

```
1 # Licensing Information: You are free to use or extend this project for
2 # educational purposes provided that (1) you do not distribute or publish
3 # solutions, (2) you retain this notice, and (3) you provide clear
4 # attribution to The Georgia Institute of Technology, including a link to https://aritter.github.io/CS-7650/
5
6 # Attribution Information:
7 # This Project was developed at the Georgia Institute of Technology by Ashutosh Baheti (ashutosh.baheti@cc.gatech.edu),
8 # adapted from the Neural Machine Translation Project (Project 2)
9 # of the UC Berkeley NLP course https://cal-cs288.github.io/sp20/
```

▼ Project #3: Neural Chatbot

Welcome to the third and final programming assignment for CS 4650!

Neural Dialog Model are Sequence-to-Sequence (Seq2Seq) models that produce conversational response given the dialog history. State-of-the-art dialog models are trained on millions of multi-turn conversations. However, in this assignment we will narrow our scope to single turn conversations to make the problem easier.

In this assignment you will implement,

1. Seq2Seq encoder-decoder model
2. Seq2Seq model with attention mechanism
3. Greedy and Beam search decoding algorithms
4. Fine-tune and Evaluate BERT on disaster tweets

▼ Part 0: Setup

First, we'll import the various libraries needed for this project and define some of the utility functions to help with loading and manipulating the dataset. Since you've had experience in the previous project with splitting and tokenizing the dataset this is done for you in this project.

First import libraries required for the implementation

```
1 from __future__ import absolute_import
2 from __future__ import division
3 from __future__ import print_function
4 from __future__ import unicode_literals
5
6 import torch
7 from torch.jit import script, trace
8 import torch.nn as nn
9 from torch import optim
10 import torch.nn.functional as F
11 import numpy as np
12 import csv
13 import random
14 import re
15 import os
16 import unicodedata
17 import codecs
18 from io import open
19 import itertools
20 import math
21 import pickle
22 import statistics
23
24 from torch.utils.data import Dataset, DataLoader
25 from torch.nn.utils.rnn import pad_sequence
26 import tqdm
27 import nltk
28 from google.colab import files
```

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

```
Mounted at /content/drive
```

Then we implement some standard util functions that will be useful in the rest of the code.

```
1 # General util functions
2 def make_dir_if_not_exists(directory):
3     if not os.path.exists(directory):
4         logging.info("Creating new directory: {}".format(directory))
5         os.makedirs(directory)
6
7 def print_list(l, K=None):
8     # If K is given then only print first K
9     for i, e in enumerate(l):
10         if i == K:
11             break
```

```

12     print(e)
13     print()
14
15 def remove_multiple_spaces(string):
16     return re.sub(r'\s+', ' ', string).strip()
17
18 def save_in_pickle(save_object, save_file):
19     with open(save_file, "wb") as pickle_out:
20         pickle.dump(save_object, pickle_out)
21
22 def load_from_pickle(pickle_file):
23     with open(pickle_file, "rb") as pickle_in:
24         return pickle.load(pickle_in)
25
26 def save_in_txt(list_of_strings, save_file):
27     with open(save_file, "w") as writer:
28         for line in list_of_strings:
29             line = line.strip()
30             writer.write(f"{line}\n")
31
32 def load_from_txt(txt_file):
33     with open(txt_file, "r") as reader:
34         all_lines = list()
35         for line in reader:
36             line = line.strip()
37             all_lines.append(line)
38     return all_lines

```

Finally we will check if GPU is available and set the device accordingly.

Tip: While debugging use `CPU` to get clearer stack traces and change the runtime type to `GPU` when you are ready to train your models efficiently

```

1 print(torch.cuda.is_available())
2 if torch.cuda.is_available():
3     device = torch.device("cuda")
4 else:
5     device = torch.device("cpu")
6 print("Using device:", device)

True
Using device: cuda

```

▼ Dataset

For the dataset we will be using a small sample of single turn input and response pairs from [Cornell Movie Dialog Corpus](#). We filter conversational pairs with sentences > 10 tokens. We have already created a sample of tokenized, lowercased single turn conversations from Cornell Movie Dialog Corpus. The preprocessed dataset sample is stored in pickle format and can be downloaded from [this link](#). Please download the `processed_CMDC.pkl` file from the link and upload it in colab.

```

1 # Loading the pre-processed conversational exchanges (source-target pairs) from pickle data files
2 all_conversations = load_from_pickle('/content/drive/MyDrive/processed_CMDC.pkl')
3 # Extract 100 conversations from the end for evaluation and keep the rest for training
4 eval_conversations = all_conversations[-100:]
5 all_conversations = all_conversations[:-100]
6
7 # Logging data stats
8 print(f"Number of Training Conversation Pairs = {len(all_conversations)}")
9 print(f"Number of Evaluation Conversation Pairs = {len(eval_conversations)}")

Number of Training Conversation Pairs = 53065
Number of Evaluation Conversation Pairs = 100

```

Let's print a couple of conversations to check if they are loaded properly.

```

1 print_list(all_conversations, 5)

('there .', 'where ?')
('you have my word . as a gentleman', 'you re sweet .')
('hi .', 'looks like things worked out tonight huh ?')
('have fun tonight ?', 'tons')
('well no . . .', 'then that s all you had to say .')

```

▼ Vocabulary

The words in the sentences need to be converted into integer tokens so that the neural model can operate on them. For this purpose, we will create a vocabulary which will convert the input strings into model recognizable integer tokens.

```

1 pad_word = "<pad>"
2 bos_word = "<s>"

```

```

3 eos_word = "</s>"
4 unk_word = "<unk>"
5 pad_id = 0
6 bos_id = 1
7 eos_id = 2
8 unk_id = 3
9
10 def normalize_sentence(s):
11     s = re.sub(r"([!?!])", r" \1", s)
12     s = re.sub(r"^[a-zA-Z.!?!]+", r" ", s)
13     s = re.sub(r"\s+", r" ", s).strip()
14     return s
15
16 class Vocabulary:
17     def __init__(self):
18         self.word_to_id = {pad_word: pad_id, bos_word: bos_id, eos_word: eos_id, unk_word: unk_id}
19         self.word_count = {}
20         self.id_to_word = {pad_id: pad_word, bos_id: bos_word, eos_id: eos_word, unk_id: unk_word}
21         self.num_words = 4
22
23     def get_ids_from_sentence(self, sentence):
24         sentence = normalize_sentence(sentence)
25         sent_ids = [bos_id] + [self.word_to_id[word] if word in self.word_to_id \
26                               else unk_id for word in sentence.split()] + \
27                     [eos_id]
28         return sent_ids
29
30     def tokenized_sentence(self, sentence):
31         sent_ids = self.get_ids_from_sentence(sentence)
32         return [self.id_to_word[word_id] for word_id in sent_ids]
33
34     def decode_sentence_from_ids(self, sent_ids):
35         words = list()
36         for i, word_id in enumerate(sent_ids):
37             if word_id in [bos_id, eos_id, pad_id]:
38                 # Skip these words
39                 continue
40             else:
41                 words.append(self.id_to_word[word_id])
42         return ' '.join(words)
43
44     def add_words_from_sentence(self, sentence):
45         sentence = normalize_sentence(sentence)
46         for word in sentence.split():
47             if word not in self.word_to_id:
48                 # add this word to the vocabulary
49                 self.word_to_id[word] = self.num_words
50                 self.id_to_word[self.num_words] = word
51                 self.word_count[word] = 1
52                 self.num_words += 1
53             else:
54                 # update the word count
55                 self.word_count[word] += 1
56
57 vocab = Vocabulary()
58 for src, tgt in all_conversations:
59     vocab.add_words_from_sentence(src)
60     vocab.add_words_from_sentence(tgt)
61 print(f"Total words in the vocabulary = {vocab.num_words}")

```

Total words in the vocabulary = 7727

Let's print the top 30 vocab words:

```

1 print_list(sorted(vocab.word_count.items(), key=lambda item: item[1], reverse=True), 30)

```

('.', 84255)
 ('?', 36822)
 ('you', 25093)
 ('i', 18946)
 ('what', 10765)
 ('s', 10089)
 ('it', 9668)
 ('!', 8872)
 ('the', 8011)
 ('t', 7411)
 ('to', 6929)
 ('a', 6582)
 ('that', 5992)
 ('no', 4931)
 ('me', 4839)
 ('do', 4745)
 ('is', 4434)
 ('don', 3577)
 ('are', 3503)
 ('he', 3413)
 ('yes', 3384)
 ('m', 3382)
 ('not', 3252)
 ('we', 3252)
 ('know', 3171)
 ('re', 2965)

```
( 'your', 2809)
( 'this', 2726)
( 'yeah', 2708)
( 'in', 2678)
```

We can also print a couple of sentences to verify that the vocabulary is working as intended, as well as ensure our encoding/decoding process works as expected.

```
1 for src, tgt in all_conversations[:3]:
2     sentence = tgt
3     word_tokens = vocab.tokenized_sentence(sentence)
4     # Automatically adds bos_id and eos_id before and after sentence ids respectively
5     word_ids = vocab.get_ids_from_sentence(sentence)
6     print(sentence)
7     print(word_tokens)
8     print(word_ids)
9     print(vocab.decode_sentence_from_ids(word_ids))
10    print()
11
12 word = "the"
13 word_id = vocab.word_to_id[word]
14 print(f"Word = {word}")
15 print(f"Word ID = {word_id}")
16 print(f"Word decoded from ID = {vocab.decode_sentence_from_ids([word_id])}")

where ?
['<s>', 'where', '?', '</s>']
[1, 6, 7, 2]
where ?

you re sweet .
['<s>', 'you', 're', 'sweet', '.', '</s>']
[1, 8, 15, 16, 5, 2]
you re sweet .

looks like things worked out tonight huh ?
['<s>', 'looks', 'like', 'things', 'worked', 'out', 'tonight', 'huh', '?', '</s>']
[1, 18, 19, 20, 21, 22, 23, 24, 7, 2]
looks like things worked out tonight huh ?

Word = the
Word ID = 47
Word decoded from ID = the
```

▼ Part 1: Dataset Preparation (5 points)

We will use built-in dataset utilities, `torch.utils.data.Dataset` and `torch.utils.data.DataLoader`, to get batched data readily useful for training like what you saw in Project 1.

Most of the dataset has been filled out for you, however the `collate_fn` needs to be finished.

```
1 class SingleTurnMovieDialog_dataset(Dataset):
2     """Single-Turn version of Cornell Movie Dialog Crops dataset."""
3
4     def __init__(self, conversations, vocab, device):
5         """
6         Args:
7             conversations: list of tuple (src_string, tgt_string)
8                           - src_string: String of the source sentence
9                           - tgt_string: String of the target sentence
10            vocab: Vocabulary object that contains the mapping of
11                  words to indices
12            device: cpu or cuda
13        """
14        self.conversations = conversations
15        self.vocab = vocab
16        self.device = device
17
18        def encode(src, tgt):
19            src_ids = self.vocab.get_ids_from_sentence(src)
20            tgt_ids = self.vocab.get_ids_from_sentence(tgt)
21            return (src_ids, tgt_ids)
22
23        # We will pre-tokenize the conversations and save in id lists for later use
24        self.tokenized_conversations = [encode(src, tgt) for src, tgt in self.conversations]
25
26        def __len__(self):
27            return len(self.conversations)
28
29        def __getitem__(self, idx):
30            if torch.is_tensor(idx):
31                idx = idx.tolist()
32
33            return {"conv_ids":self.tokenized_conversations[idx], "conv":self.conversations[idx]}
34
35 def collate_fn(data):
36     """Creates mini-batch tensors from the list of tuples (src_seq, trg_seq).
```

```

37 We should build a custom collate_fn rather than using default collate_fn,
38 because merging sequences (including padding) is not supported in default.
39 Sequences are padded to the maximum length of mini-batch sequences (dynamic padding).
40 Args:
41     data: list of dicts {"conv_ids":(src_ids, tgt_ids), "conv":(src_str, trg_str)}.
42         - src_ids: list of src piece ids; variable length.
43         - tgt_ids: list of tgt piece ids; variable length.
44         - src_str: String of src
45         - tgt_str: String of tgt
46 Returns: dict { "conv_ids":      (src_ids, tgt_ids),
47                "conv":          (src_str, tgt_str),
48                "conv_tensors":  (src_seqs, tgt_seqs)}
49     src_seqs: torch tensor of shape (src_padded_length, batch_size).
50     tgt_seqs: torch tensor of shape (tgt_padded_length, batch_size).
51     src_padded_length = length of the longest src sequence from src_ids
52     tgt_padded_length = length of the longest tgt sequence from tgt_ids
53 """
54 # Sort conv_ids based on decreasing order of the src_lengths.
55 # This is required for efficient GPU computations.
56 src_ids = [torch.LongTensor(e["conv_ids"][0]) for e in data]
57 tgt_ids = [torch.LongTensor(e["conv_ids"][1]) for e in data]
58 src_str = [e["conv"][0] for e in data]
59 tgt_str = [e["conv"][1] for e in data]
60 data = list(zip(src_ids, tgt_ids, src_str, tgt_str))
61 data.sort(key=lambda x: len(x[0]), reverse=True)
62 src_ids, tgt_ids, src_str, tgt_str = zip(*data)
63
64 # Pad the src_ids and tgt_ids using token pad_id to create src_seqs and tgt_seqs
65
66 # HINT: You can use the nn.utils.rnn.pad_sequence utility
67 # function to combine a list of variable-length sequences with padding.
68
69 # YOUR CODE HERE
70
71 src_seqs = nn.utils.rnn.pad_sequence(src_ids, padding_value = pad_id)
72 tgt_seqs = nn.utils.rnn.pad_sequence(tgt_ids, padding_value = pad_id)
73
74 return {"conv_ids":(src_ids, tgt_ids), "conv":(src_str, tgt_str), "conv_tensors":(src_seqs.to(device), tgt_seqs.to(device))}

1 # Create the DataLoader for all_conversations
2 dataset = SingleTurnMovieDialog_dataset(all_conversations, vocab, device)
3
4 batch_size = 5
5
6 data_loader = DataLoader(dataset=dataset, batch_size=batch_size,
7                           shuffle=True, collate_fn=collate_fn)

```

Let's test a batch of data to make sure everything is working as intended

HINT: If you've padded the targets correctly, each column should start with the beginning of sequence ID (i.e. 1) and should follow the end of sequence ID with some number of the pad ID (i.e. 0) if the sequence in that column is shorter than the max in the minibatch.

```

1 # Test one batch of training data
2 first_batch = next(iter(data_loader))
3 print(f"Testing first training batch of size {len(first_batch['conv'][0])}")
4 print(f"List of source strings:")
5 print_list(first_batch["conv"][0])
6 print(f"Tokenized source ids:")
7 print_list(first_batch["conv_ids"][0])
8 print(f"Padded source ids as tensor (shape {first_batch['conv_tensors'][0].size()}:")
9 print(first_batch["conv_tensors"][0])

```

```

Testing first training batch of size 5
List of source strings:
i don t know i think . . .
but i get a red light .
i will do no such thing .
pardon ?
okay

```

```

Tokenized source ids:
tensor([ 1, 54, 198, 103, 97, 54, 66, 5, 5, 5, 2])
tensor([ 1, 36, 54, 142, 13, 667, 1492, 5, 2])
tensor([ 1, 54, 728, 41, 28, 590, 169, 5, 2])
tensor([ 1, 1274, 7, 2])
tensor([ 1, 53, 2])

```

```

Padded source ids as tensor (shape torch.Size([11, 5])):
tensor([[ 1, 1, 1, 1, 1],
        [ 54, 36, 54, 1274, 53],
        [ 198, 54, 728, 7, 2],
        [ 103, 142, 41, 2, 0],
        [ 97, 13, 28, 0, 0],
        [ 54, 667, 590, 0, 0],
        [ 66, 1492, 169, 0, 0],
        [ 5, 5, 5, 0, 0],
        [ 5, 2, 2, 0, 0],
        [ 5, 0, 0, 0, 0],
        [ 2, 0, 0, 0, 0]], device='cuda:0')

```

▼ Part 2: Baseline Seq2Seq model (25 points)

In this section you will initialize the layers needed for your Seq2Seq model, define the encode and decode functions of your model, and define a loss function to handle the padded tokens when training your model.

With the training `Dataset` and `DataLoader` ready, we can implement our Seq2Seq baseline model.

The model will consist of

1. Shared embedding layer between encoder and decoder that converts the input sequence of word ids to dense embedding representations
2. Bidirectional GRU encoder that encodes the embedded source sequence into hidden representation
3. Unidirectional GRU decoder that predicts target sequence using final encoder hidden representation

```
1 class Seq2seqBaseline(nn.Module):
2     def __init__(self, vocab, emb_dim = 300, hidden_dim = 300, num_layers = 2, dropout=0.1):
3         super().__init__()
4
5         # Initialize your model's parameters here. To get started, we suggest
6         # setting all embedding and hidden dimensions to 300, using encoder and
7         # decoder GRUs with 2 layers each, and using a dropout rate of 0.1.
8
9         # HINT: To create a bidirectional GRU, you don't need to create two GRU
10        # networks, instead use the bidirectional flag when initializing the layer.
11
12        self.num_words = num_words = vocab.num_words
13        self.emb_dim = emb_dim
14        self.hidden_dim = hidden_dim
15        self.num_layers = num_layers
16        # YOUR CODE HERE
17        self.embedding = nn.Embedding(num_words, emb_dim)
18        self.encoder = nn.GRU(emb_dim, hidden_dim, num_layers, dropout = dropout, bidirectional = True)
19        self.decoder = nn.GRU(emb_dim, hidden_dim, num_layers, dropout = dropout, bidirectional = False)
20        self.fc = nn.Linear(hidden_dim, num_words)
21
22    def encode(self, source):
23        """Encode the source batch using a bidirectional GRU encoder.
24
25        Args:
26            source: An integer tensor with shape (max_src_sequence_length,
27                  batch_size) containing subword indices for the source sentences.
28
29        Returns:
30            A tuple with three elements:
31                encoder_output: The output hidden representation of the encoder
32                               with shape (max_src_sequence_length, batch_size, hidden_size).
33                               Can be obtained by adding the hidden representations of both
34                               directions of the encoder bidirectional GRU.
35                encoder_mask: A boolean tensor with shape (max_src_sequence_length,
36                    batch_size) indicating which encoder outputs correspond to padding
37                    tokens. Its elements should be True at positions corresponding to
38                    padding tokens and False elsewhere.
39                encoder_hidden: The final hidden states of the bidirectional GRU
40                               (after a suitable projection) that will be used to initialize
41                               the decoder. This should be a tensor h_n with shape
42                               (num_layers, batch_size, hidden_size). Note that the hidden
43                               state returned by the bi-GRU cannot be used directly. Its
44                               initial dimension is twice the required size because it
45                               contains state from two directions.
46
47        The first two return values are not required for the baseline model and will
48        only be used later in the attention model. If desired, they can be replaced
49        with None for the initial implementation.
50        """
51
52        # Implementation tip: consider using packed sequences to more easily work
53        # with the variable-length sequences represented by the source tensor.
54        # See https://pytorch.org/docs/stable/nn.html#torch.nn.utils.rnn.PackedSequence.
55
56        # https://stackoverflow.com/questions/51030782/why-do-we-pack-the-sequences-in-pytorch
57
58        # HINT: there are many simple ways to combine the forward
59        # and backward portions of the final hidden state, e.g. addition, averaging,
60        # or a linear transformation of the appropriate size. Any of these
61        # should let you reach the required performance.
62
63        # Compute a tensor containing the length of each source sequence.
64        source_lengths = torch.sum(source != pad_id, axis=0).cpu()
65
66        # YOUR CODE HERE
67        embd = self.embedding(source).to(device)
68        packed_source = torch.nn.utils.rnn.pack_padded_sequence(embd, source_lengths)
69        packed_encoder_output, hn = self.encoder(packed_source)
70        encoder_output, _ = nn.utils.rnn.pad_packed_sequence(packed_encoder_output)
71        new_encoder_output = torch.zeros(encoder_output.size()[0], encoder_output.size()[1], self.hidden_dim).to(device)
72        for i in range(self.hidden_dim):
73            new_encoder_output[:, :, i] = encoder_output[:, :, 2*i] + encoder_output[:, :, 2*i+1]
74        encoder_output = new_encoder_output
```

```

75
76 # encoder_hidden = encoder_hidden.view(self.num_layers, 2, -1, self.hidden_dim)
77 # encoder_1 = (encoder_hidden[0, :, :] + encoder_hidden[1, :, :]).unsqueeze(0)
78 # encoder_2 = (encoder_hidden[2, :, :] + encoder_hidden[3, :, :]).unsqueeze(0)
79 # encoder_hidden = torch.cat((encoder_1, encoder_2), 0)
80 encoder_hidden = torch.zeros(self.num_layers, hn.size()[1], self.hidden_dim).to(device)
81 for i in range(self.num_layers):
82     encoder_hidden[i, :, :] = hn[2*i, :, :] + hn[2*i+1, :, :]
83 # print(encoder_hidden.shape)
84 # encoder_output = encoder_output[:, :, :self.hidden_dim] + encoder_output[:, :, self.hidden_dim:]
85 encoder_mask = (source == pad_id)
86 return encoder_output, encoder_mask, encoder_hidden
87
88 def decode(self, decoder_input, last_hidden, encoder_output, encoder_mask):
89     """Run the decoder GRU for one decoding step from the last hidden state.
90
91     The third and fourth arguments are not used in the baseline model, but are
92     included for compatibility with the attention model in the next section.
93
94     Args:
95         decoder_input: An integer tensor with shape (1, batch_size) containing
96             the subword indices for the current decoder input.
97         last_hidden: A tensor h_{t-1} representing the last hidden
98             state of the decoder, has the shape (num_layers, batch_size,
99             hidden_size). For the first decoding step the last_hidden will be
100             encoder's final hidden representation.
101         encoder_output: The output of the encoder with shape
102             (max_src_sequence_length, batch_size, hidden_size).
103         encoder_mask: The output mask from the encoder with shape
104             (max_src_sequence_length, batch_size). Encoder outputs at positions
105             with a True value correspond to padding tokens and should be ignored.
106
107     Returns:
108         A tuple with three elements:
109             logits: A tensor with shape (batch_size,
110                 vocab_size) containing unnormalized scores for the next-word
111                 predictions at each position.
112             decoder_hidden: tensor h_n with the same shape as last_hidden
113                 representing the updated decoder state after processing the
114                 decoder input.
115             attention_weights: This will be implemented later in the attention
116                 model, but in order to maintain compatible type signatures, we also
117                 include it here. This can be None or any other placeholder value.
118     """
119
120     # These arguments are not used in the baseline model.
121     del encoder_output
122     del encoder_mask
123     # YOUR CODE HERE
124     embd = self.embedding(decoder_input).to(device)
125     output, decoder_hidden = self.decoder(embd, last_hidden)
126     logits = self.fc(output.squeeze())
127     return logits, decoder_hidden, None
128
129
130 def compute_loss(self, source, target):
131     """Run the model on the source and compute the loss on the target.
132     The loss for this project should use teacher forcing, where the
133     output of the model is used only to compute loss and not passed
134     back in to get the next predicted token.
135
136     Args:
137         source: An integer tensor with shape (max_source_sequence_length,
138             batch_size) containing subword indices for the source sentences.
139         target: An integer tensor with shape (max_target_sequence_length,
140             batch_size) containing subword indices for the target sentences.
141
142     Returns:
143         A scalar float tensor representing cross-entropy loss on the current batch
144         divided by the number of target tokens in the batch.
145         Many of the target tokens will be pad tokens. You should mask the loss
146         from these tokens using appropriate mask on the target tokens loss.
147     """
148
149     # Hint: don't feed the target tensor directly to the decoder.
150     # To see why, note that for a target sequence like <s> A B C </s>, you would
151     # want to run the decoder on the prefix <s> A B C and have it predict the
152     # suffix A B C </s>.
153
154     # You may run self.encode() on the source only once and decode the target
155     # one step at a time.
156
157     total_loss = 0
158     # YOUR CODE HERE
159     batch_size = source.shape[1]
160     encoder_output, encoder_mask, encoder_hidden = self.encode(source)
161     decoder_hidden = encoder_hidden
162     decoder_input = target[0, :].unsqueeze(0)
163
164     # decoder_input = target[0].unsqueeze(0)
165     for i in range(target.shape[0]-1):

```

```

166         # print(target[i].shape)
167         logits, decoder_hidden, _ = self.decode(decoder_input, decoder_hidden, encoder_output, encoder_mask)
168         # print(decoder_hidden.shape)
169         # print(decoder_input.shape)
170         loss = nn.functional.cross_entropy(logits, target[i+1])
171         mask = (target[i+1] != pad_id)
172         # print(mask)
173         total_loss += (loss*mask).sum()
174         decoder_input = target[i+1, :].unsqueeze(0)
175
176     return total_loss / torch.sum(target != pad_id).float()
177

```

▼ Training

We provide a training loop for training the model. You are welcome to modify the training loop by adjusting the learning rate or changing optimization settings.

Important: During our testing we found that training the encoder and decoder with different learning rates is crucial for getting good performance over the small dialog corpus. Specifically, the decoder parameter learning rate should be 5 times the encoder parameter learning rate. Hence, add the encoder parameter variable names in the `encoder_parameter_names` as a list. For example, if encoder is using `self.embedding_layer` and `self.encoder_gru` layer then the `encoder_parameter_names` should be `['embedding_layer', 'encoder_gru']`

```

1 from tqdm.notebook import trange, tqdm
2 def train(model, data_loader, num_epochs, model_file):
3     """Train the model for given number of epochs and save the trained model in
4     the final model_file.
5     """
6
7     # feel free to edit these values!
8     decoder_learning_ratio = 5.0
9     learning_rate = 0.0001
10
11     encoder_parameter_names = ['embedding', 'encoder']
12
13     encoder_named_params = list(filter(lambda kv: any(key in kv[0] for key in encoder_parameter_names), model.named_parameters()))
14     decoder_named_params = list(filter(lambda kv: not any(key in kv[0] for key in encoder_parameter_names), model.named_parameters()))
15     encoder_params = [e[1] for e in encoder_named_params]
16     decoder_params = [e[1] for e in decoder_named_params]
17     optimizer = torch.optim.AdamW({'params': encoder_params,
18                                     {'params': decoder_params, 'lr': learning_rate * decoder_learning_ratio}}, lr=learning_rate)
19
20     clip = 50.0
21     for epoch in trange(num_epochs, desc="training", unit="epoch"):
22         # print(f"Total training instances = {len(train_dataset)}")
23         # print(f"train_data_loader = {len(train_data_loader)} {1180 > len(train_data_loader)/20}")
24         with tqdm(
25             data_loader,
26             desc="epoch {}".format(epoch + 1),
27             unit="batch",
28             total=len(data_loader)) as batch_iterator:
29             model.train()
30             total_loss = 0.0
31             for i, batch_data in enumerate(batch_iterator, start=1):
32                 source, target = batch_data["conv_tensors"]
33                 optimizer.zero_grad()
34                 loss = model.compute_loss(source, target)
35                 total_loss += loss.item()
36                 loss.backward()
37                 # Gradient clipping before taking the step
38                 _ = nn.utils.clip_grad_norm_(model.parameters(), clip)
39                 optimizer.step()
40
41             batch_iterator.set_postfix(mean_loss=total_loss / i, current_loss=loss.item())
42     # Save the model after training
43     torch.save(model.state_dict(), model_file)

```

We can now train the baseline model. This should take about 5 minutes with a GPU and will take >40 minutes on just the CPU, so we highly recommend using a Colab Pro account.

A correct implementation should get a average train loss of < 3.00, however be aware, as this may not be the best sign your model will behave as desired. While the loss will give you some idea concerning the correctness of your implementation, you should also "talk" with it to confirm. Please check Piazza (specifically, the pinned post on Part 2) to see an example of a correct implementation.

The code will automatically save and download the model at the end of training, that way you won't have to retrain if you come back to the notebook later.

```

1 # You are welcome to adjust these parameters based on your model implementation.
2 num_epochs = 6
3 batch_size = 64
4 # Reloading the data_loader to increase batch_size
5 data_loader = DataLoader(dataset=dataset, batch_size=batch_size,
6                           shuffle=True, collate_fn=collate_fn)
7

```



```

8 baseline_model = Seq2seqBaseline(vocab).to(device)
9 train(baseline_model, data_loader, num_epochs, "baseline_model.pt")
10 # Download the trained model to local for future use
11 files.download('baseline_model.pt')

training: 100% 6/6 [04:17<00:00, 43.08s/epoch]

epoch 1: 100% 830/830 [00:42<00:00, 22.03batch/s, current_loss=1.79, mean_loss=2.99]

epoch 2: 100% 830/830 [00:42<00:00, 21.52batch/s, current_loss=2.43, mean_loss=2.58]

epoch 3: 100% 830/830 [00:43<00:00, 21.76batch/s, current_loss=2.79, mean_loss=2.44]

epoch 4: 100% 830/830 [00:42<00:00, 21.20batch/s, current_loss=2.79, mean_loss=2.33]

epoch 5: 100% 830/830 [00:43<00:00, 19.78batch/s, current_loss=2.54, mean_loss=2.23]

epoch 6: 100% 830/830 [00:43<00:00, 15.58batch/s, current_loss=2.04, mean_loss=2.14]

1 # Reload the model from the model file.
2 # Useful when you have already trained and saved the model
3 baseline_model = Seq2seqBaseline(vocab).to(device)
4 baseline_model.load_state_dict(torch.load("baseline_model.pt", map_location=device))

<All keys matched successfully>

```

▼ Part 3: Greedy Search (10 points)

For evaluation, we also need to be able to generate entire strings from the model. We'll first define a greedy inference procedure here. Later on, we'll implement beam search. *Hint:* Use the **normalize_sentence** and **vocab.get_ids_from_sentence** functions to prepare your input.

```

1 def predict_greedy(model, sentence, max_length=100):
2     """Make predictions for the given input using greedy inference.
3
4     Args:
5         model: A sequence-to-sequence model.
6         sentence: A input string.
7         max_length: The maximum length at which to truncate outputs in order to
8             avoid non-terminating inference.
9
10    Returns:
11        Model's predicted greedy response for the input, represented as string.
12
13    HINT: Make sure to terminate your models prediction when it outputs the end of
14    sequence ID, even if the models reponse hasn't reached the max length.
15    """
16
17    # You should make only one call to model.encode() at the start of the function,
18    # and make only one call to model.decode() per inference step.
19    model.eval()
20
21    # YOUR CODE HERE
22    model = model.to(device)
23    input = torch.tensor(vocab.get_ids_from_sentence(normalize_sentence(sentence))).unsqueeze(1).to(device)
24    encode_output, encode_mask, encode_hidden = model.encode(input)
25    decode_hidden = encode_hidden
26    output = []
27    decoder_input = torch.tensor([bos_id]).unsqueeze(0).to(device)
28    for i in range(max_length):
29        logits, decode_hidden, _ = model.decode(decoder_input, decode_hidden, encode_output, encode_mask)
30        predicted_id = int(logits.argmax())
31        output.append(predicted_id)
32        decoder_input = torch.tensor([predicted_id]).unsqueeze(0).to(device)
33    return vocab.decode_sentence_from_ids(output)
34

```

Let's chat interactively with our trained baseline Seq2Seq dialog model and save the generated conversations for submission (please make sure to keep the conversations in your submission ["PG-13"](#)). We will reuse the conversational inputs while testing Seq2Seq + Attention model.

The output of your model isn't likely to be very colorful given the simplicity of the dataset we're working on. Instead, you should expect responses that are generally grammatically correct and do not degrade (i.e. your model keeps repeating the same word(s) over and over).

IMPORTANT: FOR YOUR FINAL SUBMISSION TO GRADESCOPE, PLEASE "TALK" WITH YOUR CHATBOT IN THE CELLS BELOW FOR ABOUT FIVE TURNS AND MAKE SURE THE RESPONSES ARE VISIBLE IN YOUR UPLOADED NOTEBOOK.

Note: enter "q" or "quit" to end the interactive chat.

```

1 def chat_with_model(model, mode="greedy"):
2     if mode == "beam":
3         predict_f = predict_beam
4     else:
5         predict_f = predict_greedy
6     chat_log = list()
7     input_sentence = ''
8     while(1):
9         # Get input sentence

```

```

10     input_sentence = input('Input > ')
11     # Check if it is quit case
12     if input_sentence == 'q' or input_sentence == 'quit': break
13
14     generation = predict_f(model, input_sentence)
15     if mode == "beam":
16         generation = generation[0]
17     print('Greedy Response:', generation)
18     print()
19     chat_log.append((input_sentence, generation))
20     return chat_log

```

```

1 baseline_chat = chat_with_model(baseline_model)

```

```

Input > how are you?
Greedy Response: fine .

Input > what do you want for lunch?
Greedy Response: i m going to go .

Input > did you like our time together?
Greedy Response: i don t know .

Input > who is your favorite actor?
Greedy Response: i don t know .

Input > what do you want to eat for dinner?
Greedy Response: i m going to go .

Input > What's your favorite movie?
Greedy Response: yes .

Input > thanks.
Greedy Response: you re a good man .

Input > did you like our time together?
Greedy Response: i don t know .

Input > why not?
Greedy Response: i don t know .

Input > q

```

▼ Part 4: Seq2Seq + Attention Model (15 points)

Next, we extend the baseline model to include an attention mechanism in the decoder. With attention mechanism, the model doesn't need to encode the input into a fixed dimensional hidden representation. Rather, it creates a new context vector for each turn that is a weighted sum of encoder hidden representation.

Your implementation can use any attention mechanism to get weight distribution over the source words. One simple way to include attention in decoder goes as follows (reminder: the decoder processed one token at a time),

1. Process the current decoder_input through embedding layer and decoder GRU layer.
2. Use the current decoder token representation, d of shape $(1 * b * h)$ and encoder representation, e_1, \dots, e_n of shape $(n * b * h)$, where n is max_src_length after padding) to compute attention score matrix of shape $(b * n)$. There are multiple options to compute this score matrix. A few of such options are available in [the table provided in this blog](#). Please leave a comment in your code with the name of the method you choose to implement
3. Normalize the attention scores $(b * n)$ so that they sum up to 1.0 by taking a `softmax` over the second dimension.

After computing the normalized attention distribution, take a weighted sum of the encoder outputs to obtain the attention context

$c = \sum_i w_i e_i$, and add this to the decoder output d to obtain the final representation to be passed to the vocabulary projection layer (you may need another linear layer to make the sizes match before adding c and d).

```

1 class Seq2seqAttention(Seq2seqBaseline):
2     def __init__(self, vocab):
3         super().__init__(vocab)
4
5         # Initialize any additional parameters needed for this model that are not
6         # already included in the baseline model.
7
8         # YOUR CODE HERE
9         self.attention = nn.Linear(self.hidden_dim*2, self.hidden_dim)
10
11
12     def decode(self, decoder_input, last_hidden, encoder_output, encoder_mask):
13         """Run the decoder GRU for one decoding step from the last hidden state.
14
15         The third and fourth arguments are not used in the baseline model, but are
16         included for compatibility with the attention model in the next section.
17
18         Args:
19             decoder_input: An integer tensor with shape (1, batch_size) containing
20                 the subword indices for the current decoder input.
21             last_hidden: A pair of tensors h_{t-1} representing the last hidden
22                 state of the decoder, each with shape (num_layers, batch_size,
23                 hidden_size). For the first decoding step the last_hidden will be

```

```

24         encoder's final hidden representation.
25     encoder_output: The output of the encoder with shape
26         (max_src_sequence_length, batch_size, hidden_size).
27     encoder_mask: The output mask from the encoder with shape
28         (max_src_sequence_length, batch_size). Encoder outputs at positions
29         with a True value correspond to padding tokens and should be ignored.
30
31     Returns:
32         A tuple with three elements:
33         logits: A tensor with shape (batch_size, vocab_size)
34             containing unnormalized scores for the next-word
35             predictions at each position.
36         decoder_hidden: tensor h_n with the same shape as last_hidden
37             representing the updated decoder state after processing the
38             decoder input.
39         attention_weights: A tensor with shape (batch_size, max_src_sequence_length)
40             representing the normalized attention weights. This should sum to 1
41             along the last dimension.
42     """
43
44     # YOUR CODE HERE
45     batch_size = decoder_input.size(1)
46     embd = self.embedding(decoder_input)
47     decoder_out, decoder_hidden = self.decoder(embd, last_hidden)
48     top_hidden = decoder_hidden[-1].unsqueeze(1)
49     encoder_output = encoder_output.permute(1,0,2)
50
51     # use dot product scores
52     scores = torch.bmm(top_hidden, encoder_output.transpose(1,2))
53     attn_weights = nn.functional.softmax(scores, dim=1)
54     context = torch.bmm(attn_weights, encoder_output).squeeze(1)
55     decoder_output = decoder_out.squeeze(0)
56     joined = torch.cat((decoder_output, context), dim=1)
57     combined = self.attention(joined)
58     logits = self.fc(combined)
59     return logits, decoder_hidden, attn_weights
60
61
62

```

▼ Training

We can now train the attention model.

A correct implementation should also get an average train loss of < 3.00, however you should still check your models output to confirm you've implemented the attention mechanism correctly.

The code will automatically save and download the model at the end of training.

It may happen that the baseline model achieves a worse loss than attention model. This is because our dataset is very small and the attention model may be over parameterized for our toy dataset. Regardless, we would consider this as acceptable submission if the attention model generated responses look comparable to the baseline model.

```

1 # You are welcome to adjust these parameters based on your model implementation.
2 num_epochs = 8
3 batch_size = 64
4 data_loader = DataLoader(dataset=dataset, batch_size=batch_size,
5                           shuffle=True, collate_fn=collate_fn)
6
7 attention_model = Seq2seqAttention(vocab).to(device)
8 train(attention_model, data_loader, num_epochs, "attention_model.pt")
9 # Download the trained model to local for future use
10 files.download('attention_model.pt')

```

```

training: 100%                               8/8 [13:04<00:00, 97.24s/epoch]

epoch 1: 100%                                830/830 [01:40<00:00, 7.10batch/s, current_loss=2.39, mean_loss=2.92]

epoch 2: 100%                                830/830 [01:39<00:00, 8.67batch/s, current_loss=2.76, mean_loss=2.53]

epoch 3: 100%                                830/830 [01:38<00:00, 8.99batch/s, current_loss=2.37, mean_loss=2.36]

epoch 4: 100%                                830/830 [01:37<00:00, 9.60batch/s, current_loss=2.72, mean_loss=2.2]

epoch 5: 100%                                830/830 [01:37<00:00, 6.99batch/s, current_loss=1.77, mean_loss=2.05]

epoch 6: 100%                                830/830 [01:38<00:00, 8.86batch/s, current_loss=1.81, mean_loss=1.93]

epoch 7: 100%                                830/830 [01:36<00:00, 9.31batch/s, current_loss=2.38, mean_loss=1.82]

epoch 8: 100%                                830/830 [01:36<00:00, 9.62batch/s, current_loss=1.58, mean_loss=1.72]

```

```

1 # Reload the model from the model file.
2 # Useful when you have already trained and saved the model
3 attention_model = Seq2seqAttention(vocab).to(device)
4 attention_model.load_state_dict(torch.load("attention_model.pt", map_location=device))

<All keys matched successfully>

```

Let's test the attention model on the some sample inputs.

```
1 def test_conversations_with_model(model, conversational_inputs = None, include_beam = False):
2     # Some predefined conversational inputs.
3     # You may append more inputs at the end of the list, if you want to.
4     basic_conversational_inputs = [
5         "hello.",
6         "please share you bank account number with me",
7         "i have never met someone more annoying that you",
8         "i like pizza. what do you like?",
9         "give me coffee, or i'll hate you",
10        "i'm so bored. give some suggestions",
11        "stop running or you'll fall hard",
12        "what is your favorite sport?",
13        "do you believe in a miracle?",
14        "which sport team do you like?"
15    ]
16    if not conversational_inputs:
17        conversational_inputs = basic_conversational_inputs
18    for input in conversational_inputs:
19        print(f"Input > {input}")
20        generation = predict_greedy(model, input)
21        print('Greedy Response:', generation)
22        if include_beam:
23            # Also print the beam search responses from models
24            generations = predict_beam(model, input)
25            print('Beam Responses:')
26            print_list(generations)
27        print()

1 baseline_chat_inputs = [inp for inp, gen in baseline_chat]
2 attention_chat = test_conversations_with_model(attention_model, baseline_chat_inputs)
```

```
Input > how are you?
Greedy Response: i m fine .
```

```
Input > what do you want for lunch?
Greedy Response: i want to go to cash .
```

```
Input > did you like our time together?
Greedy Response: no .
```

```
Input > who is your favorite actor?
Greedy Response: my employer .
```

```
Input > what do you want to eat for dinner?
Greedy Response: i can t take a walk .
```

```
Input > What's your favorite movie?
Greedy Response: yeah .
```

```
Input > thanks.
Greedy Response: i m glad you d like me .
```

```
Input > did you like our time together?
Greedy Response: no .
```

```
Input > why not?
Greedy Response: i don t know .
```

▼ Part 5: Automatic Evaluation (5 points)

Automatic evaluation of chatbots is an active research area. For this assignment we are going to use 3 very simple evaluation metrics.

1. Average Length of the Responses
2. Distinct1 = proportion of unique unigrams / total unigrams
3. Distinct2 = proportion of unique bigrams / total bigrams

Length in this case refers to the number of tokens in the models response. You will evaluate your baseline and attention models by running the cells below.

```
1 # Evaluate diversity of the models
2 def evaluate_diversity(model, mode="greedy"):
3     """Evaluates the model's greedy or beam responses on eval_conversations
4
5     Args:
6         model: A sequence-to-sequence model.
7         mode: "greedy" or "beam"
8
9     Returns: avg_length, distinct1, distinct2
10        avg_length: average length of the model responses
11        distinct1: proportion of unique unigrams / total unigrams
12        distinct2: proportion of unique bigrams / total bigrams
13    """
14    if mode == "beam":
15        predict_f = predict_beam
```

```

16     else:
17         predict_f = predict_greedy
18         generations = list()
19         for src, tgt in eval_conversations:
20             generation = predict_f(model, src)
21             if mode == "beam":
22                 generation = generation[0]
23             generations.append(generation)
24         # Calculate average length, distinct unigrams and bigrams from generations
25         avg_length, distinct1, distinct2 = 0, 0, 0
26
27         # YOUR CODE HERE
28         total_length = sum(len(g) for g in generations)
29         num_g = len(generations)
30         avg_length = total_length/num_g
31
32         all_tokens = [token for generation in generations for token in vocab.tokenized_sentence(generation)]
33         num_tokens = len(all_tokens)
34         distinct_tokens = set(all_tokens)
35         num_distinct_tokens = len(distinct_tokens)
36
37         all_bigrams = list(nltk.bigrams(all_tokens))
38         num_bigrams = len(all_bigrams)
39         distinct_bigrams = set(all_bigrams)
40         num_distinct_bigrams = len(distinct_bigrams)
41
42         distinct1 = num_distinct_tokens / num_tokens
43         distinct2 = num_distinct_bigrams / num_bigrams
44
45
46         return avg_length, distinct1, distinct2

1 print(f"Baseline Model evaluation:")
2 avg_length, distinct1, distinct2 = evaluate_diversity(baseline_model)
3 print(f"Greedy decoding:")
4 print(f"Avg Response Length = {avg_length}")
5 print(f"Distinct1 = {distinct1}")
6 print(f"Distinct2 = {distinct2}")
7 print(f"Attention Model evaluation:")
8 avg_length, distinct1, distinct2 = evaluate_diversity(attention_model)
9 print(f"Greedy decoding:")
10 print(f"Avg Response Length = {avg_length}")
11 print(f"Distinct1 = {distinct1}")
12 print(f"Distinct2 = {distinct2}")

Baseline Model evaluation:
Greedy decoding:
Avg Response Length = 11.56
Distinct1 = 0.0777027027027027
Distinct2 = 0.1404399323181049
Attention Model evaluation:
Greedy decoding:
Avg Response Length = 14.89
Distinct1 = 0.13574660633484162
Distinct2 = 0.26132930513595165

```

▼ Part 6: BERT Finetuning (5 points)

Introduced in the paper "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding"

(<https://arxiv.org/pdf/1810.04805.pdf>), the pretrained transformer model BERT is heavily used within NLP research and engineering. This section will walk you through the use of the popular Huggingface Transformers library so that you can utilize it for your final projects and any research you may pursue.

The HuggingFace documentation can be found here: <https://huggingface.co/transformers/>. You will need to refer to the documentation frequently through this section.

The Dataset preparation and Model Helpers subsections contain utility code to setup this portion of the project. **Your first task begins in the second cell in the Model Setup subsection** where you will download the pretrained model. After this, you will add a classification head to the model so that we can classify disaster tweets.

▼ Dataset Preparation

Kaggle is a popular machine learning website that runs competitions for machine learning datasets. We will be using the Kaggle dataset "Natural Language Processing with Disaster Tweets" for this assignment. This dataset contains tweets that were sent in response to an actual disaster or that merely contain language similar to that used to describe a disaster. The goal of this challenge, and of this section, is to train a model that can classify tweets as either disaster related or non disaster related. For the following section, we are using the data from <https://www.kaggle.com/c/nlp-getting-started/overview>. Feel free to create a Kaggle account and look at the competition in more depth; for this project, however, we will download the training data directly from the class repository.

```

1 import pandas as pd
2 import numpy as np
3 import sys

```

```

4 from functools import partial
5 import time

1 #load the data into a pandas dataframe
2 !wget https://raw.githubusercontent.com/cocoxu/CS4650_projects_spring2023/master/p3_bert_train.csv
3 full_df = pd.read_csv('p3_bert_train.csv', header=0)

--2023-04-06 16:22:59-- https://raw.githubusercontent.com/cocoxu/CS4650_projects_spring2023/master/p3_bert_train.csv
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.110.133, 185.199.109.133, 185.199.108.133, ...
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.110.133|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 987712 (965K) [text/plain]
Saving to: 'p3_bert_train.csv.1'

p3_bert_train.csv.1 100%[=====>] 964.56K --.-KB/s in 0.004s

2023-04-06 16:23:00 (221 MB/s) - 'p3_bert_train.csv.1' saved [987712/987712]

1 #divide data into train, validation, and test datasets
2 num_tweets = len(full_df)
3 idxs = list(range(num_tweets))
4 print('Total tweets in dataset: ', num_tweets)
5 test_idx = idxs[:int(0.1*num_tweets)]
6 val_idx = idxs[int(0.1*num_tweets):int(0.2*num_tweets)]
7 train_idx = idxs[int(0.2*num_tweets):]
8
9 train_df = full_df.iloc[train_idx].reset_index(drop=True)
10 val_df = full_df.iloc[val_idx].reset_index(drop=True)
11 test_df = full_df.iloc[test_idx].reset_index(drop=True)
12
13 train_data = train_df[['id', 'text', 'target']]
14 val_data = val_df[['id', 'text', 'target']]
15 test_data = test_df[['id', 'text', 'target']]

Total tweets in dataset: 7613

1 #Defining torch dataset class for disaster tweet dataset
2 class TweetDataset(Dataset):
3     def __init__(self, df):
4         self.df = df
5
6     def __len__(self):
7         return len(self.df)
8
9     def __getitem__(self, idx):
10        return self.df.iloc[idx]

1 #set up train, validation, and testing datasets
2 train_dataset = TweetDataset(train_data)
3 val_dataset = TweetDataset(val_data)
4 test_dataset = TweetDataset(test_data)

```

The following code creates a collate function for our tweet dataset that will tokenize the input tweets for use with our BERT models.

```

1 def transformer_collate_fn(batch, tokenizer):
2     bert_vocab = tokenizer.get_vocab()
3     bert_pad_token = bert_vocab['[PAD]']
4     bert_unk_token = bert_vocab['[UNK]']
5     bert_cls_token = bert_vocab['[CLS]']
6
7     sentences, labels, masks = [], [], []
8     for data in batch:
9         tokenizer_output = tokenizer([data['text']])
10        tokenized_sent = tokenizer_output['input_ids'][0]
11        mask = tokenizer_output['attention_mask'][0]
12        sentences.append(torch.tensor(tokenized_sent))
13        labels.append(torch.tensor(data['target']))
14        masks.append(torch.tensor(mask))
15    sentences = pad_sequence(sentences, batch_first=True, padding_value=bert_pad_token)
16    labels = torch.stack(labels, dim=0)
17    masks = pad_sequence(masks, batch_first=True, padding_value=0.0)
18    return sentences, labels, masks

```

▼ Model Helpers

This section defines helper functions for model training, evaluation, and inspection. You do not need to modify any code in the Model Helpers section.

```

1 #computes the amount of time that a training epoch took and displays it in human readable form
2 def epoch_time(start_time: int,
3               end_time: int):
4     elapsed_time = end_time - start_time
5     elapsed_mins = int(elapsed_time / 60)
6     elapsed_secs = int(elapsed_time - (elapsed_mins * 60))
7     return elapsed_mins, elapsed_secs

1 #count the number of trainable parameters in the model
2 def count_parameters(model: nn.Module):
3     return sum(p.numel() for p in model.parameters() if p.requires_grad)

1 #train a given model, using a pytorch dataloader, optimizer, and scheduler (if provided)
2 def train(model,
3          dataloader,
4          optimizer,
5          device,
6          clip: float,
7          scheduler = None):
8
9     model.train()
10
11     epoch_loss = 0
12
13     for batch in dataloader:
14         sentences, labels, masks = batch[0], batch[1], batch[2]
15
16         optimizer.zero_grad()
17
18         output = model(sentences.to(device), masks.to(device))
19         loss = F.cross_entropy(output, labels.to(device))
20         loss.backward()
21         torch.nn.utils.clip_grad_norm_(model.parameters(), clip)
22
23         optimizer.step()
24         if scheduler is not None:
25             scheduler.step()
26
27         epoch_loss += loss.item()
28     return epoch_loss / len(dataloader)

1 #calculate the loss from the model on the provided dataloader
2 def evaluate(model,
3            dataloader,
4            device):
5
6     model.eval()
7
8     epoch_loss = 0
9     with torch.no_grad():
10         for batch in dataloader:
11             sentences, labels, masks = batch[0], batch[1], batch[2]
12             output = model(sentences.to(device), masks.to(device))
13             loss = F.cross_entropy(output, labels.to(device))
14
15             epoch_loss += loss.item()
16     return epoch_loss / len(dataloader)

1 #calculate the prediction accuracy on the provided dataloader
2 def evaluate_acc(model,
3                dataloader,
4                device):
5
6     model.eval()
7
8     epoch_loss = 0
9     with torch.no_grad():
10         total_correct = 0
11         total = 0
12         for i, batch in enumerate(dataloader):
13
14             sentences, labels, masks = batch[0], batch[1], batch[2]
15             output = model(sentences.to(device), masks.to(device))
16             output = F.softmax(output, dim=1)
17             output_class = torch.argmax(output, dim=1)
18             total_correct += torch.sum(torch.where(output_class == labels.to(device), 1, 0))
19             total += sentences.size()[0]
20
21     return total_correct / total

```

▼ Model Setup

```

1 #first, install the hugging face transformer package in your colab
2 !pip install -q transformers
3 from transformers import get_linear_schedule_with_warmup
4 from tokenizers.processors import BertProcessing

```

Having prepared our datasets, we now need to load in a BERT model for use as an encoder. Fortunately, the Hugging Face Library makes this easy for us. Use the hugging face AutoClass functionality to set up a pretrained Distil BERT Model and its corresponding tokenizer (1 Point). You will need to import functionality from the Hugging Face library for this question. If you are curious about the differences between BERT and Distil Bert, please see this page within the Huggingface Documentation: https://huggingface.co/transformers/model_summary.html

```
1 # Do not change this line, as it sets the model the model that Hugging Face will load
2 # If you are interested in what other models are available, you can find the list of model names here:
3 # https://huggingface.co/transformers/pretrained_models.html
4 bert_model_name = 'distilbert-base-uncased'
5
6 ##YOUR CODE HERE##
7
8 from transformers import DistilBertTokenizer, DistilBertModel
9 bert_model = DistilBertModel.from_pretrained(bert_model_name)
10 tokenizer = DistilBertTokenizer.from_pretrained(bert_model_name)
```

Some weights of the model checkpoint at distilbert-base-uncased were not used when initializing DistilBertModel: ['vocab_transform.bias']
- This IS expected if you are initializing DistilBertModel from the checkpoint of a model trained on another task or with another architecture
- This IS NOT expected if you are initializing DistilBertModel from the checkpoint of a model that you expect to be exactly identical (

If you've loaded the architecture correctly, the displayed name of the model below should be "DistilBertModel"

```
1 #print the loaded model architecture
2 bert_model

DistilBertModel(
  (embeddings): Embeddings(
    (word_embeddings): Embedding(30522, 768, padding_idx=0)
    (position_embeddings): Embedding(512, 768)
    (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
    (dropout): Dropout(p=0.1, inplace=False)
  )
  (transformer): Transformer(
    (layer): ModuleList(
      (0-5): 6 x TransformerBlock(
        (attention): MultiHeadSelfAttention(
          (dropout): Dropout(p=0.1, inplace=False)
          (q_lin): Linear(in_features=768, out_features=768, bias=True)
          (k_lin): Linear(in_features=768, out_features=768, bias=True)
          (v_lin): Linear(in_features=768, out_features=768, bias=True)
          (out_lin): Linear(in_features=768, out_features=768, bias=True)
        )
        (sa_layer_norm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
        (ffn): FFN(
          (dropout): Dropout(p=0.1, inplace=False)
          (lin1): Linear(in_features=768, out_features=3072, bias=True)
          (lin2): Linear(in_features=3072, out_features=768, bias=True)
          (activation): GELUActivation()
        )
        (output_layer_norm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)
      )
    )
  )
)
```

After loading the pretrained Distil BERT Model, we need to add our own classification head that we can train for our task. Assuming that the BERT encoder is a pretrained DistilBert model, add a BERT sequence classification head to architecture below. The classification head should take the encoded classification token as an input and output raw, unnormalized classification scores for each input sentence in the batch. You will need to look at the Huggingface documentation for DistilBert to complete this question, and you may want to look at the DistilBertForSequenceClassification architecture for guidance on creating a bert sequence classification head. Both can be found here: https://huggingface.co/transformers/model_doc/distilbert.html . (2 Points)

Please note that we are not allowing you to directly use the DistilBertForSequenceClassification architecture, as we want you to implement the BERT sequence classification head yourself.

```
1 class TweetClassifier(nn.Module):
2     def __init__(self,
3                 bert_encoder: nn.Module,
4                 enc_hid_dim=768, #default embedding size
5                 outputs=2,
6                 dropout=0.1):
7         super().__init__()
8
9         self.bert_encoder = bert_encoder
10
11         self.enc_hid_dim = enc_hid_dim
12
13
14     ### YOUR CODE HERE ###
15     # Define a linear layer to map the output from BERT to the classification layer
16     self.fc = nn.Linear(enc_hid_dim, outputs)
17
18     # Define a dropout layer to prevent overfitting
19     self.dropout = nn.Dropout(dropout)
20
21
```



```

22     def forward(self,
23                 src,
24                 mask):
25         bert_output = self.bert_encoder(src, mask)
26
27         ### YOUR CODE HERE ###
28         last_hidden_state = bert_output.last_hidden_state
29
30         # Average pool across tokens to get a single vector representation
31         avg_pool = torch.mean(last_hidden_state, 1)
32
33         # Apply dropout to avoid overfitting
34         x = self.dropout(avg_pool)
35
36         # Pass through fully connected layer to get logits
37         x = self.fc(x)
38         return x
39
40

```

Finally, we want to initialize the weights of our classification head without overwriting the weights within the DistilBert encoder. The `init_weights` function below will overwrite all weights within the model. Fill in the `init_classification_head_weights` function so that it will only overwrite weights in the classification head (using the same initialization scheme as the `init_weights` function). It may be helpful to refer to the PyTorch documentation on `nn.module.named_parameters()` while working on this question (1 point)

It should be noted that the weight initialization scheme utilized here is automatically implemented by PyTorch Linear layers. The goal of this question is to show how to change aspects of your model's set up at the parameter level basis, not just to initialize the correct weights for this architecture. As such, stating that the PyTorch Linear layer already implements this initialization scheme is not sufficient to earn points for this question.

```

1 def init_weights(m: nn.Module, hidden_size=768):
2     k = 1/hidden_size
3     for name, param in m.named_parameters():
4         if 'weight' in name:
5             print(name)
6             nn.init.uniform_(param.data, a=-1*k**0.5, b=k**0.5)
7         else:
8             print(name)
9             nn.init.uniform_(param.data, 0)

1 def init_classification_head_weights(m: nn.Module, hidden_size=768):
2     ### YOUR CODE STARTS HERE ###
3     k = 1/hidden_size
4     for name, param in m.named_parameters():
5         if 'classifier.weight' in name:
6             nn.init.uniform_(param.data, a=-1*k**0.5, b=k**0.5)
7         elif 'classifier.bias' in name:
8             nn.init.zeros_(param.data)

```

▼ Model Training

Once you have written the `init_classification_head_weights` function, you are done coding for this question. Run the following cells to initialize your model, to set up training, validation, and test dataloaders, and to train/evaluate the model. If you have completed the previous steps correctly, your model should achieve a test accuracy of 80% or greater without any hyperparameter tuning. Please note that if you need to train your model more than once, you will need to reload the BERT model to ensure that you are starting with fresh weights. Make sure that your submitted colab notebook file for includes the printed test accuracy to receive full credit for this question. (1 Point)

```

1 #define hyperparameters
2 BATCH_SIZE = 10
3 LR = 1e-5
4 WEIGHT_DECAY = 0
5 N_EPOCHS = 3
6 CLIP = 1.0
7
8 #define models, move to device, and initialize weights
9 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
10
11 model = TweetClassifier(bert_model).to(device)
12 model.apply(init_classification_head_weights)
13 model.to(device)
14 print('Model Initialized')

Model Initialized

1 #create pytorch dataloaders from train_dataset, val_dataset, and test_dataset
2 train_dataloader = DataLoader(train_dataset, batch_size=BATCH_SIZE, collate_fn=partial(transformer_collate_fn, tokenizer=tokenizer), shuffle=True)
3 val_dataloader = DataLoader(val_dataset, batch_size=BATCH_SIZE, collate_fn=partial(transformer_collate_fn, tokenizer=tokenizer))
4 test_dataloader = DataLoader(test_dataset, batch_size=BATCH_SIZE, collate_fn=partial(transformer_collate_fn, tokenizer=tokenizer))

```

```

1 optimizer = optim.Adam(model.parameters(), lr=LR)
2
3 scheduler = get_linear_schedule_with_warmup(optimizer, num_warmup_steps=10, num_training_steps=N_EPOCHS*len(train_dataloader))
4
5 print(f'The model has {count_parameters(model):,} trainable parameters')
6
7 train_loss = evaluate(model, train_dataloader, device)
8 train_acc = evaluate_acc(model, train_dataloader, device)
9
10 valid_loss = evaluate(model, val_dataloader, device)
11 valid_acc = evaluate_acc(model, val_dataloader, device)
12
13 print(f'Initial Train Loss: {train_loss:.3f}')
14 print(f'Initial Train Acc: {train_acc:.3f}')
15 print(f'Initial Valid Loss: {valid_loss:.3f}')
16 print(f'Initial Valid Acc: {valid_acc:.3f}')
17
18 for epoch in range(N_EPOCHS):
19     start_time = time.time()
20     train_loss = train(model, train_dataloader, optimizer, device, CLIP, scheduler)
21     end_time = time.time()
22     train_acc = evaluate_acc(model, train_dataloader, device)
23     valid_loss = evaluate(model, val_dataloader, device)
24     valid_acc = evaluate_acc(model, val_dataloader, device)
25     epoch_mins, epoch_secs = epoch_time(start_time, end_time)
26
27     print(f'Epoch: {epoch+1:02} | Time: {epoch_mins}m {epoch_secs}s')
28     print(f'\tTrain Loss: {train_loss:.3f}')
29     print(f'\tTrain Acc: {train_acc:.3f}')
30     print(f'\tValid Loss: {valid_loss:.3f}')
31     print(f'\tValid Acc: {valid_acc:.3f}')

The model has 66,364,418 trainable parameters
Initial Train Loss: 0.691
Initial Train Acc: 0.497
Initial Valid Loss: 0.696
Initial Valid Acc: 0.470
Epoch: 01 | Time: 0m 50s
    Train Loss: 0.446
    Train Acc: 0.866
    Valid Loss: 0.359
    Valid Acc: 0.854
Epoch: 02 | Time: 0m 47s
    Train Loss: 0.344
    Train Acc: 0.895
    Valid Loss: 0.364
    Valid Acc: 0.846
Epoch: 03 | Time: 0m 49s
    Train Loss: 0.295
    Train Acc: 0.906
    Valid Loss: 0.382
    Valid Acc: 0.852

1 #run this cell and save its outputs. A 75% test accuracy is needed for full credit.
2 test_loss = evaluate(model, test_dataloader, device)
3 test_acc = evaluate_acc(model, test_dataloader, device)
4 print(f'Test Loss: {test_loss:.3f}')
5 print(f'Test Acc: {test_acc:.3f}')

Test Loss: 0.511
Test Acc: 0.808

```

▼ Part 7: Beam Search (Extra Credit, 10 points)

Similar to greedy search, beam search generates one token at a time. However, rather than keeping only the single best hypothesis, we instead keep the top k candidates at each time step. This is accomplished by computing the set of next-token extensions for each item on the beam and finding the top k across all candidates according to total log-probability.

Candidates that are finished should be extracted in a final list of `generations` and removed from the beam. This strategy is useful for doing re-ranking the beam candidates using alternate scorers (example, Maximum Mutual Information Objective from [Li et al. 2015](#)). For this assignment, you will re-rank the beam generations as follows,

$$final_score_i = \frac{score_i}{|generation_i|^\alpha}, \text{ where } \alpha \in [0.5, 2].$$

Terminate the search process once you have k items in the `generations` list.

HINT: Given the simplicity of the dataset we're working with, it's likely that the responses from your model will be similar to each other but they should not be the exact same.

```

1 def predict_beam(model, sentence, k=5, max_length=100):
2     """Make predictions for the given inputs using beam search.
3
4     Args:
5         model: A sequence-to-sequence model.
6         sentence: An input sentence, represented as string.
7         k: The size of the beam.
8         max_length: The maximum length at which to truncate outputs in order to
9             avoid non-terminating inference.
10

```

```

11 Returns:
12     A list of k beam predictions. Each element in the list should be a string
13     corresponding to one of the top k predictions for the corresponding input,
14     sorted in descending order by its final score.
15     """
16
17 # Implementation tip: once an eos_token has been generated for any beam,
18 # remove its subsequent predictions from that beam by adding a small negative
19 # number like -1e9 to the appropriate logits. This will ensure that the
20 # candidates are removed from the beam, as its probability will be very close
21 # to 0. Using this method, you will be able to reuse the beam of an already
22 # finished candidate
23
24 # Implementation tip: while you are encouraged to keep your tensor dimensions
25 # constant for simplicity (aside from the sequence length), some special care
26 # will need to be taken on the first iteration to ensure that your beam
27 # doesn't fill up with k identical copies of the same candidate.
28
29 # You are welcome to tweak alpha
30 alpha = 0.7
31 model.eval()
32
33 # YOUR CODE HERE
34 # Tokenize input sentence
35 input = torch.tensor(vocab.get_ids_from_sentence(normalize_sentence(sentence))).unsqueeze(1).to(device)
36 inputs, encode_mask, encode_hidden = model.encode(input)
37 hidden = encode_hidden
38
39 # Number of input tokens
40 num_tokens = inputs.size(1)
41
42 # Expand inputs to size k
43 inputs = inputs.expand(k, num_tokens)
44
45 # Initialize scores and output sequences
46 seq_scores = torch.zeros(k, 1)
47 seq_outputs = inputs.clone()
48
49 # Initialize the hidden state and cell state of the decoder with zeros
50 # hidden = torch.zeros(model.decoder.num_layers, k, model.decoder.hidden_size)
51 cell = torch.zeros(model.decoder.num_layers, k, model.decoder.hidden_size)
52
53 # The first input to the decoder is the <sos> token
54 decoder_input = torch.tensor([[bos_id]])
55
56 # List to store completed sequences and their scores
57 completed_seqs = []
58 completed_seq_scores = []
59
60 for i in range(max_length):
61
62     # Pass the inputs and the decoder state through the decoder to get
63     # the logits and the new decoder state
64     logits, hidden, cell = model.decode(decoder_input, hidden)
65
66     # Apply softmax to the logits to get the probabilities over the vocabulary
67     probs = F.softmax(logits, dim=-1)
68
69     # Multiply the probabilities by the scores of the corresponding sequence
70     # and take the sum over the sequence dimension to get the new scores
71     scores = seq_scores.expand_as(probs) * probs
72     scores = scores.reshape(-1, k)
73
74     # Keep the top k scores and their corresponding indices
75     top_scores, top_indices = torch.topk(scores, k, dim=1)
76
77     # Convert the flattened indices to the indices within the sequence and
78     # the beam
79     beam_indices = top_indices // len(tokenizer)
80     token_indices = top_indices % len(tokenizer)
81
82     # Append the new tokens and scores to the output sequences and scores
83     new_seq_outputs = torch.cat([seq_outputs[beam_indices, :], token_indices.unsqueeze(-1)], dim=-1)
84     new_seq_scores = top_scores.view(-1, 1)
85
86     # Check if any of the sequences have reached the end token
87     eos_mask = token_indices == eos_id
88     if eos_mask.any():
89
90         # Remove the completed sequences from the current sequences
91         new_seq_outputs = new_seq_outputs[~eos_mask, :]
92         new_seq_scores = new_seq_scores[~eos_mask, :]
93
94         # Get the completed sequences and their scores
95         completed_seqs.extend(new_seq_outputs[eos_mask, :])
96         completed_seq_scores.extend(new_seq_scores[eos_mask, :])
97
98     # If all sequences have been completed, break out of the loop
99     if len(completed_seqs) == k:
100         break
101

```

```

102         # Otherwise, reduce k to the number of incomplete sequences
103         k -= eos_mask.sum().item()

```

Now let's test both baseline and attention models on some predefined inputs and compare their greedy and beam responses side by side.

```

1 test_conversations_with_model(baseline_model, include_beam=False)

```

```

Input > hello.
Greedy Response: hello .

Input > please share you bank account number with me
Greedy Response: i m not sure .

Input > i have never met someone more annoying that you
Greedy Response: i know .

Input > i like pizza. what do you like?
Greedy Response: i don t know .

Input > give me coffee, or i'll hate you
Greedy Response: you re not going to get my car ?

Input > i'm so bored. give some suggestions
Greedy Response: i m sorry .

Input > stop running or you'll fall hard
Greedy Response: i don t know .

Input > what is your favorite sport?
Greedy Response: i don t know .

Input > do you believe in a miracle?
Greedy Response: no .

Input > which sport team do you like?
Greedy Response: i don t know .

```

```

1 test_conversations_with_model(baseline_model, include_beam=True)

```

```

❏ Input > hello.
Greedy Response: hello .

-----
RuntimeError                                Traceback (most recent call last)
<ipython-input-206-4c9801536792> in <cell line: 1>()
----> 1 test_conversations_with_model(baseline_model, include_beam=True)

----- 1 frames -----
<ipython-input-205-609fb518387a> in predict_beam(model, sentence, k, max_length)
    41
    42     # Expand inputs to size k
----> 43     inputs = inputs.expand(k, num_tokens)
    44
    45     # Initialize scores and output sequences

RuntimeError: expand(torch.cuda.FloatTensor[[4, 1, 300]], size=[5, 1]): the number of sizes provided (2) must be greater or equal to
the number of dimensions in the tensor (3)

SEARCH STACK OVERFLOW

```

+ 代码 + 文本

```

1 test_conversations_with_model(attention_model, include_beam=False)

```

```

1 test_conversations_with_model(attention_model, include_beam=True)

```

Let's also check how our models do using our automatic evaluation metrics.

```

1 print(f"Baseline Model evaluation:")
2 avg_length, distinct1, distinct2 = evaluate_diversity(baseline_model)
3 print(f"Greedy decoding:")
4 print(f"Avg Response Length = {avg_length}")
5 print(f"Distinct1 = {distinct1}")
6 print(f"Distinct2 = {distinct2}")
7 avg_length, distinct1, distinct2 = evaluate_diversity(baseline_model, mode='beam')
8 print(f"Beam search decoding:")
9 print(f"Avg Response Length = {avg_length}")
10 print(f"Distinct1 = {distinct1}")
11 print(f"Distinct2 = {distinct2}")
12 print(f"Attention Model evaluation:")
13 avg_length, distinct1, distinct2 = evaluate_diversity(attention_model,)
14 print(f"Greedy decoding:")
15 print(f"Avg Response Length = {avg_length}")
16 print(f"Distinct1 = {distinct1}")
17 print(f"Distinct2 = {distinct2}")
18 avg_length, distinct1, distinct2 = evaluate_diversity(attention_model, mode='beam')
19 print(f"Beam decoding:")
20 print(f"Avg Response Length = {avg_length}")
21 print(f"Distinct1 = {distinct1}")
22 print(f"Distinct2 = {distinct2}")

```

What to turn in?

This is the end. Congratulations!

Now, follow the steps below to submit your homework in [Gradescope](#):

1. Rename this ipynb file to 'CS4650_p2_GTusername.ipynb'. We recommend ensuring you have removed any extraneous cells & print statements, clearing all outputs, and using the Runtime --> Run all tool to make sure all output is update to date. Additionally, leaving comments in your code to help us understand your operations will assist the teaching staff in grading. It is not a requirement, but is recommended.
2. Click on the menu 'File' --> 'Download' --> 'Download .py'.
3. Click on the menu 'File' --> 'Download' --> 'Download .ipynb'.
4. Download the notebook as a .pdf document. Make sure the output from your training loops are captured so we can see how the loss and accuracy changes while training.
5. Upload all 3 files to GradeScope.