the instructions, complete the code, and answer the questions from this Google Colab notebook¹: <https://colab.research.google.com/drive/1WIUACyCAgrPiuKzNBwXNCh0zWrecLnCF?usp=sharing>

3 Perplexity

- (a) Show that perplexity calculated using the natural logarithm $\ln(x)$ is equal to perplexity calculated using $\log_2(x)$. i.e:

$$2^{-\frac{1}{M} \sum_{i=1}^M \log_2 p(s_i | s_1, \dots, s_{i-1})} = e^{-\frac{1}{M} \sum_{i=1}^M \ln p(s_i | s_1, \dots, s_{i-1})}$$

- (b) In this section you will be computing the perplexity of your previous trained models on two different passages. Please provide your results in the PDF file, as well as attach the code to your code files. The two different passages appear in the .zip file you've got. Their names are: `shakespeare_for_perplexity.txt` which contains a subset from the Shakespeare dataset, and `wikipedia_for_perplexity.txt` which contains a certain passage from Wikipedia. Please compute the perplexity of the bi-gram LM, and the model from section 3, on both these passages.
- (c) Try to explain the results you've got. Particularly, why there might be large gaps in perplexity, while looking on different passages.

4 Deep Averaging Networks

In this question we will implement [Deep Averaging Networks \(DAN\)](#), and work on the IMDB dataset. The [IMDB dataset](#) consists of positive and negative reviews for movies.

This exercise will also introduce you to the popular `transformers` package from [Hugging Face](#). It will also require reading some of the documentation of `pytorch`. The total number of lines of code you need to write for implementing the model itself is roughly 10 or less (i.e., not a lot). Complete the code, and answer the following questions: <https://colab.research.google.com/drive/141W2qonpIYahv0Wj-iJ-xqAhRFEa1muV?usp=sharing>

Note: The model in the paper uses GloVe embeddings, in this exercise, you will implement this model using GloVe embedding trained on slightly less data, so you should expect different results than the ones shown in the paper.

- (a) Implement the DAN model as described in section 3 in the paper. In a nutshell, the DAN model proposes to average the GloVe word embeddings to represent the sentence, and then pass this sentence representation through a multi-layer feed-forward network (or multi-layer perceptron). Your best model should get accuracy of at least 83.5% on the evaluation set (anything below will result in partial credit). Add feed-forward layers, and tune the learning rate and batch size as necessary. Include a plot of the evaluation accuracy as a function of the number of epochs.

¹Feel free to comment inside the notebook.

- (b) Word dropout is a popular method for regularizing the model and avoiding over-fitting. Read section 3.1 of the paper and add word dropout as described (see [pytorch](#) documentation), and include a plot of the accuracy of the model across different values of the dropout rate (See Figure 2 in the paper).
- (c) Train 4 models with an increasing number of hidden layers (from 0 to 3 hidden layers), and compare the accuracy as the number of layers increases. Before you start training, think about when will we start seeing the effect of diminishing returns, are the results the way you expected them to be? Did the linear model outperform the model with 4 hidden layers? Include a plot of the accuracy as a function of the number of layers.
- (d) Use `nn.ReLU` and 2 other activation functions (of your choosing) from the [torch documentation](#), include a plot of the accuracy across epochs. What have you learned from this experiment?
- (e) For your best model, sample 5 examples from the evaluation set that the model classified incorrectly and for each example try to explain why the model classified it incorrectly.

5 Sentiment Analysis (practical)

In this exercise, you will be asked to browse solutions for specific downstream tasks in NLP. The goal of this exercise is to help you become comfortable with using other researchers' code and common tools. We hope that this practical experience will aid you in your final project. We encourage you to continue diving deeper with critical thinking and creative ideas!

[Paper with code](#) is free and open resource platform with Machine Learning papers, code, datasets, methods and evaluation tables (by Meta AI Research). You can browse state-of-the-art solutions for various tasks according to categories, get familiar with thousands of benchmark datasets, compare models using different evaluation metrics, and much more.

Explore the [sentiment analysis on text challenge](#). Go to the task lead-board to see whether the article code is available in the description table. Currently, the best model (in the leadboard) for this task is T5-11B (We can see it's not the most active research topic, but it's nevertheless a good way to introduce you to some practical tools! Also, note that you will not load T5-11B in this exercise but T5-Small, due to limited resources).

Include the code you used for this section in the Colab Notebook, and answer to the questions in the pdf.

1. Download the article "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer" (top accuracy on this task). It's OK if you don't fully understand the meaning of this article and its solution! You will learn more about it during the course. Find the link to their GitHub repository and copy it here. Once you've opened their GitHub, answer the following questions: Which model have they made publicly available? Which dataset is used to benchmark sentiment analysis? How do you evaluate success in this task?
2. In the article, T5 was pre-trained on multiple datasets and fine-tuned for specific downstream tasks. The original repository only published the weights for the pre-trained version (and not for the fine-tuned version used in the sentiment analysis task).

Search for a T5 model fine-tuned on the SST2 dataset on the HuggingFace platform, and use it as demonstrated in [the Hugging face platform tutorial](#).

LLMs (I bet most of you are familiar with the term; you will learn more about it very soon) come in different sizes. Use the **T5-Small** pre-trained version.

Please state exactly which model will you use and load.

Now, use the model to predict the following sentences (print each result in a different cell in the Colab Notebook):

- "This movie is awesome"
 - "I didn't like the movie so much"
 - "I'm not sure what I think about this movie."
 - "Did you like the movie?"
3. Load the SST2 dataset and evaluate the accuracy of this model on it. Does it match the results declared on Papers with Code?
 4. Is the data in the SST2 dataset balanced? Why is this an important question when evaluating the model's performance with respect to this dataset?
 5. Can you think on some properties of sentiment analysis evaluation that human-evaluators may notice but are not considered when evaluating the accuracy of this dataset??
 6. (Bonus:) Imagine you are consulting with friends who work in a healthcare organization that has recently approved the use of an alternative supplier for one of their drugs. They want to monitor client satisfaction with both suppliers, ensuring they are equally satisfactory to patients. To achieve this, they have collected a dataset consisting of summaries from visits to a family doctor. They have confirmed that the dataset includes patients treated with each of the two drug suppliers. Now, they would like to classify the visit summaries as positive or negative to analyze the results. Since this dataset is quite unique, they seek your advice on how to approach this task.

Please suggest some general high-level approaches to solve the task. Explain to them which factors they should take into account and what the advantages and disadvantages are of the different approaches.