

▼ Fall 2023 Applied NLP Homework 4

Instructors: Dr. Mahdi Roozbahani, Wafa Louhichi, Dr. Nimisha Roy

Deadline: December 1st, 11:59PM AoE

Honor Code and Assignment Deadline

- No unapproved extension of the deadline is allowed. Late submission will lead to 0 credit.
- Discussion is encouraged on Ed as part of the Q/A. However, all assignments should be done individually.
- **Plagiarism is a serious offense.** You are responsible for completing your own work. You are not allowed to copy and paste, or paraphrase, or submit materials created or published by others, as if you created the materials. All materials submitted must be your own.
- **All incidents of suspected dishonesty, plagiarism, or violations of the Georgia Tech Honor Code will be subject to the institute's Academic Integrity procedures. If we observe any (even small) similarities/plagiarisms detected by Gradescope or our TAs, WE WILL DIRECTLY REPORT ALL CASES TO OSI, which may, unfortunately, lead to a very harsh outcome. Consequences can be severe, e.g., academic probation or dismissal, grade penalties, a 0 grade for assignments concerned, and prohibition from withdrawing from the class.**

Instructions for the assignment

- This entire assignment will be autograded through Gradescope. There is only one Gradescope submissions for this assignment:
 - **Homework 4:** Submit all files to this section
- We provided you different .py files and we added libraries in those files please DO NOT remove those lines and add your code after those lines. Note that these are the only allowed libraries that you can use for the homework.
- You will submit your implemented .py files to the corresponding homework section on Gradescope.
- You are allowed to make as many submissions until the deadline as you like. Additionally, note that the autograder tests each function separately, therefore it can serve as a useful tool to help you debug your code if you are not sure of what part of your implementation might have an issue.

Using the local tests

- For some of the programming questions we have included a local test using a small toy dataset to aid in debugging. The local test sample data and outputs are stored in .py files in the **local_tests** folder
- There are no points associated with passing or failing the local tests, you must still pass the autograder to get points.
- **It is possible to fail the local test and pass the autograder** since the autograder has a certain allowed error tolerance while the local test allowed error may be smaller. Likewise, passing the local tests does not guarantee passing the autograder.
- **You do not need to pass both local and autograder tests to get points, passing the Gradescope autograder is what is required for credit.**
- It might be helpful to comment out the tests for functions that have not been completed yet.
- It is recommended to test the functions as it gets completed instead of completing the whole class and then testing. This may help in isolating errors. Do not solely rely on the local tests, continue to test on the autograder regularly as well.

▼ Google Colab Setup (Optional for running on Colab)

If you choose to work on the assignment on Google Colab, the following cell may help get you set up. You may need to right click on the Applied NLP folder and Add shortcut to Drive. You do not have to run this cell if you are working on the notebook locally.

```
1 # # Mount google drive
2 # from google.colab import drive
3 # drive.mount('/content/drive/')
4 # # You may need to create an Applied_NLP/HW#/hw#_code/ folder
5 # %cd '/content/drive/MyDrive/Applied_NLP/HW#/hw#_code/'
6
7 ### If no GPU selected it will ask for GPU to be selected
8 # gpu_info = !nvidia-smi
9 # gpu_info = '\n'.join(gpu_info)
10 # if gpu_info.find('failed') >= 0:
11 #     print('Select the Runtime > "Change runtime type" menu to enable a GPU accelerator, ')
12 #     print('and then re-execute this cell.')
13 # else:
14 #     print(gpu_info)
15
16
17 # ## This wraps output text according to the window size
18 # from IPython.display import HTML, display
19
20 # def set_css():
21 #     display(HTML('''
22 # <style>
23 #     pre {
24 #         white-space: pre-wrap;
25 #     }
26 # </style>
27 # '''))
28 # get_ipython().events.register('pre_run_cell', set_css)
```

▼ Assignment Overview

In this homework we will explore non-linear text classification algorithms using deep neural networks.

We will reuse the datasets from HW3 and a new dataset for this exploration:

- The first dataset is a subset of a [Clickbait Dataset](#) that has article headlines and a binary label on whether the headline is considered clickbait.
- The second dataset is a subset of [Web of Science Dataset](#) that has articles and a corresponding label on the domain of the articles.

- The third dataset is the [CoNLL-2003](#) that has text corresponding with Part of Speech and Name Entity Recognition labels.

We will first explore LSTM and LSTM with Attention. Then we will be looking into attention based architectures : Transformers and we will apply it to sequence labeling task : Part of Speech Tagging (POS) and Named Entity Recognition (NER). We will also look into Topic Modeling.

Deliverables and Points Distribution

Q1: LSTM [15pts]

- **1.1 Implementing the LSTM Model** [10pts] Deliverables: [lstm.py](#)
 - [4pts] `__init__`
 - [6pts] forward
- **1.2 Classifying Clickbait Dataset using LSTM** [5pts] Deliverables: [lstm.py](#), [best_lstm.model.pt](#)

Q2: LSTM with Attention [15pts]

- **2.1 Implementing the LSTM + Attention Model** [10pts] Deliverables: [attention.py](#)
 - [2pts] `__init__`
 - [2pts] forward_lstm
 - [2pts] forward_attention
 - [2pts] forward_context
 - [2pts] forward
- **2.2 Classifying Web of Science Dataset using LSTM + Attention** [5pts] Deliverables: [attention.py](#), [best_attention.model.pt](#)

Q3: Transformers [15pts]

- **3.1 Implementing the BERT Classifier** [5pts] Deliverables: [bert.py](#)
 - [2pts] `__init__`
 - [3pts] forward
- **3.4 Classifying Web of Science Dataset using BERT Classifier** [5pts] Deliverables: [bert_clickbait.pkl](#)
- **3.5 Classifying Web of Science Dataset using BERT Classifier** [5pts] Deliverables: [bert_wos.pkl](#)

Q4: Sequence Labeling [15pts]

- **4.1 Implementing Sequence Labeling** [5pts] Deliverables: [sequenceLabeling.py](#)
 - [2pts] `__init__`
 - [3pts] forward
- **4.5 POS Tagging and NER using CoNLL-2003 Dataset** [10pts] Deliverables: [bert_pos.pkl](#), [bert_ner.pkl](#)

Q5: Topic Modeling [20pts]

- **Topic modeling with Latent Dirichlet Allocation (LDA)** [20pts] Deliverables: [lda.py](#)
 - [5pts] `tokenize_words`
 - [5pts] `remove_stopwords`
 - [5pts] `create_dictionary`
 - [5pts] `build_LDAModel`

Setup

Please checkout the `environment_setup.md` file to create the environment for this homework. This notebook is tested under the package versions noted in the Library Imports cell output below, and the corresponding packages can be downloaded from [miniconda](#). You may also want to get yourself familiar with several packages:

- [jupyter notebook](#)
- [numpy](#)
- [sklearn](#)
- [pytorch](#)

In the .py files please implement the functions that have `raise NotImplementedError`, and after you finish the coding, please delete or comment out `raise NotImplementedError`.

Library imports

```
1 ## This cell only needs to be run once after the environment is created and then it can be commented out
2 !pip install transformers
3 !pip install datasets
4 !pip install pyLDavis
5 # !pip install numpy==1.24 # may be necessary if running on colab only
```

```
1 #Import the necessary libraries
2 import pandas as pd
3 import pickle
4 import numpy as np
5 import scipy as sp
6 import sys
7 import re
8 from copy import deepcopy
9 import random
10 from sklearn.metrics import accuracy_score
11 import torch
12 import torch.nn as nn
13 from torch import optim
14 torch.manual_seed(10)
15 from torch.autograd import Variable
```

```

16 import torch.nn.functional as F
17 from torch.utils.data import DataLoader
18 import transformers
19 import nltk
20 nltk.download('stopwords')
21 from nltk.corpus import stopwords
22 import gensim
23 import sklearn
24
25
26 import pyLDAvis.gensim_models
27 import pickle
28 import pyLDAvis
29
30 import torchtext
31 import datasets
32
33 import warnings
34 warnings.filterwarnings("ignore")
35
36 %load_ext autoreload
37 %autoreload 2
38 %reload_ext autoreload
39
40 # If set to True then all N_EPOCHS will be set to run for one epoch only
41 # Some compute intensive training loop will break after the first batch
42 SHORT = False
43
44 print('Version information')
45
46 print('python: {}'.format(sys.version))
47 print('numpy: {}'.format(np.__version__))
48 print('pd: {}'.format(pd.__version__))
49 print('torch: {}'.format(torch.__version__))
50 print('transformers: {}'.format(transformers.__version__))
51 print('gensim: {}'.format(gensim.__version__))
52 print('pyLDAvis: {}'.format(pyLDAvis.__version__))
53 print('torchtext: {}'.format(torchtext.__version__))
54 print('datasets: {}'.format(datasets.__version__))

```

```

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
Version information
python: 3.10.12 (main, Jun 11 2023, 05:26:28) [GCC 11.4.0]
numpy: 1.24.0
pd: 2.1.2
torch: 2.1.0+cu118
transformers: 4.35.0
gensim: 4.3.2
pyLDAvis: 3.4.0
torchtext: 0.16.0+cpu
datasets: 2.14.6

```

▼ Load Dataset

We start by loading both data sets already split into an 80/20 train and test set.

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 df_train = pd.read_csv('./data/train.csv')
6 df_test = pd.read_csv('./data/test.csv')
7
8 # Separate dataframes into train and test lists
9 x_train, y_train = list(df_train['headline']), list(df_train['label'])
10 x_test, y_test = list(df_test['headline']), list(df_test['label'])

```

Below is the number of headlines in the train and test set as well as a sample of the article headlines and its binary label, where 0 is considered not clickbait and 1 is clickbait.

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 print(f'Number of Train Headlines: {len(x_train)}')
6 print(f'Number of Test Headlines: {len(x_test)}')
7
8 print('\n\nSample Label and Headlines:')
9 x = 105
10 for label, line in zip(y_train[x:x+5], x_train[x:x+5]):
11     print(f'{label}: {line}')
12
13 print('\n\nOutput of Sample Headlines without Print Statement:')
14 x_train[x:x+5]

```

```

Number of Train Headlines: 19200
Number of Test Headlines: 4800

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 # Save test and train as csv
6 df_train_wos = pd.read_csv('./data/train_wos.csv')
7 df_test_wos = pd.read_csv('./data/test_wos.csv')
8
9 # Separate dataframes into train and test lists
10 x_train_wos, y_train_wos = list(df_train_wos['article']), list(df_train_wos['label'])
11 x_test_wos, y_test_wos = list(df_test_wos['article']), list(df_test_wos['label'])
12
13 # Numerical label to domain mapping
14 wos_label = {0:'CS', 1:'ECE', 2:'Civil', 3:'Medical'}
15 # Numerical label to Numerical mapping
16 label_mapping = {0:0, 1:1, 4:2, 5:3}
17
18 for i, label in enumerate(y_train_wos):
19     y_train_wos[i] = label_mapping[label]
20 for i, label in enumerate(y_test_wos):
21     y_test_wos[i] = label_mapping[label]

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 print(f'Number of Train Articles: {len(x_train_wos)}')
6 print(f'Number of Test Articles: {len(x_test_wos)}')
7
8 print('\nLabel Key:', wos_label)
9
10 print('\nSample Label and Articles:\n')
11 x = 107
12 for label, line in zip(y_train_wos[x:x+3], x_train_wos[x:x+3]):
13     print(f'{label} - {wos_label[label]}: {line}')

Number of Train Articles: 1600
Number of Test Articles: 400

Label Key: {0: 'CS', 1: 'ECE', 2: 'Civil', 3: 'Medical'}

Sample Label and Articles:

0 - CS: An efficient procedure for calculating the electromagnetic fields in multilayered cylindrical structures is reported in this paper. Using symbolic computation, spectral Green's functions, suitable for numerical implementations are determined in compact and closed forms. Applications are presented for structures with two dielectric layers.

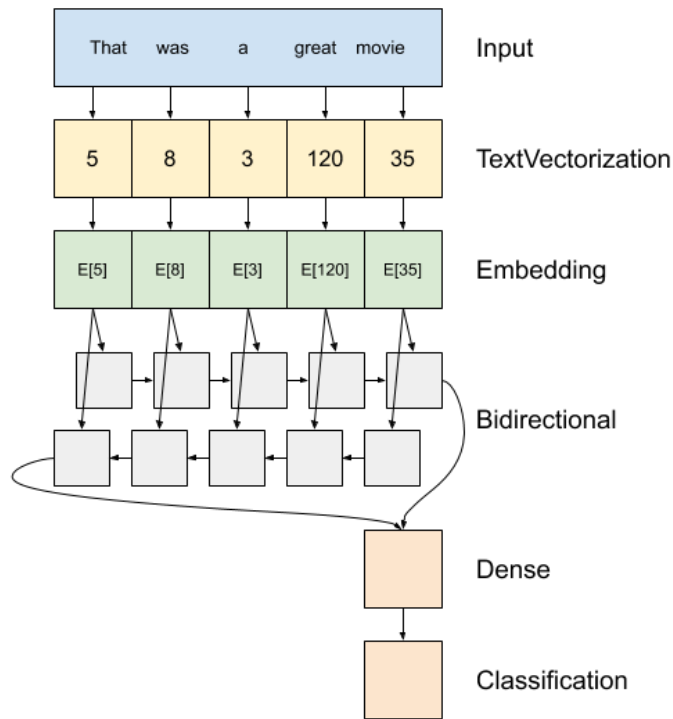
1 - ECE: A multifunctional platform based on the microhotplate was developed for applications including a Pirani vacuum gauge, temperature, and gas sensor. It consisted of a tungsten microhotplate and an on-chip operational amplifier. The platform was fabricated in a standard complementary metal oxide semiconductor (CMOS) process. A tungsten plug in standard CMOS process was specially designed as the serpentine resistor for the microhotplate, acting as both heater and thermister. With the sacrificial layer technology, the microhotplate was suspended over the silicon substrate with a 340 nm gap. The on-chip operational amplifier provided a bias current for the microhotplate. This platform has been used to develop different kinds of sensors. The first one was a Pirani vacuum gauge ranging from 10(-1) to 10(5) Pa. The second one was a temperature sensor ranging from -20 to 70 degrees C. The third one was a thermal-conductivity gas sensor, which could distinguish gases with different thermal conductivities in constant gas pressure and environment temperature. In the fourth application, with extra fabrication processes including the deposition of gas-sensitive film, the platform was used as a metal-oxide gas sensor for the detection of gas concentration.

2 - Civil: Artificial neural networks have been effectively used in various civil engineering fields, including construction management and labour productivity. In this study, the performance of the feed forward neural network (FFNN) was compared with radial basis neural network (RBNN) in modelling the productivity of masonry crews. A variety of input factors were incorporated and analysed. Mean absolute percentage error (MAPE) and correlation coefficient (R) were used to evaluate model performance. Research results indicated that the neural computing techniques could be

```

Q1: Classification with LSTM [15pts]

We will be using an LSTM (Long Short-Term Memory Network) for classification. The architecture of our model looks like :



For more details on LSTM, please refer to class lectures.

Use an Embedding layer, followed by a LSTM layer, and a linear layer.

We will then classify the Clickbait and Web of science dataset for this task.

▼ 1.1 : Implementing the LSTM Model [10 pts]

In the `lstm.py` file complete the following functions:

- `__init__`
- `forward`

We have included local tests in 1.2.2 for you to test your implementation

▼ 1.1.1 : Pre-Processing Data [No Points]

Run the cell below to load functions for building vocabulary and tokenizing the sentences.

```
1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from torchtext.data import get_tokenizer
6 from torchtext.vocab import build_vocab_from_iterator
7 import torchtext
8
9 tokenizer = get_tokenizer("basic_english")
10
11 def build_vocabulary(datasets):
12     for dataset in datasets:
13         for text in dataset:
14             yield tokenizer(text)
15
16 vocab = build_vocab_from_iterator(build_vocabulary([x_train]), min_freq=1, specials=["<UNK>"])
17 vocab.set_default_index(vocab["<UNK>"])
18
19 vocab_wos = build_vocab_from_iterator(build_vocabulary([x_train_wos]), min_freq=1, specials=["<UNK>"])
20 vocab_wos.set_default_index(vocab["<UNK>"])
```

▼ 1.2 : Classifying Clickbait Dataset using LSTM [5 Pts]

Run the cells below to classify the Clickbait train and test datasets using the LSTM functions that you have implemented.

An accuracy of more than 85% is acceptable.

▼ 1.2.1 : Create the Dataloaders [No Points]

```
1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from torch.utils.data import DataLoader
6
7 max_words = max(map(len, x_train))
8
9 def vectorize_batch(batch):
10     Y, X = list(zip(*batch))
11     X = [vocab(tokenizer(text)) for text in X] ## Tokenize and map tokens to indexes
```

```

12     X_len = [len(text) for text in X]
13     X = [tokens+([0]*(max_words-len(tokens))) if len(tokens)<max_words else tokens[:max_words] for tokens in X] ## Bringing all samples to max_words length.
14     return torch.tensor(X, dtype=torch.int32), torch.tensor(X_len), torch.tensor(Y)
15

```

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 train_dataset = list(map(lambda y, x: (y, x), y_train, x_train))
6 test_dataset = list(map(lambda y, x: (y, x), y_test, x_test))
7
8 train_loader = DataLoader(train_dataset, batch_size=1024, collate_fn=vectorize_batch, shuffle=True)
9 test_loader = DataLoader(test_dataset, batch_size=1024, collate_fn=vectorize_batch)

```

▼ 1.2.2: Local Tests for LSTM Functions [No Points]

You may test your implementation of the LSTM functions contained in **lstm.py** in the cell below. Feel free to comment out tests for functions that have not been completed yet. See [Using the Local Tests](#) for more details.

```

1 # #####
2 # ### DO NOT CHANGE THIS CELL ###
3 # #####
4
5 from lstm import LSTM
6 from local_tests.lstm_test import LSTM_Test
7 torch.manual_seed(10)
8
9 local_test = LSTM_Test()
10 lstm_model = LSTM(vocab, num_classes=2)
11 lstm_model.load_state_dict(torch.load('./local_tests/basic_lstm_model.pt')) # upload default weights. If this errors out, make sure you have initialized the layers correctly.
12 lstm_model.eval()
13
14 # creating a vectorized batch from the sample datapoints
15 local_test_dataset = list(map(lambda y, x: (y, x), local_test.output_labels, local_test.input_sequences))
16 X_localtest, X_len_localtest, _ = vectorize_batch(local_test_dataset)
17
18 print('Local Tests for LSTM Implementation \n')
19
20 # Local test for forward
21 output = lstm_model.forward(X_localtest, X_len_localtest)
22 forward_test = (output.shape == local_test.output.shape) and (torch.allclose(output, local_test.output, rtol=0.0001, atol=0.0001))
23 print('Your forward works as expected:', forward_test)

```

Local Tests for LSTM Implementation

Your forward works as expected: True

▼ 1.2.3 : Train and Evaluate on the Clickbait dataset [5 Pts]

You must reach an accuracy of >= 85% to receive credit for this section..

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 # Use cuda if present
6 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
7 print("Device available for running: ", device)

```

Device available for running: cuda

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from lstm import LSTM
6 from tqdm import tqdm
7
8 NUM_CLASSES = 2
9
10 model = LSTM(vocab, num_classes=NUM_CLASSES)
11 model.to(device)
12 criterion = nn.CrossEntropyLoss()
13 optimizer = optim.Adam(model.parameters(), lr=0.01)
14 N_EPOCHS = 5 if not SHORT else 1
15
16 model.train()
17 for epoch in range(N_EPOCHS):
18     total_loss = 0.0
19     for X, X_len, Y in tqdm(train_loader):
20         X = X.to(device)
21         Y = Y.to(device)
22         outputs = model(X, X_len)
23         loss = criterion(outputs, Y)
24
25         optimizer.zero_grad()
26         loss.backward()
27         optimizer.step()
28
29     total_loss += loss.item()
30
31     print("loss on epoch %i: %f" % (epoch, total_loss))

```

```

100%|██████████| 19/19 [00:01<00:00, 9.78it/s]
loss on epoch 0: 6.547718

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from sklearn.metrics import accuracy_score
6
7 with torch.no_grad():
8     Y_truth, Y_preds = [], []
9     for X, X_len, Y in test_loader:
10         X = X.to(device)
11         outputs = model(X, X_len)
12
13         Y_truth.append(Y)
14         Y_preds.append(outputs)
15
16 Y_truth = torch.cat(Y_truth)
17 Y_preds = torch.cat(Y_preds)
18
19 print("Test Accuracy on Clickbait Dataset using LSTM : {:.3f}".format(accuracy_score(Y_truth.cpu().detach().numpy(), F.softmax(Y_preds, dim=-1).argmax(dim=-1).cpu().detach().numpy())

Test Accuracy on Clickbait Dataset using LSTM : 0.967

```

▼ Submit these files to Gradescope

Run the cell below to save the state of your model. You will be required to upload the `best_lstm_model.pt` file that is saved in the `outputs` folder and the `lstn.py` file to Gradescope for accuracy evaluation. You must reach an accuracy of $\geq 85\%$ on the above test dataset to receive credit for this section.

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 torch.save(model.state_dict(), './outputs/best_lstm_model.pt')

```

▼ 1.3 : Classifying Web of Science Dataset using LSTM [No Points]

Run the cells below to classify the Web of Science train and test datasets using the `lstn` functions that you have implemented. You should observe an accuracy of more than 40%

▼ 1.3.1 : Create the Dataloaders [No Points]

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 max_words = max(map(len, x_train_wos))
6
7 def vectorize_batch(batch):
8     Y, X = list(zip(*batch))
9     X = [vocab_wos(tokenizer(text)) for text in X] ## Tokenize and map tokens to indexes
10    X_len = [len(text) for text in X]
11    X = [tokens+([0]*(max_words-len(tokens))) if len(tokens)<max_words else tokens[:max_words] for tokens in X]
12    return torch.tensor(X, dtype=torch.int32), torch.tensor(X_len), torch.tensor(Y)

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 train_dataset = list(map(lambda y, x: (y, x), y_train_wos, x_train_wos))
6 test_dataset = list(map(lambda y, x: (y, x), y_test_wos, x_test_wos))
7
8 train_loader = DataLoader(train_dataset, batch_size=128, collate_fn=vectorize_batch, shuffle=True)
9 test_loader = DataLoader(test_dataset, batch_size=128, collate_fn=vectorize_batch)

```

▼ 1.3.2 Train and Evaluate on the Web of Science Dataset [No Points]

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from lstn import LSTM
6 from tqdm import tqdm
7
8 NUM_CLASSES = 4
9
10 model = LSTM(vocab_wos, num_classes=NUM_CLASSES)
11 model.to(device)
12 criterion = nn.CrossEntropyLoss()
13 optimizer = optim.Adam(model.parameters(), lr=0.001)
14 N_EPOCHS = 20 if not SHORT else 1
15
16 model.train()
17 for epoch in range(N_EPOCHS):
18     total_loss = 0.0
19     for X, X_len, Y in tqdm(train_loader):
20         X = X.to(device)
21         Y = Y.to(device)
22         outputs = model(X, X_len)
23         loss = criterion(outputs, Y)
24
25         optimizer.zero_grad()
26         loss.backward()

```

```

27         optimizer.step()
28
29         total_loss += loss.item()
30
31         if SHORT:
32             break
33
34         print("loss on epoch %i: %f" % (epoch, total_loss))

```

```

100%|██████████| 13/13 [00:02<00:00, 4.73it/s]
loss on epoch 0: 17.730795
100%|██████████| 13/13 [00:02<00:00, 6.23it/s]
loss on epoch 1: 17.067045
100%|██████████| 13/13 [00:02<00:00, 6.13it/s]
loss on epoch 2: 16.460056
100%|██████████| 13/13 [00:02<00:00, 6.49it/s]
loss on epoch 3: 15.760631
100%|██████████| 13/13 [00:02<00:00, 6.43it/s]
loss on epoch 4: 15.091935
100%|██████████| 13/13 [00:02<00:00, 6.12it/s]
loss on epoch 5: 14.114826
100%|██████████| 13/13 [00:02<00:00, 4.59it/s]
loss on epoch 6: 13.049757
100%|██████████| 13/13 [00:02<00:00, 4.45it/s]
loss on epoch 7: 11.780611
100%|██████████| 13/13 [00:02<00:00, 5.32it/s]
loss on epoch 8: 10.286781
100%|██████████| 13/13 [00:02<00:00, 6.14it/s]
loss on epoch 9: 9.146650
100%|██████████| 13/13 [00:02<00:00, 6.31it/s]
loss on epoch 10: 8.352035
100%|██████████| 13/13 [00:02<00:00, 6.33it/s]
loss on epoch 11: 7.511066
100%|██████████| 13/13 [00:02<00:00, 6.26it/s]
loss on epoch 12: 6.068389
100%|██████████| 13/13 [00:02<00:00, 5.31it/s]
loss on epoch 13: 5.045795
100%|██████████| 13/13 [00:02<00:00, 4.49it/s]
loss on epoch 14: 4.757475
100%|██████████| 13/13 [00:02<00:00, 4.77it/s]
loss on epoch 15: 4.005289
100%|██████████| 13/13 [00:02<00:00, 5.78it/s]
loss on epoch 16: 3.193915
100%|██████████| 13/13 [00:02<00:00, 6.38it/s]
loss on epoch 17: 2.588529
100%|██████████| 13/13 [00:02<00:00, 4.41it/s]
loss on epoch 18: 2.277907
100%|██████████| 13/13 [00:02<00:00, 4.61it/s]loss on epoch 19: 1.821465

```

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from sklearn.metrics import accuracy_score
6
7 with torch.no_grad():
8     Y_truth, Y_preds = [], []
9     for X, X_len, Y in test_loader:
10         X = X.to(device)
11         outputs = model(X, X_len)
12
13         Y_truth.append(Y)
14         Y_preds.append(outputs)
15
16         if SHORT:
17             break
18
19     Y_truth = torch.cat(Y_truth)
20     Y_preds = torch.cat(Y_preds)
21
22 print("Test Accuracy on WoS Dataset using LSTM : {:.3f}".format(accuracy_score(Y_truth.cpu().detach().numpy(), F.softmax(Y_preds, dim=-1).argmax(dim=-1))

```

```

Test Accuracy on WoS Dataset using LSTM : 0.562

```

NOTE: LSTM alone is not able to perform good on the WoS dataset and that can be attributed to the very limited data with large vocabulary and lack of embedding structure.

▼ Q2: Classification with LSTM + Attention [15pts]

A potential issue with a vanilla LSTM approach is that a neural network needs to be able to compress all of the necessary information of a source sentence into a fixed-length vector. This may make it difficult for the neural network to cope with long sentences, especially those that are longer than the sentences in the training corpus. The Attention mechanism helps to look at all of the hidden states from the sequence for making predictions unlike the vanilla approach.

In this task, we will be implementing LSTM with Attention.



Please refer to lecture for more details.

You will be extending the LSTM model and incorporating attention on top of it.

We will then classify the Clickbait and Web of science dataset for this task.

2.1 : Implementing the LSTM + Attention Model [10 Points]

In the **attention.py** file complete the following functions:

- **__init__**
- **forward_lstm**
- **forward_attention**
- **forward_context**
- **forward**

We have included local tests in 2.2.2 for you to test your implementation

▼ 2.2 : Classifying Clickbait Dataset using LSTM with Attention [No Points]

Run the cells below to classify the Clickbait train and test datasets using the attention functions that you have implemented. You should observe an accuracy of more than 88%.

▼ 2.2.1 : Create the Dataloaders [No Points]

```
1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 max_words = max(map(len, x_train))
6
7 def vectorize_batch(batch):
8     Y, X = list(zip(*batch))
9     X = [vocab(tokenizer(text)) for text in X] ## Tokenize and map tokens to indexes
10    X_len = [len(text) for text in X]
11    X = [tokens+[0]*(max_words-len(tokens)) if len(tokens)<max_words else tokens[:max_words] for tokens in X] ## Bringing all samples to max_words len
12    return torch.tensor(X, dtype=torch.int32), torch.tensor(X_len), torch.tensor(Y)
```

```
1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 train_dataset = list(map(lambda y, x: (y, x), y_train, x_train))
6 test_dataset = list(map(lambda y, x: (y, x), y_test, x_test))
7
8 train_loader = DataLoader(train_dataset, batch_size=1024, collate_fn=vectorize_batch, shuffle=True)
9 test_loader = DataLoader(test_dataset, batch_size=1024, collate_fn=vectorize_batch)
```

▼ 2.2.2 : Local Tests for LSTM with Attention Functions [No Points]

You may test your implementation of the LSTM with Attention functions contained in **attention.py** in the cell below. Feel free to comment out tests for functions that have not been completed yet. See [Using the Local Tests](#) for more details.

```
1 # #####
2 # ### DO NOT CHANGE THIS CELL ###
3 # #####
4
5 from attention import Attention
6 from local_tests.attention_test import Attention_Test
7 torch.manual_seed(10)
8
9 local_test = Attention_Test()
10 attention_model = Attention(vocab, num_classes=2)
11 attention_model.load_state_dict(torch.load('./local_tests/basic_attention_model.pt')) # upload default weights. If this errors out, make sure you have
12 attention_model.eval()
13
14 # creating a vectorized batch from the sample datapoints
15 local_test_dataset = list(map(lambda y, x: (y, x), local_test.output_labels, local_test.input_sequences))
16 X_localtest, X_len_localtest, _ = vectorize_batch(local_test_dataset)
17 # print(X_localtest.shape, X_len_localtest)
18
19 print('Local Tests for LSTM + Attention Implementation \n')
20
21 # Local test for forward_lstm
22 output, hidden = attention_model.forward_lstm(X_localtest, X_len_localtest)
23 forward_lstm_shape_test = (output.shape == local_test.lstm_output.shape) and (hidden.shape == local_test.lstm_hidden.shape)
24 forward_lstm_value_test = (torch.allclose(output, local_test.lstm_output, rtol=0.001, atol=0.001)) and (torch.allclose(hidden, local_test.lstm_hidden, rt
25 print('Your forward_lstm works as expected: ', (forward_lstm_shape_test and forward_lstm_value_test))
26
27 # Local test for forward_attention
28 attn_weights = attention_model.forward_attention(local_test.lstm_output, local_test.lstm_hidden)
29 forward_attention_test = (attn_weights.shape == local_test.attention_weights.shape) and (torch.allclose(attn_weights, local_test.attention_weights, rtol=
30 print('Your forward_attention works as expected: ', forward_attention_test)
31
32 # Local test for forward_context
33 context_output = attention_model.forward_context(local_test.lstm_output, local_test.attention_weights, local_test.lstm_hidden)
34 forward_context_test = (context_output.shape == local_test.context_output.shape) and (torch.allclose(context_output, local_test.context_output, rtol=0.00
35 print('Your forward_context works as expected: ', forward_context_test)
36
37 # Local test for forward
38 output = attention_model.forward(X_localtest, X_len_localtest)
39 forward_test = (output.shape == local_test.output.shape) and (torch.allclose(output, local_test.output, rtol=0.0001, atol=0.0001))
40 print('Your forward works as expected: ', forward_test)

Local Tests for LSTM + Attention Implementation

Your forward_lstm works as expected: True
Your forward_attention works as expected: True
Your forward_context works as expected: True
Your forward works as expected: True
```

▼ 2.2.3: Train and Evaluate on the Clickbait Dataset [No Points]

```
1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 # Use cuda if present
6 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
7 print("Device available for running: ", device)

Device available for running:  cuda

1 #####
2 ### DO NOT CHANGE THIS CELL ###
```

```

3 #####
4
5 from attention import Attention
6 from tqdm import tqdm
7
8 NUM_CLASSES = 2
9
10 model = Attention(vocab, num_classes=NUM_CLASSES)
11 model.to(device)
12 criterion = nn.CrossEntropyLoss()
13 optimizer = optim.Adam(model.parameters(), lr=0.01)
14 N_EPOCHS = 5 if not SHORT else 1
15
16 model.train()
17 for epoch in range(N_EPOCHS):
18     total_loss = 0.0
19     for X, X_len, Y in tqdm(train_loader):
20         X = X.to(device)
21         Y = Y.to(device)
22         outputs = model(X, X_len)
23         loss = criterion(outputs, Y)
24
25         optimizer.zero_grad()
26         loss.backward()
27         optimizer.step()
28
29     total_loss += loss.item()
30
31     print("loss on epoch %i: %f" % (epoch, total_loss))

100%|██████████| 19/19 [00:01<00:00, 14.02it/s]
loss on epoch 0: 5.596366
100%|██████████| 19/19 [00:01<00:00, 14.65it/s]
loss on epoch 1: 1.509814
100%|██████████| 19/19 [00:01<00:00, 13.38it/s]
loss on epoch 2: 0.468158
100%|██████████| 19/19 [00:01<00:00, 12.86it/s]
loss on epoch 3: 0.147131
100%|██████████| 19/19 [00:02<00:00, 7.43it/s]loss on epoch 4: 0.031196

```

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from sklearn.metrics import accuracy_score
6
7 with torch.no_grad():
8     Y_truth, Y_preds = [], []
9     for X, X_len, Y in test_loader:
10         X = X.to(device)
11         outputs = model(X, X_len)
12
13         Y_truth.append(Y)
14         Y_preds.append(outputs)
15
16     Y_truth = torch.cat(Y_truth)
17     Y_preds = torch.cat(Y_preds)
18
19 print("Test Accuracy on Clickbait Dataset using LSTM with Attention : {:.3f}".format(accuracy_score(Y_truth.cpu().detach().numpy(), F.softmax(Y_preds, dim=-1).argmax(dim=-1).cpu().data.numpy())))

Test Accuracy on Clickbait Dataset using LSTM with Attention : 0.957

```

2.3 : Classifying Web of Science Dataset using LSTM with Attention [5 Pts]

Run the cells below to classify the Web of Science train and test datasets using the attention functions that you have implemented.

You must reach an accuracy of >=50% to receive credit for this section.

2.3.1 : Create the Dataloaders [No Points]

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 max_words = max(map(len, x_train_wos))
6
7 def vectorize_batch(batch):
8     Y, X = list(zip(*batch))
9     X = [vocab_wos(tokenizer(text)) for text in X] ## Tokenize and map tokens to indexes
10    X_len = [len(text) for text in X]
11    X = [tokens+([0]*(max_words-len(tokens))) if len(tokens)<max_words else tokens[:max_words] for tokens in X] ## Bringing all samples to max_words length.
12    return torch.tensor(X, dtype=torch.int32), torch.tensor(X_len), torch.tensor(Y)

```

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 train_dataset = list(map(lambda y, x: (y, x), y_train_wos, x_train_wos))
6 test_dataset = list(map(lambda y, x: (y, x), y_test_wos, x_test_wos))
7
8 train_loader = DataLoader(train_dataset, batch_size=128, collate_fn=vectorize_batch, shuffle=True)
9 test_loader = DataLoader(test_dataset, batch_size=128, collate_fn=vectorize_batch)

```

2.3.2 : Train and Evaluate on the Web of Science dataset [5 Pts]

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4

```

```

5 # Use cuda if present
6 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
7 print("Device available for running: ", device)

    Device available for running:  cuda

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from attention import Attention
6 from tqdm import tqdm
7
8 NUM_CLASSES = 4
9
10 model = Attention(vocab_wos, num_classes=NUM_CLASSES)
11 model.to(device)
12 criterion = nn.CrossEntropyLoss()
13 optimizer = optim.Adam(model.parameters(), lr=0.001)
14 N_EPOCHS = 25 if not SHORT else 1
15
16 model.train()
17 for epoch in range(N_EPOCHS):
18     total_loss = 0.0
19     for X, X_len, Y in tqdm(train_loader):
20         X = X.to(device)
21         Y = Y.to(device)
22         outputs = model(X, X_len)
23         loss = criterion(outputs, Y)
24
25         optimizer.zero_grad()
26         loss.backward()
27         optimizer.step()
28
29         total_loss += loss.item()
30
31     if SHORT:
32         break
33
34     print("loss on epoch %i: %f" % (epoch, total_loss))

100%|██████████| 13/13 [00:02<00:00, 6.25it/s]
loss on epoch 0: 17.639267
100%|██████████| 13/13 [00:02<00:00, 6.39it/s]
loss on epoch 1: 16.107268
100%|██████████| 13/13 [00:02<00:00, 6.20it/s]
loss on epoch 2: 14.302524
100%|██████████| 13/13 [00:01<00:00, 6.53it/s]
loss on epoch 3: 12.863941
100%|██████████| 13/13 [00:02<00:00, 4.74it/s]
loss on epoch 4: 11.127718
100%|██████████| 13/13 [00:02<00:00, 4.48it/s]
loss on epoch 5: 9.498385
100%|██████████| 13/13 [00:02<00:00, 4.60it/s]
loss on epoch 6: 8.060333
100%|██████████| 13/13 [00:02<00:00, 6.19it/s]
loss on epoch 7: 6.614249
100%|██████████| 13/13 [00:02<00:00, 6.43it/s]
loss on epoch 8: 5.804770
100%|██████████| 13/13 [00:02<00:00, 6.44it/s]
loss on epoch 9: 4.907240
100%|██████████| 13/13 [00:02<00:00, 6.36it/s]
loss on epoch 10: 3.515091
100%|██████████| 13/13 [00:02<00:00, 5.93it/s]
loss on epoch 11: 2.663215
100%|██████████| 13/13 [00:02<00:00, 4.58it/s]
loss on epoch 12: 1.798664
100%|██████████| 13/13 [00:02<00:00, 4.67it/s]
loss on epoch 13: 1.200081
100%|██████████| 13/13 [00:02<00:00, 5.31it/s]
loss on epoch 14: 0.735106
100%|██████████| 13/13 [00:02<00:00, 6.38it/s]
loss on epoch 15: 0.411879
100%|██████████| 13/13 [00:02<00:00, 6.40it/s]
loss on epoch 16: 0.252141
100%|██████████| 13/13 [00:02<00:00, 6.31it/s]
loss on epoch 17: 0.199521
100%|██████████| 13/13 [00:02<00:00, 6.27it/s]
loss on epoch 18: 0.141240
100%|██████████| 13/13 [00:02<00:00, 5.11it/s]
loss on epoch 19: 0.130441
100%|██████████| 13/13 [00:02<00:00, 4.68it/s]
loss on epoch 20: 0.075788
100%|██████████| 13/13 [00:02<00:00, 4.78it/s]
loss on epoch 21: 0.046381
100%|██████████| 13/13 [00:02<00:00, 5.93it/s]
loss on epoch 22: 0.032549
100%|██████████| 13/13 [00:02<00:00, 6.36it/s]
loss on epoch 23: 0.024646
100%|██████████| 13/13 [00:02<00:00, 6.44it/s]loss on epoch 24: 0.021094

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from sklearn.metrics import accuracy_score
6
7 with torch.no_grad():
8     Y_truth, Y_preds = [], []
9     for X, X_len, Y in test_loader:
10         X = X.to(device)
11         outputs = model(X, X_len)
12
13         Y_truth.append(Y)
14         Y_preds.append(outputs)
15
16     if SHORT:
17         break
18
19     Y_truth = torch.cat(Y_truth)
20     Y_preds = torch.cat(Y_preds)

```

```

21
22 print("Test Accuracy on WoS Dataset using LSTM with Attention : {:.3f}".format(accuracy_score(Y_truth.cpu().detach().numpy(), F.softmax(Y_preds, dim=-1).argmax(dim=-1).cpu().detach().numpy())))

Test Accuracy on WoS Dataset using LSTM with Attention : 0.547

```

Submit these files to Gradescope

Run the cell below to save the state of your model. You will be required to upload the `best_attention_model.pt` file that is saved in the `outputs` folder and the `attention.py` file to Gradescope for accuracy evaluation. You must reach an accuracy of $\geq 50\%$ on the above test dataset to receive credit for this section.

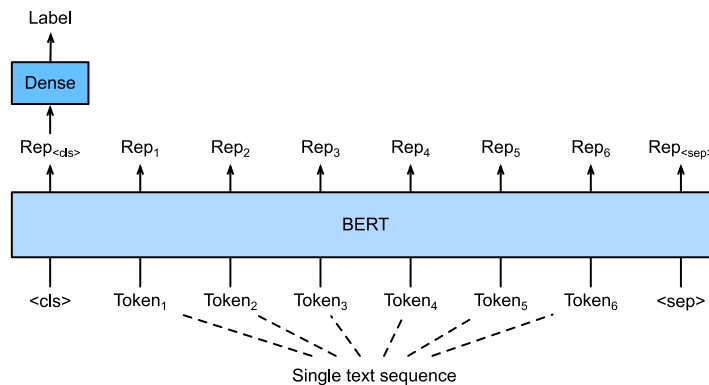
```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 torch.save(model.state_dict(), './outputs/best_attention_model.pt')

```

Q3: Classification with BERT [15pts]

The transformer neural network is a novel architecture that aims to solve sequence-to-sequence tasks while handling long-range dependencies with ease. In a transformer, we can pass all the words of a sentence and determine the word embedding simultaneously.



We will be using BERT (Bidirectional Encoder Representations from Transformers) pre-trained models for embeddings. BERT architecture consists of several Transformer encoders stacked together. Each Transformer encoder encapsulates two sub-layers: a self-attention layer and a feed-forward layer.

The details on BERT can be referred from the paper : [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#).

We will be using the BERT embeddings and a fully connected linear layer to perform classification.

We will then classify the Clickbait and Web of science dataset for this task.

3.1 : Implementing the BERT Classifier [10pts]

In the `bert.py` file complete the following functions:

- `__init__`
- `forward`

3.2: Local Tests for BERT Functions [No Points]

You may test your implementation of the BERT functions contained in `bert.py` in the cell below. Feel free to comment out tests for functions that have not been completed yet. See [Using the Local Tests](#) for more details.

```

1 # #####
2 # ### DO NOT CHANGE THIS CELL ###
3 # #####
4
5 from bert import BERTClassifier
6 from local_tests.bert_test import BERT_Test
7
8 local_test = BERT_Test()
9
10 print('Local Tests for BERT Functions \n')
11
12 # Instantiate BERTClassifier and load linear weights for local tests
13 torch.manual_seed(10)
14 student_bert = BERTClassifier(2)
15 basic_bert_linear = torch.load('./local_tests/basic_bert_linear.pt')
16
17 # If loading the state dict causes an error, then there may be an issue with your layer name or shape
18 student_bert.linear.load_state_dict(basic_bert_linear)
19
20 # Local test for forward
21 student_bert.eval()
22 outputs = student_bert(local_test.inputs, local_test.masks)
23 forward_test = (torch.allclose(outputs, local_test.outputs, rtol=0.0001, atol=0.0001) and (outputs.shape == local_test.outputs.shape))
24 print('Your forward works as expected:', forward_test)

```

Local Tests for BERT Functions

Downloading (...)ve/main/config.json: 100%

570/570 [00:00<00:00, 21.5kB/s]

Downloading model.safetensors: 100%

440M/440M [00:03<00:00, 81.0MB/s]

Your forward works as expected: True

▼ 3.3 : Initialize Tokenizer [No Points]

Run the below cells to initialize tokenizer for the pretrained BERT model.

```
1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from transformers import BertTokenizer
6
7 tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')

Downloading (...)okenizer_config.json: 100%          28.0/28.0 [00:00<00:00, 1.24kB/s]

Downloading (...)solve/main/vocab.txt: 100%         232k/232k [00:00<00:00, 3.98MB/s]

Downloading (...)main/tokenizer.json: 100%          466k/466k [00:00<00:00, 9.81MB/s]

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 def vectorize_batch(batch):
6     Y, X = list(zip(*batch))
7     X = tokenizer(X, padding='max_length', max_length = 128, truncation=True, return_tensors="pt")
8     input_ids, attention_mask = X['input_ids'], X['attention_mask']
9     return input_ids, attention_mask, torch.tensor(Y)
```

▼ 3.4 : Classifying Clickbait Dataset using BERT [5pts]

Run the below cell to classify the Clickbait train and test dataset using the bert functions that you have already implemented in 3.

You must reach an accuracy of >= 90% to receive credit for this section.

```
1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 train_dataset = list(map(lambda y, x: (y, x), y_train, x_train))
6 test_dataset = list(map(lambda y, x: (y, x), y_test, x_test))
7
8 train_loader = DataLoader(train_dataset, batch_size=32, collate_fn=vectorize_batch, shuffle=True)
9 test_loader = DataLoader(test_dataset, batch_size=32, collate_fn=vectorize_batch)

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 # Use cuda if present
6 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
7 print("Device available for running: ", device)

Device available for running:  cuda

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from bert import BERTClassifier
6 from tqdm import tqdm
7
8 NUM_CLASSES = 2
9
10 model = BERTClassifier(num_classes=NUM_CLASSES)
11 model.to(device)
12 criterion = nn.CrossEntropyLoss()
13 optimizer = optim.Adam(model.parameters(), lr=0.0001)
14 N_EPOCHS = 3 if not SHORT else 1
15
16 model.train()
17 for epoch in range(N_EPOCHS):
18     total_loss = 0.0
19     for X, X_mask, Y in tqdm(train_loader):
20         X = X.to(device)
21         X_mask = X_mask.to(device)
22         Y = Y.to(device)
23         outputs = model(X, X_mask)
24
25         loss = criterion(outputs, Y)
26
27         optimizer.zero_grad()
28         loss.backward()
29         optimizer.step()
30
31     total_loss += loss.item()
32
33     if SHORT:
34         break
35
36     print("loss on epoch %i: %f" % (epoch, total_loss))

100%|██████████| 600/600 [06:23<00:00, 1.57it/s]
loss on epoch 0: 45.047180
100%|██████████| 600/600 [06:30<00:00, 1.54it/s]
loss on epoch 1: 18.853489
100%|██████████| 600/600 [06:30<00:00, 1.54it/s]loss on epoch 2: 11.036154

1 #####
2 ### DO NOT CHANGE THIS CELL ###
```

```

3 #####
4
5 from sklearn.metrics import accuracy_score
6
7 with torch.no_grad():
8     Y_truth, Y_preds = [], []
9     for X, X_mask, Y in test_loader:
10         X = X.to(device)
11         X_mask = X_mask.to(device)
12         outputs = model(X, X_mask)
13
14         Y_truth.append(Y)
15         Y_preds.append(outputs)
16
17         if SHORT:
18             break
19
20 Y_truth = torch.cat(Y_truth)
21 Y_preds = torch.cat(Y_preds)
22
23 print("Test Accuracy on Clickbait Dataset using BERT : {:.3f}".format(accuracy_score(Y_truth.cpu().detach().numpy(), F.softmax(Y_preds, dim=-1).argmax(dim=-1).cpu().detach().numpy())))

```

Test Accuracy on Clickbait Dataset using BERT : 0.985

Submit this file to Gradescope

Run the cell below to save the results of your model. You will be required to upload the `bert_clickbait.pkl` file that is saved in the `outputs` folder to Gradescope for accuracy evaluation. You must reach an accuracy of $\geq 90\%$ on the above test dataset to receive credit for this section.

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 preds = F.softmax(Y_preds, dim=-1).argmax(dim=-1).cpu().detach().numpy()
6
7 with open('./outputs/bert_clickbait.pkl', 'wb') as fp:
8     pickle.dump(preds, fp)

```

3.5 : Classifying Web of Science Dataset using BERT [5pts]

Run the below cell to classify the Web of Science train and test dataset using the bert functions that you have already implemented in 3.

You must reach an accuracy of $\geq 70\%$ to receive credit for this section.

[] ↪ 6 cells hidden

Q4: Sequence Labeling [15pts]

Part-of-speech (POS) tagging is a popular Natural Language Processing process which refers to categorizing words in a text (corpus) in correspondence with a particular part of speech, depending on the definition of the word and its context.

Named entity recognition (NER) seeks to locate and classify named entities in text into pre-defined categories such as the names of persons, organizations, locations, expressions of times, quantities, monetary values, percentages, etc.

We will be using BERT (Bidirectional Encoder Representations from Transformers) for sequence labeling. The architecture of the model is shown above in the diagram.

We will be using the BERT embeddings and a fully connected linear layer to perform classification.

We will then classify using the conll2003 dataset for this task.

4.1 : Implementing Sequence Labeling [5pts]

In the `sequenceLabeling.py` file complete the following functions:

- `__init__`
- `forward`

4.2 : Local Tests for Sequence Labeling Functions [No Points]

You may test your implementation of the Sequence Labeling functions contained in `sequenceLabeling.py` in the cell below. Feel free to comment out tests for functions that have not been completed yet. See [Using the Local Tests](#) for more details.

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from sequenceLabeling import SequenceLabeling
6 from local_tests.sequenceLabeling_test import SequenceLabeling_Test
7
8 local_test = SequenceLabeling_Test()
9
10 print('Local Tests for Sequence Labeling Functions \n')
11
12 # Instantiate SequenceLabeling and load linear weights for local tests
13 torch.manual_seed(10)
14 student_seq = SequenceLabeling(9)
15 basic_seq_linear = torch.load('./local_tests/basic_seq_linear.pt')
16
17 # If loading the state dict causes an error, then there may be an issue with your layer name or shape
18 student_seq.linear.load_state_dict(basic_seq_linear)
19
20 # Local test for forward
21 student_seq.eval()
22 outputs = student_seq(local_test.input_ids, local_test.masks, local_test.token_type_ids)
23 forward_test = torch.allclose(outputs, local_test.outputs, rtol=0.0001, atol=0.0001)
24 print('Your forward outputs the expected shape:', outputs.shape == local_test.outputs.shape)

```

```

25 print('Your forward outputs the expected values:', forward_test)
26 if not forward_test:
27     print('If you forward values do not match the expected values but you are matching in shape and are getting the desired scores in 4.5, then no need to worry about the discrepancy.

```

Local Tests for Sequence Labeling Functions

```

Your forward outputs the expected shape: True
Your forward outputs the expected values: True

```

► 4.3 : Loading Dataset & Tokenizer [No Points]

Run the below cell to download the conll2003 dataset. Each word in the dataset has been put on a separate line and there is an empty line after each sentence. The first item on each line is a word, the second a part-of-speech (POS) tag, the third a syntactic chunk tag and the fourth the named entity tag. The chunk tags and the named entity tags have the format I-TYPE which means that the word is inside a phrase of type TYPE.

For more details on the structure and the supported labels, please refer : <https://huggingface.co/datasets/conll2003>.

```
[ ] ↩ 1 cell hidden
```

▼ 4.4 : Pre-process Dataset [No Points]

Run the below cells to vectorize the dataset for tokenizers and to initialize the dataloaders.

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 ### The code is used to vectorize the data for each batch as defined in the data loader.
6 ### For data in the batch, the input tokens are encoded and
7 ### then special tokens [CLS] (101) is added at the beginning and [SEP] (102) is added at the end.
8 ### These tokens are added because the BERT model was pretrained with these tokens. So to get the same results for inference we need to add them.
9 ### The input is padded with 0 if it is lesser than Max length.
10
11 MAX_LEN = 128
12
13 def vectorize_batch(batch):
14     batch_input_ids = []
15     batch_mask = []
16     batch_token_type_ids = []
17     batch_pos = []
18     batch_ner = []
19
20     for data in batch:
21         target_pos = []
22         target_ner = []
23         inputs = []
24         tokens = data['tokens']
25         pos_tags = data['pos_tags']
26         ner_tags = data['ner_tags']
27         for i in range(len(tokens)):
28             input = tokenizer.encode(tokens[i], add_special_tokens=False)
29             input_len = len(input)
30             target_pos.extend([pos_tags[i]] * input_len)
31             target_ner.extend([ner_tags[i]] * input_len)
32             inputs.extend(input)
33         inputs = inputs[:MAX_LEN - 2]
34         target_pos = target_pos[:MAX_LEN - 2]
35         target_ner = target_ner[:MAX_LEN - 2]
36
37         inputs = [101] + inputs + [102]
38         target_pos = [0] + target_pos + [0]
39         target_ner = [0] + target_ner + [0]
40
41         mask = [1] * len(inputs)
42         token_type_ids = [0] * len(inputs)
43
44         padding_len = MAX_LEN - len(inputs)
45         inputs = inputs + ([0] * padding_len)
46         mask = mask + ([0] * padding_len)
47
48         token_type_ids = token_type_ids + ([0] * padding_len)
49         target_pos = target_pos + ([0] * padding_len)
50         target_ner = target_ner + ([0] * padding_len)
51
52         batch_input_ids.append(inputs)
53         batch_mask.append(mask)
54         batch_token_type_ids.append(token_type_ids)
55         batch_pos.append(target_pos)
56         batch_ner.append(target_ner)
57
58     return torch.tensor(batch_input_ids, dtype=torch.long), torch.tensor(batch_mask, dtype=torch.long), torch.tensor(batch_token_type_ids, dtype=torch.long), torch.tensor(batch_pos, d

```

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 train_loader = DataLoader(train, batch_size=8, collate_fn=vectorize_batch, shuffle=True)
6 test_loader = DataLoader(test, batch_size=1, collate_fn=vectorize_batch)

```

▼ 4.5 POS Tagging and NER using CoNLL-2003 dataset [10pts]

Run the below cell to classify the conll2003 dataset for POS Tagging and NER using the sequenceLabeling functions that you have already implemented in 4.

You must reach an F1-score of >= 90% in both tasks to receive credit for this section.

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4

```

```

5 # Use cuda if present
6 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
7 print("Device available for running: ", device)

    Device available for running:  cuda

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from sequenceLabeling import SequenceLabeling
6 from tqdm import tqdm
7
8 NUM_CLASSES_POS = 47
9 NUM_CLASSES_NER = 9
10
11 model_pos = SequenceLabeling(NUM_CLASSES_POS)
12 model_pos.to(device)
13
14 model_ner = SequenceLabeling(NUM_CLASSES_NER)
15 model_ner.to(device)
16
17 criterion = nn.CrossEntropyLoss()
18 optimizer_pos = optim.Adam(model_pos.parameters(), lr=0.0001)
19 optimizer_ner = optim.Adam(model_ner.parameters(), lr=0.0001)
20 N_EPOCHS = 3 if not SHORT else 1
21
22 model_pos.train()
23 model_ner.train()
24
25 for epoch in range(N_EPOCHS):
26     total_loss_pos = 0.0
27     total_loss_ner = 0.0
28     for input_ids, mask, token_type_ids, target_pos, target_ner in tqdm(train_loader):
29         input_ids = input_ids.to(device)
30         mask = mask.to(device)
31         token_type_ids = token_type_ids.to(device)
32         target_pos = target_pos.to(device)
33         target_ner = target_ner.to(device)
34
35         outputs_pos = model_pos(input_ids, mask, token_type_ids)
36         outputs_ner = model_ner(input_ids, mask, token_type_ids)
37
38         active_loss_pos = mask.view(-1) == 1
39         active_logits_pos = outputs_pos.view(-1, NUM_CLASSES_POS)
40
41         active_labels_pos = torch.where(
42             active_loss_pos,
43             target_pos.view(-1),
44             torch.tensor(criterion.ignore_index).type_as(target_pos)
45         )
46
47         loss_pos = criterion(active_logits_pos, active_labels_pos)
48
49         optimizer_pos.zero_grad()
50         loss_pos.backward()
51         optimizer_pos.step()
52
53         active_loss_ner = mask.view(-1) == 1
54         active_logits_ner = outputs_ner.view(-1, NUM_CLASSES_NER)
55
56         active_labels_ner = torch.where(
57             active_loss_ner,
58             target_ner.view(-1),
59             torch.tensor(criterion.ignore_index).type_as(target_ner)
60         )
61
62         loss_ner = criterion(active_logits_ner, active_labels_ner)
63
64         optimizer_ner.zero_grad()
65         loss_ner.backward()
66         optimizer_ner.step()
67
68         total_loss_pos += loss_pos.item()
69         total_loss_ner += loss_ner.item()
70
71         if SHORT:
72             break
73
74     print("POS Tagging loss on epoch %i: %f" % (epoch, total_loss_pos))
75     print("NER loss on epoch %i: %f" % (epoch, total_loss_pos))

100%|██████████| 1756/1756 [11:40<00:00, 2.51it/s]
POS Tagging loss on epoch 0: 583.109677
NER loss on epoch 0: 583.109677
100%|██████████| 1756/1756 [11:42<00:00, 2.50it/s]
POS Tagging loss on epoch 1: 334.972757
NER loss on epoch 1: 334.972757
100%|██████████| 1756/1756 [11:48<00:00, 2.48it/s]POS Tagging loss on epoch 2: 284.867697
NER loss on epoch 2: 284.867697

```

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 with torch.no_grad():
6     y_true_pos = []
7     y_pred_pos = []
8     y_true_ner = []
9     y_pred_ner = []
10    for input_ids, mask, token_type_ids, target_pos, target_ner in tqdm(test_loader):
11        input_ids = input_ids.to(device)
12        mask = mask.to(device)
13        token_type_ids = token_type_ids.to(device)
14        target_pos = target_pos.to(device).view(-1)
15        target_ner = target_ner.to(device).view(-1)
16
17        outputs_pos = model_pos(input_ids, mask, token_type_ids)
18        predicted_pos = torch.argmax(outputs_pos, dim=-1)

```



```

19
20     outputs_ner = model_ner(input_ids, mask, token_type_ids)
21     predicted_ner = torch.argmax(outputs_ner, dim=-1)
22
23     active_loss = mask == 1
24     active_loss = active_loss.view(-1)
25     predicted_pos = predicted_pos.view(-1)
26     predicted_ner = predicted_ner.view(-1)
27
28     for i in range(len(active_loss)):
29         if not active_loss[i]:
30             break
31         y_true_pos.append(target_pos[i].cpu().detach().numpy())
32         y_pred_pos.append(predicted_pos[i].cpu().detach().numpy())
33         y_true_ner.append(target_ner[i].cpu().detach().numpy())
34         y_pred_ner.append(predicted_ner[i].cpu().detach().numpy())
35
36     if SHORT:
37         break
38 from sklearn.metrics import f1_score
39
40 print("Test F1-score on Conll2003 Dataset for POS Tagging : {:.3f}".format(f1_score(y_true_pos, y_pred_pos, average='micro')))
41 print("Test F1-score on Conll2003 Dataset for NER : {:.3f}".format(f1_score(y_true_ner, y_pred_ner, average='micro')))

```

100% | ██████████ | 3453/3453 [01:58<00:00, 29.26it/s]
Test F1-score on Conll2003 Dataset for POS Tagging : 0.920
Test F1-score on Conll2003 Dataset for NER : 0.967

📁 📌 Submit these files to Gradescope

Run the cell below to save the results of your model. You will be required to upload the `bert_pos.pkl` and `bert_ner.pkl` files that are saved in the `outputs` folder to Gradescope for accuracy evaluation. You must reach an accuracy of $\geq 90\%$ on both tasks to receive credit for this section.

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 with open('./outputs/bert_pos.pkl', 'wb') as fp:
6     pickle.dump(np.array(y_pred_pos), fp)
7
8 with open('./outputs/bert_ner.pkl', 'wb') as fp:
9     pickle.dump(np.array(y_pred_ner), fp)

```

Q5: Topic Modeling [20pts]

Topic Models are a type of statistical language models used for uncovering hidden structure in a collection of texts. There are several existing algorithms you can use to perform the topic modeling. In this HW, we will be exploring LDA (Latent Dirichlet Allocation) method. For more details please refer to class lectures and slides.

5.1: Latent Dirichlet Allocation [20pts]

In the `lda.py` file complete the following functions:

- `tokenize_words`
- `remove_stopwords`
- `create_dictionary`
- `build_LDAModel`

Latent Dirichlet Allocation (LDA) is a generative statistical model that explains a set of observations through unobserved groups, and each group explains why some parts of the data are similar. The LDA is an example of a topic model. In this, observations (e.g., words) are collected into documents, and each word's presence is attributable to one of the document's topics. Each document will contain a small number of topics.

5.2: Local Tests for LDA Functions [No Points]

You may test your implementation of the LDA functions contained in `lda.py` in the cell below. Feel free to comment out tests for functions that have not been completed yet. See [Using the Local Tests](#) for more details.

```

1 # #####
2 # ### DO NOT CHANGE THIS CELL ###
3 # #####
4
5 from lda import LDA
6 from local_tests.lda_test import LDA_Test
7
8 local_test = LDA_Test()
9
10 sample_data = local_test.input_sequences
11 stop_words = stopwords.words('english')
12 lda = LDA()
13
14 print('Local Tests for LDA Functions \n')
15
16 # Local test for tokenize_words
17 output = list(lda.tokenize_words(sample_data))
18 tokens_test = (output == local_test.output_tokenized)
19 print('Your tokenize_words works as expected:', tokens_test)
20
21 # Local test for remove_stopwords
22 filtered_output = lda.remove_stopwords(local_test.output_tokenized, stop_words)
23 remove_stopwords_test = (filtered_output == local_test.output_stopword_removal)
24 print('Your remove_stopwords works as expected:', remove_stopwords_test)
25
26 # Local test for create dictionary
27 id2word_test, corpus_test = lda.create_dictionary(local_test.output_stopword_removal)
28 id2word_corpus_test = ((id2word_test == local_test.output_id2word_dictionary) and (corpus_test == local_test.output_sample_corpus))
29 print('Your create_dictionary works as expected:', id2word_corpus_test)
30

```

```

31 # Instantiate build_LDAModel
32 lda_model_build = lda.build_LDAModel(id2word_test, corpus_test)
33 lda_model_test = isinstance(lda_model_build, gensim.models.LdaMulticore)
34 print("Your build_LDAModel works as expected:", lda_model_test)

WARNING:gensim.models.ldamulticore:too few updates, training might not converge; consider increasing the number of passes or iterations to improve accuracy
Local Tests for LDA Functions

Your tokenize_words works as expected: True
Your remove_stopwords works as expected: True
Your create_dictionary works as expected: True
Your build_LDAModel works as expected: True

```

5.3: Topic Modeling on Clickbait dataset using LDA [No Points]

Run the below cell to build LDA Model and visualize topics using gensim library.

```

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 from lda import LDA
6
7 stop_words = stopwords.words('english')
8 lda = LDA()
9
10 data = list(lda.tokenize_words(x_train))
11 data = lda.remove_stopwords(data, stop_words)
12 id2word, corpus = lda.create_dictionary(data)
13 lda_model = lda.build_LDAModel(id2word, corpus)

WARNING:gensim.models.ldamulticore:too few updates, training might not converge; consider increasing the number of passes or iterations to improve accuracy

1 #####
2 ### DO NOT CHANGE THIS CELL ###
3 #####
4
5 # Visualize the topics
6 pyLDavis.enable_notebook()
7
8 LDavis_prepared = pyLDavis.gensim_models.prepare(lda_model, corpus, id2word)
9 LDavis_prepared

```

Selected Topic:

Slide to adjust relevance metric:⁽²⁾
 $\lambda = 1$

