

By:

Tom Kessous 206018749

Dan Ben Ami 316333079

## Assignment 4: Language Processing with RNN-Based Autoencoders

**Deadline:** Sunday, June 15th, by 9pm.

**Submission:** Submit a PDF export of the completed notebook as well as the ipynb file.

In this assignment, we will practice the application of deep learning to natural language processing. We will be working with a subset of Reuters news headlines that are collected over 15 months, covering all of 2019, plus a few months in 2018 and in a few months of this year.

In particular, we will be building an **autoencoder** of news headlines. The idea is similar to the kind of image autoencoder we built in lecture: we will have an **encoder** that maps a news headline to a vector embedding, and then a **decoder** that reconstructs the news headline. Both our encoder and decoder networks will be Recurrent Neural Networks, so that you have a chance to practice building

- a neural network that takes a sequence as an input
- a neural network that generates a sequence as an output

This assignment is organized as follows:

- Question 1. Exploring the data
- Question 2. Building the autoencoder
- Question 3. Training the autoencoder using *data augmentation*
- Question 4. Analyzing the embeddings (interpolating between headlines)

Furthermore, we'll be introducing the idea of **data augmentation** for improving of the robustness of the autoencoder, as proposed by Shen et al [1] in ICML 2020.

[1] Shen, Tianxiao, Jonas Mueller, Regina Barzilay, and Tommi Jaakkola. "Educating text autoencoders: Latent representation guidance via denoising." In International Conference on Machine Learning, pp. 8719-8729. PMLR, 2020.

In [ ]:

```
import torch
import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim

import matplotlib.pyplot as plt
import numpy as np
import random
```

### Question 1. Data (20 %)

Download the files `reuters_train.txt` and `reuters_valid.txt`, and upload them to Google Drive.

Then, mount Google Drive from your Google Colab notebook:

In [ ]:

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

train_path = '/content/gdrive/My Drive/Colab Notebooks/reuters_train.txt' # Update me
```

```
valid_path = '/content/gdrive/My Drive/Colab Notebooks/reuters_valid.txt' # Update me
```

Drive already mounted at /content/gdrive; to attempt to forcibly remount, call `drive.mount('/content/gdrive', force_remount=True)`.

As we did in some of our examples (e.g., training transformers on IMDB reviews) will be using PyTorch's `torchtext` utilities to help us load, process, and batch the data. We'll be using a `TabularDataset` to load our data, which works well on structured CSV data with fixed columns (e.g. a column for the sequence, a column for the label). Our tabular dataset is even simpler: we have no labels, just some text. So, we are treating our data as a table with one field representing our sequence.

In [ ]:

```
import torchtext.legacy.data as data

# Tokenization function to separate a headline into words
def tokenize_headline(headline):
    """Returns the sequence of words in the string headline. We also
    prepend the "<bos>" or beginning-of-string token, and append the
    "<eos>" or end-of-string token to the headline.
    """
    return ("<bos> " + headline + " <eos>").split()

# Data field (column) representing our *text*.
text_field = data.Field(
    sequential=True,           # this field consists of a sequence
    tokenize=tokenize_headline, # how to split sequences into words
    include_lengths=True,      # to track the length of sequences, for batching
    batch_first=True,         # similar to batch_first=True used in nn.RNN demonstrate
    d in lecture
    use_vocab=True)           # to turn each character into an integer index
train_data = data.TabularDataset(
    path=train_path,          # data file path
    format="tsv",              # fields are separated by a tab
    fields=[('title', text_field)]) # list of fields (we have only one)
```

## Part (a) -- 5%

Draw histograms of the number of words per headline in our training set. Excluding the `<bos>` and `<eos>` tags in your computation. Explain why we would be interested in such histograms.

In [ ]:

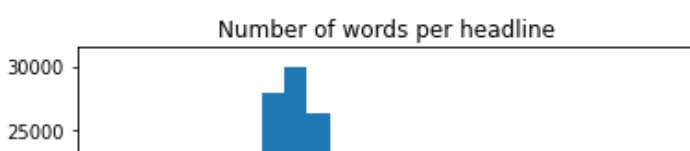
```
# Include your histogram and your written explanations

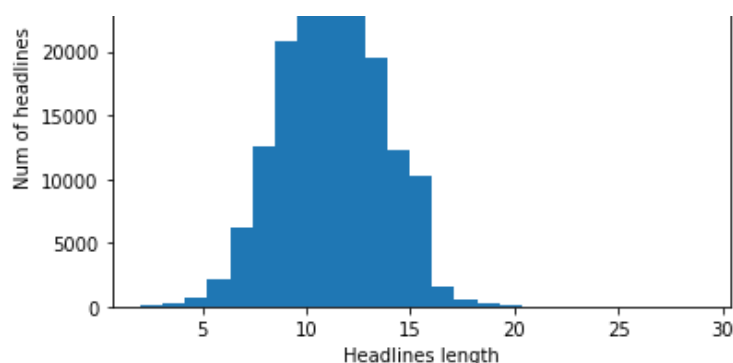
headlines_len = [len(headline.title)-2 for headline in train_data]

plt.hist(headlines_len, bins=25)
plt.title("Number of words per headline")
plt.xlabel("Headlines length")
plt.ylabel("Num of headlines")

# Here are some sample code that uses the train_data object:
print(train_data[5].title)
for example in train_data:
    print(example.title)
    break
```

```
[ '<bos>', 'u.s.', 'navy', 'pursuing', 'block', 'buy', 'of', 'two', 'aircraft', 'carriers'
, '-', 'senator', '<eos>' ]
[ '<bos>', 'dems', 'move', 'to', 'end', 'shutdown', ',', 'without', 'wall', 'money', '<eos>' ]
```





**Write your explanation here:**

we are interested in such histograms because our auto-encoder is a seq-to-seq model, and it contains RNN, so the histogram of the input's length will help us to know how much we need to take in counts the past (the previous words). Moreover, we can see that the variance of the headlines length's distribution is not high, therefore the padding won't be so dominant. This information helps us to batch the data into similar batch sizes, without padding too many words. If our model will contain much dominant sentence length, it will affect the model training and may bias the model to generate sentences with this length.

## Part (b) -- 5%

How many distinct words appear in the training data? Exclude the `<bos>` and `<eos>` tags in your computation.

In [ ]:

```
# Report your values here. Make sure that you report the actual values,
# and not just the code used to get those values

# You might find the python class Counter from the collections package useful

from collections import Counter

cnt = Counter()
for headline in train_data:
    cnt.update(headline.title)

cnt.pop('<bos>')
cnt.pop('<eos>')

print("There are", len(cnt), "distinct words in the training data.")
```

There are 51298 distinct words in the training data.

## Part (c) -- 5%

The distribution of *words* will have a long tail, meaning that there are some words that will appear very often, and many words that will appear infrequently. How many words appear exactly once in the training set? Exactly twice? Print these numbers below

In [ ]:

```
# Report your values here. Make sure that you report the actual values,
# and not just the code used to get those values

print("Number of words that appear exactly once in the training set:", Counter(cnt.values()) [1])
print("Number of words that appear exactly twice in the training set:", Counter(cnt.values()) [2])
```

Number of words that appear exactly once in the training set: 19854  
 Number of words that appear exactly twice in the training set: 7193

## Part (d) -- 5%

We will replace the infrequent words with an `<unk>` tag, instead of learning embeddings for these rare words. `torchtext` also provides us with the `<pad>` tag used for padding short sequences for batching. We will thus only model the top 9995 words in the training set, excluding the tags `<bos>`, `<eos>`, `<unk>`, and `<pad>`.

What percentage of total word count(whole dataset) will be supported? Alternatively, what percentage of total word count(whole dataset) in the training set will be set to the `<unk>` tag?

In [ ]:

```
# Report your values here. Make sure that you report the actual values,
# and not just the code used to get those values

support = sum([freq for _, freq in cnt.most_common(9995)])

print("percentage of total word count(whole dataset) will be supported: "+str(100 * support/sum(headlines_len))[:5]+"%")
print("percentage of total word count(whole dataset) in the training set will be set to the <unk> tag: "+str(100*(1 - support/sum(headlines_len)))[:5]+"%")
```

percentage of total word count(whole dataset) will be supported: 93.97%  
percentage of total word count(whole dataset) in the training set will be set to the `<unk>` tag: 6.021%

The `torchtext` package will help us keep track of our list of unique words, known as a **vocabulary**. A vocabulary also assigns a unique integer index to each word.

In [ ]:

```
# Build the vocabulary based on the training data. The vocabulary
# can have at most 9997 words (9995 words + the <bos> and <eos> token)
text_field.build_vocab(train_data, max_size=9997)

# This vocabulary object will be helpful for us
vocab = text_field.vocab
print(vocab.stoi["hello"]) # for instances, we can convert from string to (unique) index
print(vocab.itos[10])     # ... and from word index to string

# The size of our vocabulary
vocab_size = len(text_field.vocab.stoi)

# Here are the two tokens that torchtext adds for us:
print(vocab.itos[0]) # <unk> represents an unknown word not in our vocabulary
print(vocab.itos[1]) # <pad> will be used to pad short sequences for batching
```

0  
on  
<unk>  
<pad>

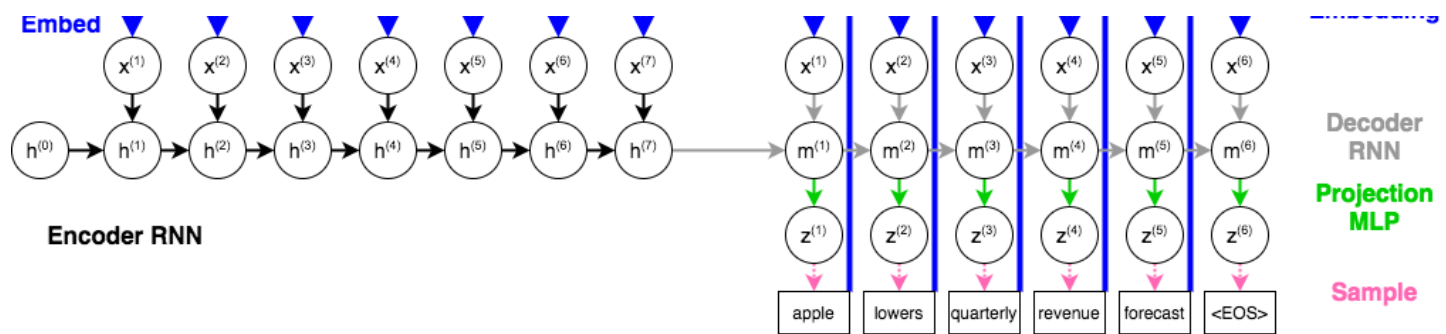
## Question 2. Text Autoencoder (40%)

Building a text autoencoder is a little more complicated than an image autoencoder like we did in class. So we will need to thoroughly understand the model that we want to build before actually building it. Note that the best and fastest way to complete this assignment is to spend time upfront understanding the architecture. The explanations are quite dense, but it is important to understand the operation of this model. The rationale here is similar in nature to the `seq2seq` RNN model we discussed in class, only we are dealing with unsupervised learning here rather than machine translation.

## Architecture description

Here is a diagram showing our desired architecture:





There are two main components to the model: the **encoder** and the **decoder**. As always with neural networks, we'll first describe how to make **predictions** with of these components. Let's get started:

The **encoder** will take a sequence of words (a headline) as *input*, and produce an embedding (a vector) that represents the entire headline. In the diagram above, the vector  $h^{(7)}$  is the vector embedding containing information about the entire headline. This portion is very similar to the sentiment analysis RNN that we discussed in lecture (but without the fully-connected layer that makes a prediction).

The **decoder** will take an embedding (in the diagram, the vector  $h^{(7)}$ ) as input, and uses a separate RNN to **generate a sequence of words**. To generate a sequence of words, the decoder needs to do the following:

1. Determine the previous word that was generated. This previous word will act as  $x^{(t)}$  to our RNN, and will be used to update the hidden state  $m^{(t)}$ . Since each of our sequences begin with the `<bos>` token, we'll set  $x^{(1)}$  to be the `<bos>` token.
2. Compute the updates to the hidden state  $m^{(t)}$  based on the previous hidden state  $m^{(t-1)}$  and  $x^{(t)}$ . Intuitively, this hidden state vector  $m^{(t)}$  is a representation of *all the words we still need to generate*.
3. We'll use a fully-connected layer to take a hidden state  $m^{(t)}$ , and determine *what the next word should be*. This fully-connected layer solves a *classification problem*, since we are trying to choose a word out of  $K = \text{vocab\_size}$  distinct words. As in a classification problem, the fully-connected neural network will compute a *probability distribution* over these `vocab_size` words. In the diagram, we are using  $z^{(t)}$  to represent the logits, or the pre-softmax activation values representing the probability distribution.
4. We will need to *sample* an actual word from this probability distribution  $z^{(t)}$ . We can do this in a number of ways, which we'll discuss in question 3. For now, you can imagine your favourite way of picking a word given a distribution over words.
5. This word we choose will become the next input  $x^{(t+1)}$  to our RNN, which is used to update our hidden state  $m^{(t+1)}$ , i.e., to determine what are the remaining words to be generated.

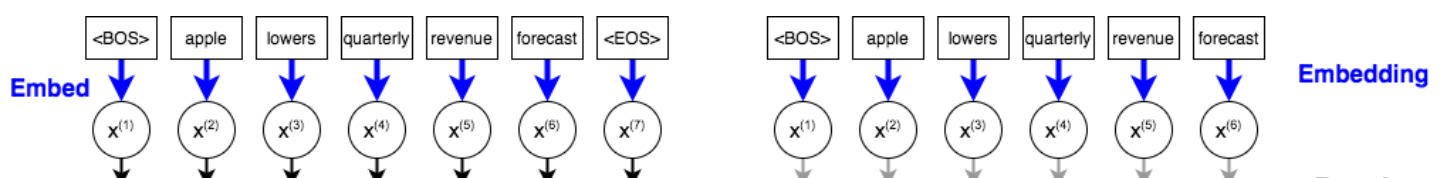
We can repeat this process until we see an `<eos>` token generated, or until the generated sequence becomes too long.

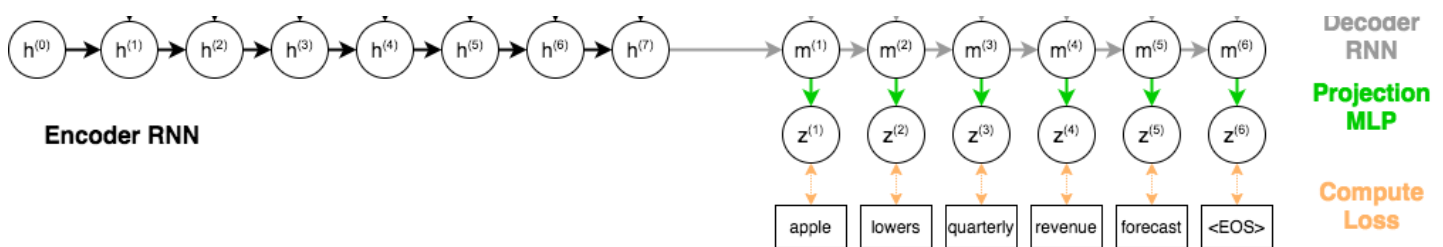
## Training the architecture

While our autoencoder produces a sequence, computing the loss by comparing the complete generated sequence to the ground truth (the encoder input) gives rise to multiple challenges. One is that the generated sequence might be longer or shorter than the actual sequence, meaning that there may be more/fewer  $z^{(t)}$ s than ground-truth words. Another more insidious issue is that the **gradients will become very high-variance and unstable**, because **early mistakes will easily throw the model off-track**. Early in training, our model is unlikely to produce the right answer in step  $t = 1$ , so the gradients we obtain based on the other time steps will not be very useful.

At this point, you might have some ideas about "hacks" we can use to make training work. Fortunately, there is one very well-established solution called **teacher forcing** which we can use for training: instead of *sampling* the next word based on  $z^{(t)}$ , we will forget sampling, and use the **ground truth**  $x^{(t)}$  as the input in the next step.

Here is a diagram showing how we can use **teacher forcing** to train our model:





We will use the RNN generator to compute the logits  $z^{(1)}, z^{(2)}, \dots, z^{(T)}$ . These distributions can be compared to the ground-truth words using the cross-entropy loss. The loss function for this model will be the sum of the losses across each  $t \in \{1, \dots, T\}$ .

We'll train the encoder and decoder model simultaneously. There are several components to our model that contain tunable weights:

- The word embedding that maps a word to a vector representation. In theory, we could use GloVe embeddings, as we did in class. In this assignment we will not do that, but learn the word embedding from data. The word embedding component is represented with blue arrows in the diagram.
- The encoder RNN (which will use GRUs) that computes the embedding over the entire headline. The encoder RNN is represented with black arrows in the diagram.
- The decoder RNN (which will also use GRUs) that computes hidden states, which are vectors representing what words are to be generated. The decoder RNN is represented with gray arrows in the diagram.
- The **projection MLP** (a fully-connected layer) that computes a distribution over the next word to generate, given a decoder RNN hidden state. The projection is represented with green arrows

## Part (a) -- 20%

Complete the code for the AutoEncoder class below by:

1. Filling in the missing numbers in the `__init__` method using the parameters `vocab_size`, `emb_size`, and `hidden_size`.
2. Complete the `forward` method, which uses teacher forcing and computes the logits  $z^{(t)}$  of the reconstruction of the sequence.

You should first try to understand the `encode` and `decode` methods, which are written for you. The `encode` method bears much similarity to the RNN we wrote in class for sentiment analysis. The `decode` method is a bit more challenging. You might want to scroll down to the `sample_sequence` function to see how this function will be called.

You can (but don't have to) use the `encode` and `decode` method in your `forward` method. In either case, be careful of the input that you feed into either `decode` or to `self.decoder_rnn`. Refer to the teacher-forcing diagram. **Notice that `batch_first` is set to `True`, understand how deal with it.**

In [ ]:

```
class AutoEncoder(nn.Module):
    def __init__(self, vocab_size, emb_size, hidden_size):
        """
        A text autoencoder. The parameters
        - vocab_size: number of unique words/tokens in the vocabulary
        - emb_size: size of the word embeddings  $x^{(t)}$ 
        - hidden_size: size of the hidden states in both the
                      encoder RNN ( $h^{(t)}$ ) and the
                      decoder RNN ( $m^{(t)}$ )

        """
        super().__init__()
        self.embed = nn.Embedding(num_embeddings=vocab_size, # TODO
                                   embedding_dim=emb_size) # TODO
        self.encoder_rnn = nn.GRU(input_size=emb_size, #TODO
                                   hidden_size=hidden_size, #TODO
                                   batch_first=True)
        self.decoder_rnn = nn.GRU(input_size=emb_size, #TODO
                                   hidden_size=hidden_size, #TODO
                                   batch_first=True)
```

```

self.proj = nn.Linear(in_features=hidden_size, # TODO
                       out_features=vocab_size) # TODO

def encode(self, inp):
    """
    Computes the encoder output given a sequence of words.
    """
    emb = self.embed(inp)
    out, last_hidden = self.encoder_rnn(emb)
    return last_hidden

def decode(self, inp, hidden=None):
    """
    Computes the decoder output given a sequence of words, and
    (optionally) an initial hidden state.
    """
    emb = self.embed(inp)
    out, last_hidden = self.decoder_rnn(emb, hidden)
    out_seq = self.proj(out)
    return out_seq, last_hidden

def forward(self, inp):
    """
    Compute both the encoder and decoder forward pass
    given an integer input sequence inp with shape [batch_size, seq_length],
    with inp[a,b] representing the (index in our vocabulary of) the b-th word
    of the a-th training example.

    This function should return the logits  $z^{(t)}$  in a tensor of shape
    [batch_size, seq_length - 1, vocab_size], computed using *teacher forcing*.

    The (seq_length - 1) part is not a typo. If you don't understand why
    we need to subtract 1, refer to the teacher-forcing diagram above.
    """
    last_hidden = self.encode(inp)
    out_seq, last_hidden = self.decode(inp[:, :-1], last_hidden)
    # -1 because <eos> isn't input
    return out_seq

```

## Part (b) -- 10%

To check that your model is set up correctly, we'll train our autoencoder neural network for at least 300 iterations to memorize this sequence:

In [ ]:

```

headline = train_data[42].title
input_seq = torch.Tensor([vocab.stoi[w] for w in headline]).long().unsqueeze(0)

```

We are looking for the way that you set up your loss function corresponding to the figure above. Be careful of off-by-one errors here.

Note that the Cross Entropy Loss expects a rank-2 tensor as its first argument (the output of the network), and a rank-1 tensor as its second argument (the true label). You will need to properly reshape your data to be able to compute the loss.

In [ ]:

```

model = AutoEncoder(vocab_size, 128, 128)
optimizer = optim.Adam(model.parameters(), lr=0.001)
criterion = nn.CrossEntropyLoss()

for it in range(300):
    output = model(input_seq)
    loss = criterion(output.squeeze(0).double(), input_seq[0, 1:])
    optimizer.zero_grad()
    loss.backward()
    optimizer.step()

```



```

if (it+1) % 50 == 0:
    print("[Iter %d] Loss %f" % (it+1, float(loss)))

```

```

[Iter 50] Loss 0.111175
[Iter 100] Loss 0.028504
[Iter 150] Loss 0.017759
[Iter 200] Loss 0.012393
[Iter 250] Loss 0.009225
[Iter 300] Loss 0.007190

```

In [ ]:

```

print("original headline is: ", headline)
print("restored headline is: ", [vocab.itos[idx] for idx in model(input_seq).squeeze(0).max(1, keepdim=True)[1].squeeze()])

```

```

original headline is:  ['<bos>', 'zambian', 'president', 'swears', 'in', 'new', 'army', 'chief', '<eos>']
restored headline is:  ['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief', '<eos>']

```

## Part (c) -- 4%

Once you are satisfied with your model, encode your input using the RNN encoder, and sample some sequences from the decoder. The sampling code is provided to you, and performs the computation from the first diagram (without teacher forcing).

Note that we are sampling from a multi-nomial distribution described by the logits  $z^{(t)}$ . For example, if our distribution is [80%, 20%] over a vocabulary of two words, then we will choose the first word with 80% probability and the second word with 20% probability.

Call `sample_sequence` at least 5 times, with the default temperature value. Make sure to include the generated sequences in your PDF report.

In [ ]:

```

def sample_sequence(model, hidden, max_len=20, temperature=1):
    """
    Return a sequence generated from the model's decoder
    - model: an instance of the AutoEncoder model
    - hidden: a hidden state (e.g. computed by the encoder)
    - max_len: the maximum length of the generated sequence
    - temperature: described in Part (d)
    """
    # We'll store our generated sequence here
    generated_sequence = []
    # Set input to the <BOS> token
    inp = torch.Tensor([text_field.vocab.stoi["<bos>"]]).long()
    for p in range(max_len):
        # compute the output and next hidden unit
        output, hidden = model.decode(inp.unsqueeze(0), hidden)
        # Sample from the network as a multinomial distribution
        output_dist = output.data.view(-1).div(temperature).exp()
        top_i = int(torch.multinomial(output_dist, 1)[0])
        # Add predicted word to string and use as next input
        word = text_field.vocab.itos[top_i]
        # Break early if we reach <eos>
        if word == "<eos>":
            break
        generated_sequence.append(word)
        inp = torch.Tensor([top_i]).long()
    return generated_sequence

# Your solutions go here
for i in range(5):
    headline = train_data[i+30].title
    print("original seq: ", headline)
    input_seq = torch.Tensor([vocab.stoi[w] for w in headline]).long().unsqueeze(0)

```



```
hidden = model.encode(input_seq)
print("restored seq: ", sample_sequence(model, hidden, max_len=20, temperature=1), "\n"
)
```

```
original seq: ['<bos>', 'refile-macau', 'casinos', 'rake', 'in', '$', '_num_', 'bln', 'g
aming', 'revenue', 'in', '_num_', '<eos>']
restored seq: ['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief']
```

```
original seq: ['<bos>', 'murray', 'makes', 'winning', 'return', 'in', 'brisbane', 'after
', 'hip', 'injury', '<eos>']
restored seq: ['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief']
```

```
original seq: ['<bos>', 'nasa', 'probe', 'believed', 'to', 'have', 'passed', 'distant',
'space', 'rock', 'on', 'landmark', 'mission', '<eos>']
restored seq: ['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief']
```

```
original seq: ['<bos>', 'soccer-qatar', 'capable', 'of', 'going', 'deep', 'in', 'asian',
'cup', ' ', ' ', 'says', 'al-haydos', '<eos>']
restored seq: ['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief']
```

```
original seq: ['<bos>', 'thai', 'king', 'to', 'be', 'crowned', 'in', 'coronation', 'cere
monies', 'may', '_num_-6', '<eos>']
restored seq: ['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief']
```

## Part (d) -- 6%

The multi-nomial distribution can be manipulated using the `temperature` setting. This setting can be used to make the distribution "flatter" (e.g. more likely to generate different words) or "peakier" (e.g. less likely to generate different words).

Call `sample_sequence` at least 5 times each for at least 3 different temperature settings (e.g. 1.5, 2, and 5). Explain why we generally don't want the temperature setting to be too large.

**Explain:** After training the model (even for memorizing one sentence in previous section) we have a words distribution that the model is sampling from in order to generate the new sentence. We generally don't want the temperature setting to be too large because then the distribution will be "flat" that no matter what sequence we will give to the model, the predicted sequence will be merely random. As we can see from the example below, in case of temperature high (lets say 5) the predicted sequence look like a random sequence without any meaning. If the temperature is low, the sampler will be more conservative - sticking more to those very likely tokens. On the other hand, if the temperature is high - less likely tokens will have a relatively higher likelihood - therefore, there is a bigger chance they will be selected.

In [ ]:

```
# Include the generated sequences and explanation in your PDF report.
temperature = [1,1.5,3,5]
for temp in temperature:
    print("Temperature:", temp)
    for i in range(5):
        headline = train_data[i+30].title
        print("original seq:",headline)
        input_seq = torch.Tensor([vocab.stoi[w] for w in headline]).long().unsqueeze(0)
        hidden = model.encode(input_seq)
        print("restored seq:", sample_sequence(model, hidden, max_len=20, temperature=temp),
"\n")
```

```
Temperature: 1
original seq: ['<bos>', 'refile-macau', 'casinos', 'rake', 'in', '$', '_num_', 'bln', 'ga
ming', 'revenue', 'in', '_num_', '<eos>']
restored seq: ['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief']
```

```
original seq: ['<bos>', 'murray', 'makes', 'winning', 'return', 'in', 'brisbane', 'after
', 'hip', 'injury', '<eos>']
restored seq: ['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief']
```

```
original seq: ['<bos>', 'nasa', 'probe', 'believed', 'to', 'have', 'passed', 'distant', '
space', 'rock', 'on', 'landmark', 'mission', '<eos>']
```

restored seq: ['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief']

original seq: ['<bos>', 'soccer-qatar', 'capable', 'of', 'going', 'deep', 'in', 'asian', 'cup', ',', 'says', 'al-haydos', '<eos>']

restored seq: ['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief']

original seq: ['<bos>', 'thai', 'king', 'to', 'be', 'crowned', 'in', 'coronation', 'ceremonies', 'may', '\_num\_-6', '<eos>']

restored seq: ['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief']

Temperature: 1.5

original seq: ['<bos>', 'refile-macau', 'casinos', 'rake', 'in', '\$', '\_num\_', 'bln', 'gaming', 'revenue', 'in', '\_num\_', '<eos>']

restored seq: ['zambian', 'president', 'transneft', 'in', 'into', 'founding', 'swears', 'pension', 'in', 'new', 'barcelona']

original seq: ['<bos>', 'murray', 'makes', 'winning', 'return', 'in', 'brisbane', 'after', 'hip', 'injury', '<eos>']

restored seq: ['zambian', 'president', 'new', 'army', 'chief']

original seq: ['<bos>', 'nasa', 'probe', 'believed', 'to', 'have', 'passed', 'distant', 'space', 'rock', 'on', 'landmark', 'mission', '<eos>']

restored seq: ['wrongdoing', 'predicts', 'in', 'new', 'army', 'mitsubishi', 'tour', 'libor', 'ministry', 'elon', 'returns', 'children', 'fortunes', 'strategic', 'in', 'new', 'army', 'chief']

original seq: ['<bos>', 'soccer-qatar', 'capable', 'of', 'going', 'deep', 'in', 'asian', 'cup', ',', 'says', 'al-haydos', '<eos>']

restored seq: ['zambian', 'president', 'swears', 'in', 'belarus', '\_num\_-australian', 'watch', 'take', 'sound', 'bertens', 'hunter', '\_num\_-occidental', 'america', 'arthur', 'in', 'new', 'raps', 'platforms', 'committees', 'retaliate']

original seq: ['<bos>', 'thai', 'king', 'to', 'be', 'crowned', 'in', 'coronation', 'ceremonies', 'may', '\_num\_-6', '<eos>']

restored seq: ['zambian', 'president', 'swears', 'kraft', 'w', 'new', 'army', 'already', 'in', 'new', 'army', 'chief']

Temperature: 3

original seq: ['<bos>', 'refile-macau', 'casinos', 'rake', 'in', '\$', '\_num\_', 'bln', 'gaming', 'revenue', 'in', '\_num\_', '<eos>']

restored seq: ['gaming', 'xinhua', 'structural', 'exclusion', 'throws', 'cross-border', 'indigenous', 'kerr', 'renews', 'canadians', 'strongly', 'down', 'japan-south', 'ratify', 'buffalo', 'tea', 'euros', 'pitt', 'bloodshed', 'texts']

original seq: ['<bos>', 'murray', 'makes', 'winning', 'return', 'in', 'brisbane', 'after', 'hip', 'injury', '<eos>']

restored seq: ['deportation', 'hotels', 'collapsed', 'false', 'current-quarter', 'contested', 'regains', 'tours', 'roster', 'spotify', 'prominent', 'national', 'trading', 'plotting', 'bmw', 'guests', 'constitutional', 'ivory', 'thrown', 'supports']

original seq: ['<bos>', 'nasa', 'probe', 'believed', 'to', 'have', 'passed', 'distant', 'space', 'rock', 'on', 'landmark', 'mission', '<eos>']

restored seq: ['woos', 'nokia', 'observatory', 'assad', 'accusations', 'rent', 're-elected', 'irregularities', 'dow', 'offshore', 'undergo', 'approving', 'esg', '\_num\_-mexico', 'ted', 'ad', 'is', '\_num\_-bristol-myers', 'garrett', 'ey']

original seq: ['<bos>', 'soccer-qatar', 'capable', 'of', 'going', 'deep', 'in', 'asian', 'cup', ',', 'says', 'al-haydos', '<eos>']

restored seq: ['holding', 'etfs', '\_num\_-cn', 'leviathan', 'jeremy', 'solid', 'lapses', '\_num\_-brexit', 'advances', 'author', 'operations', 'bird', 'fitness', 'values', 'clouds', 'lingering', 'travellers', 'light', 'ontario', 'walk-off']

original seq: ['<bos>', 'thai', 'king', 'to', 'be', 'crowned', 'in', 'coronation', 'ceremonies', 'may', '\_num\_-6', '<eos>']

restored seq: ['data', 'intended', '3rd', 'restoration', 'relation', 'damages', 'frustration', 'exodus', 'assad', 'denver', 'loved', 'buffalo', 'admissions', 'leclerc', 'flooding', 'facing', 'shoulder', 'rogers', 'sparking', '\_num\_-at']

Temperature: 5

original seq: ['<bos>', 'refile-macau', 'casinos', 'rake', 'in', '\$', '\_num\_', 'bln', 'gaming', 'revenue', 'in', '\_num\_', '<eos>']

restored seq: ['<bos>', 'announce', 'giuliani', 'ibrahimovic', 'num -france', 'beats', '']

```

assange', 'patients', 'interior', 'anything', 'prison', 'sacrifice', 'his', 'f-35', '_num_
_china', 'wildfires', 'processing', 'privatisation', 'france-klm', 'turned']

original seq: ['<bos>', 'murray', 'makes', 'winning', 'return', 'in', 'brisbane', 'after'
, 'hip', 'injury', '<eos>']
restored seq: ['conflicts', 'medicine', 'offset', 'wynn', 'failing', 'francis', 'pins', '
ruled', 'transcript', 'assault', 'judges', '_num_chinese', 'breach', 'district', 'fold',
'confronts', 'reducing', 'aware', 'say', 'lingering']

original seq: ['<bos>', 'nasa', 'probe', 'believed', 'to', 'have', 'passed', 'distant', '
space', 'rock', 'on', 'landmark', 'mission', '<eos>']
restored seq: ['ags', 'kickback', 'pregnancy', 'comedian', 'brace', 'suffer', 'smucker',
'planes', 'spiegel', 'once', 'grinds', 'truce', 'steelmaker', 'credible', 'stoke', 'outpe
rforms', 'wide', 'advisor', 'subpoena', 'td']

original seq: ['<bos>', 'soccer-qatar', 'capable', 'of', 'going', 'deep', 'in', 'asian',
'cup', ',', 'says', 'al-haydos', '<eos>']
restored seq: ['reinsurance', 'tariff', 'aluminium', 'party', 'reactions', 'never', 'impe
achment', 'emission', 'aides', 'ramadan', 'content', 'linked', 'r.', 'choppy', 'bln-', 'n
ov.', 'frictions', 'discipline', 'oregon', 'date']

original seq: ['<bos>', 'thai', 'king', 'to', 'be', 'crowned', 'in', 'coronation', 'cerem
onies', 'may', '_num_-6', '<eos>']
restored seq: ['employee', 'bosnia', 'fourteen', 'provincial', 'lng', 'stocks-dubai', 'au
tism', 'minister', 'speak', 'orbit', 'mester', 'nielsen', 'forgotten', 'affiliate', 'butt
on', 'tape', 'hurt', 'g7', 'premature', 'legislature']

```

### Question 3. Data augmentation (20%)

It turns out that getting good results from a text auto-encoder is very difficult, and that it is very easy for our model to **overfit**. We have discussed several methods that we can use to prevent overfitting, and we'll introduce one more today: **data augmentation**.

The idea behind data augmentation is to artificially increase the number of training examples by "adding noise" to the image. For example, during AlexNet training, the authors randomly cropped  $224 \times 224$  regions of a  $256 \times 256$  pixel image to increase the amount of training data. The authors also flipped the image left/right. Machine learning practitioners can also add Gaussian noise to the image.

When we use data augmentation to train an *autoencoder*, we typically to only add the noise to the input, and expect the reconstruction to be *noise free*. This makes the task of the autoencoder even more difficult. An autoencoder trained with noisy inputs is called a **denoising auto-encoder**. For simplicity, we will *not* build a denoising autoencoder today.

#### Part (a) -- 5%

We will add noise to our headlines using a few different techniques:

1. Shuffle the words in the headline, taking care that words don't end up too far from where they were initially
2. Drop (remove) some words
3. Replace some words with a blank word (a `<pad>` token)
4. Replace some words with a random word

The code for adding these types of noise is provided for you:

In [ ]:

```

def tokenize_and_randomize(headline,
                            drop_prob=0.1, # probability of dropping a word
                            blank_prob=0.1, # probability of "blanking" out a word
                            sub_prob=0.1,   # probability of substituting a word with a r
andom one
                            shuffle_dist=3): # maximum distance to shuffle a word
    """
    Add 'noise' to a headline by slightly shuffling the word order,
    dropping some words, blanking out some words (replacing with the <pad> token)

```

```

and substituting some words with random ones.
"""
headline = [vocab.stoi[w] for w in headline.split()]
n = len(headline)
# shuffle
headline = [headline[i] for i in get_shuffle_index(n, shuffle_dist)]

new_headline = [vocab.stoi['<bos>']]
for w in headline:
    if random.random() < drop_prob:
        # drop the word
        pass
    elif random.random() < blank_prob:
        # replace with blank word
        new_headline.append(vocab.stoi["<pad>"])
    elif random.random() < sub_prob:
        # substitute word with another word
        new_headline.append(random.randint(0, vocab_size - 1))
    else:
        # keep the original word
        new_headline.append(w)
new_headline.append(vocab.stoi['<eos>'])
return new_headline

def get_shuffle_index(n, max_shuffle_distance):
    """ This is a helper function used to shuffle a headline with n words,
    where each word is moved at most max_shuffle_distance. The function does
    the following:
    1. start with the *unshuffled* index of each word, which
       is just the values [0, 1, 2, ..., n]
    2. perturb these "index" values by a random floating-point value between
       [0, max_shuffle_distance]
    3. use the sorted position of these values as our new index
    """
    index = np.arange(n)
    perturbed_index = index + np.random.rand(n) * 3
    new_index = sorted(enumerate(perturbed_index), key=lambda x: x[1])
    return [index for (index, pert) in new_index]

```

Call the function `tokenize_and_randomize` 5 times on a headline of your choice. Make sure to include both your original headline, and the five new headlines in your report.

In [ ]:

```

# Report your values here. Make sure that you report the actual values,
# and not just the code used to get those values
headline = train_data[17].title
print("Original headline:", headline)

for i in range(5):
    prob = i*0.1
    print("After tokenize_and_randomize:", [vocab.itos[x] for x in tokenize_and_randomize(
        ' '.join(headline[1:-1]), prob, prob, prob)])

Original headline: ['<bos>', 'netflix', 'poaches', 'cfo', 'from', 'activision', 'blizzard', '-', 'source', '<eos>']
After tokenize_and_randomize: ['<bos>', 'poaches', 'netflix', 'cfo', 'activision', 'from', 'blizzard', 'source', '-', '<eos>']
After tokenize_and_randomize: ['<bos>', 'netflix', 'poaches', 'cfo', 'from', 'activision', 'blizzard', '-', 'source', '<eos>']
After tokenize_and_randomize: ['<bos>', 'u.n', '<pad>', 'from', '<pad>', '<pad>', '<eos>']
After tokenize_and_randomize: ['<bos>', 'netflix', 'divers', '<pad>', 'activision', 'celebrates', '<eos>']
After tokenize_and_randomize: ['<bos>', '<pad>', 'activision', '<eos>']

```

## Part (b) -- 8%

The training code that we use to train the model is mostly provided for you. The only part we left blank are the parts from Q2(b). Complete the code, and train a new AutoEncoder model for 1 epoch. You can train your model

parts from Q2(b). Complete the code, and train a new AutoEncoder model for 1 epoch. You can train your model for longer if you want, but training tend to take a long time, so we're only checking to see that your training loss is trending down.

If you are using Google Colab, you can use a GPU for this portion. Go to "Runtime" => "Change Runtime Type" and set "Hardware acceleration" to GPU. Your Colab session will restart. You can move your model to the GPU by typing `model.cuda()` , and move other tensors to GPU (e.g. `xs = xs.cuda()` ). To move a model back to CPU, type `model.cpu()` . To move a tensor back, use `xs = xs.cpu()` . For training, your model and inputs need to be on the *same device*.

In [ ]:

```
def train_autoencoder(model, batch_size=64, learning_rate=0.001, num_epochs=10):
    optimizer = optim.Adam(model.parameters(), lr=learning_rate)
    criterion = nn.CrossEntropyLoss()

    model.cuda()          #Sending the model to GPU
    val_loss_list = []
    val_n = 0
    iter = []
    for ep in range(num_epochs):
        # We will perform data augmentation by re-reading the input each time
        field = data.Field(sequential=True,
                           tokenize=tokenize_and_randomize, # <-- data augment
                           include_lengths=True,
                           batch_first=True,
                           use_vocab=False, # <-- the tokenization function re
                           pad_token=vocab.stoi['<pad>'])
        dataset = data.TabularDataset(train_path, "tsv", [('title', field)])

        # This BucketIterator will handle padding of sequences that are not of the same l
        train_iter = data.BucketIterator(dataset,
                                         batch_size=batch_size,
                                         sort_key=lambda x: len(x.title), # t
                                         repeat=False)

        dataset_valid = data.TabularDataset(valid_path, "tsv", [('title', field)])

        # This BucketIterator will handle padding of sequences that are not of the same l
        valid_iter = data.BucketIterator(dataset_valid,
                                         batch_size=batch_size,
                                         sort_key=lambda x: len(x.title), # t
                                         repeat=False)

        for it, ((xs, lengths), _) in enumerate(train_iter):
            xs = xs.cuda()
            output = model(xs)
            loss = criterion(output.reshape(-1, vocab_size).double(), xs[:,1:].reshape(-1
            ))

            optimizer.zero_grad()
            loss.backward()
            optimizer.step()

            if (it+1) % 100 == 0:
                print("[Iter %d] Loss %f" % (it+1, float(loss)))
                val_loss = 0
                for i, ((vs, lengths), _) in enumerate(valid_iter):
                    vs = vs.cuda()
                    zs = model(vs)
                    loss = criterion(zs.reshape(-1, vocab_size).double(), vs[:,1:].reshap
                    e(-1))

                    val_loss += float(loss)
                val_loss_list.append(val_loss)
```

```
        iter.append(it)
    return val_loss_list, iter
```

*# Include your training curve or output to show that your training loss is trending down*

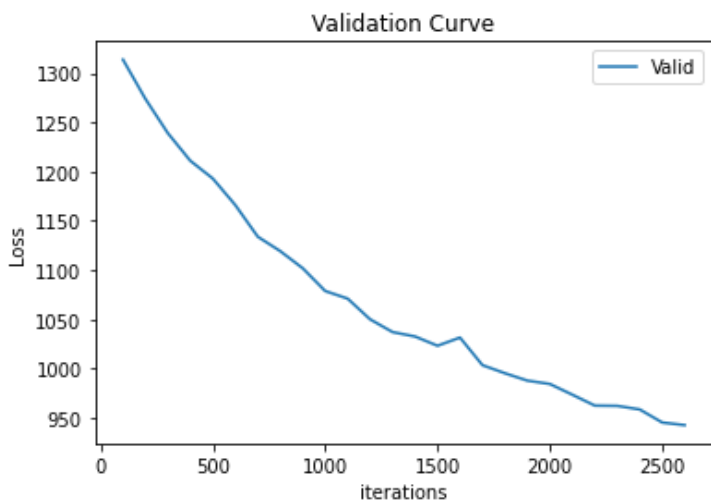
In [ ]:

```
model = AutoEncoder(10000, 128, 128)
val_loss_list, iter = train_autoencoder(model, batch_size=64, learning_rate=0.001, num_epochs=1)
```

```
[Iter 100] Loss 4.691894
[Iter 200] Loss 4.721688
[Iter 300] Loss 4.335512
[Iter 400] Loss 4.169982
[Iter 500] Loss 4.063227
[Iter 600] Loss 3.529230
[Iter 700] Loss 3.322038
[Iter 800] Loss 3.809353
[Iter 900] Loss 3.532978
[Iter 1000] Loss 3.167476
[Iter 1100] Loss 3.698731
[Iter 1200] Loss 3.877786
[Iter 1300] Loss 3.410537
[Iter 1400] Loss 3.973885
[Iter 1500] Loss 3.387378
[Iter 1600] Loss 3.217453
[Iter 1700] Loss 3.193566
[Iter 1800] Loss 3.421686
[Iter 1900] Loss 3.030582
[Iter 2000] Loss 3.316858
[Iter 2100] Loss 3.142178
[Iter 2200] Loss 3.681817
[Iter 2300] Loss 2.983117
[Iter 2400] Loss 3.527899
[Iter 2500] Loss 3.312171
[Iter 2600] Loss 3.269079
```

In [ ]:

```
plt.title("Validation Curve")
plt.plot(iter, val_loss_list, label="Valid")
plt.xlabel("iterations")
plt.ylabel("Loss")
plt.legend(loc='best')
plt.show()
```



## Part (c) -- 7%

This model requires many epochs (>50) to train, and is quite slow without using a GPU. You can train a model yourself, or you can load the model weights that we have trained, and available on the course website (AE\_RNN\_model.pk).

Assuming that your `AutoEncoder` is set up correctly, the following code should run without error.

In [ ]:

```
model = AutoEncoder(10000, 128, 128)
checkpoint_path = '/content/gdrive/My Drive/Colab Notebooks/model_parameters/AE/AE_RNN_model.pk' # Update me
model.load_state_dict(torch.load(checkpoint_path))
```

Out[ ]:

<All keys matched successfully>

Then, repeat your code from Q2(d), for `train_data[10].title` with temperature settings 0.7, 0.9, and 1.5. Explain why we generally don't want the temperature setting to be too small.

In [ ]:

```
# Include the generated sequences and explanation in your PDF report.
```

```
headline = train_data[10].title
input_seq = torch.Tensor([vocab.stoi[w] for w in headline]).unsqueeze(0).long()
print("original seq:", headline)
temperature = [0.7, 0.9, 1.5]
for temp in temperature:
    print("Temperature:", temp)
    for i in range(5):
        hidden = model.encode(input_seq)
        print("restored seq:", sample_sequence(model, hidden, max_len=20, temperature=temp),
              "\n")
```

```
original seq: ['<bos>', 'wall', 'street', 'rises', ',', 'limps', 'across', 'the', 'finish',
', 'line', 'of', 'a', 'turbulent', 'year', '<eos>']
```

```
Temperature: 0.7
```

```
restored seq: ['wall', 'street', 'rises', ',', 'limps', 'turn', 'vet', ',', 'closes', '<p',
ad>', 'after', 'euro', 'parliament']
```

```
restored seq: ['wall', 'street', 'rises', ',', 'limps', 'across', 'the', 'finish', 'line',
, 'of', 'a', 'highlights', 'note']
```

```
restored seq: ['wall', 'street', 'rises', ',', 'limps', 'die', 'win', '"s', 'employees',
'<pad>', 'one-time', 'after', 'wto']
```

```
restored seq: ['wall', 'street', 'rises', ',', 'limps', 'dead', 'cincinnati', 'to', 'pain',
t', ':', 'capital', 'election', 'fourth']
```

```
restored seq: ['wall', 'street', 'rises', ',', 'limps', 'die', 'win', 'at', 'of', 'scienc',
es', 'election', 'four']
```

```
Temperature: 0.9
```

```
restored seq: ['wall', 'street', 'rises', ',', 'limps', 'die', 'win', '"s', 'factory', 'p',
rotesters', '<pad>', 'olympics', 'war']
```

```
restored seq: ['wall', 'street', 'rises', ',', 'limps', 'dead', 'anthem', 'to', 'final',
'<pad>', 'israeli', 'year', 'angst']
```

```
restored seq: ['wall', 'street', 'rises', ',', 'limps', 'across', 'the', 'finish', 'line',
, 'of', 'a', 'turbulent', 'year']
```

```
restored seq: ['wall', 'street', 'rises', ',', 'limps', 'across', 'australia', 'request',
'peaceful', '<pad>', 'drone', 'three', 'tame']
```

```
restored seq: ['wall', 'street', 'rises', ',', 'limps', 'die', 'win', 'at', 'of', 'scienc',
es', 'election', 'four', 'daily']
```

```
Temperature: 1.5
```

```
restored seq: ['wall', 'street', 'del', ',', 'parting', 'del', 'into', 'others', 'league',
, 'bangladesh', '"s', 'investments', 'call']
```

```
restored seq: ['wall', 'street', 'rises', ',', 'jd.com', 'sector', 'lunch', 'win', 'detro',
it', '"s', 'second', 'dec', 'insults']
```



```
restored seq: ['wall', 'street', 'rises', ',', 'asahi', 'shows', 'collecting', 'at', 'for', 'stays', 'europe', 'city', 'nordstrom']
```

```
restored seq: ['wall', 'pro-democracy', 'results', 'up', 'risk', 'mario', 'butina', 'man', 'with', 'low', '<pad>', 'january', 'for']
```

```
restored seq: ['captures', 'european', 'stocks', 'curtailments', 'men', 'full-year', 'turn', 'after', 'cut', 'to', 'northeast', 'site', 'presidency']
```

**Write your explanation here:**

After training the model we have a words distribution that the model is sampling from in order to generate the new sentence. We can see that for low temperature values the distribution is "peakier", with low variance, and the model's sampler will be more conservative - sticking more to those very likely tokens.

## Question 4. Latent space manipulations (20%)

In parts 2-3, we've explored the decoder portion of the autoencoder. In this section, let's explore the **encoder**. In particular, the encoder RNN gives us embeddings of news headlines!

First, let's load the **validation data set**:

In [ ]:

```
valid_data = data.TabularDataset(
    path=valid_path,          # data file path
    format="tsv",             # fields are separated by a tab
    fields=[('title', text_field)]) # list of fields (we have only one)

print(len(valid_data))
```

19046

### Part (a) -- 4%

Compute the embeddings of every item in the validation set. Then, store the result in a single PyTorch tensor of shape `[19046, 128]`, since there are 19,046 headlines in the validation set.

In [ ]:

```
# Write your code here
# Show that your resulting PyTorch tensor has shape `[19046, 128]`
valid_iter = data.BucketIterator(valid_data,
                                batch_size=len(valid_data),
                                sort_key=lambda x: len(x.title), # to minimize padding
                                repeat=False)

for it, ((xs, lengths), _) in enumerate(valid_iter): #There is just 1 iteration.
    Emb_valid = model.encode(xs).squeeze(0)

print("Embedding valid shape:", Emb_valid.shape)
print("type of embedding valid:", type(Emb_valid))
```

```
Embedding valid shape: torch.Size([19046, 128])
type of embedding valid: <class 'torch.Tensor'>
```

### Part (b) -- 4%

Find the 5 closest headlines to the headline `valid_data[13]`. Use the cosine similarity to determine closeness. (Hint: You can use code from assignment 2)

In [ ]:

```
# Write your code here. Make sure to include the actual 5 closest headlines
```

*# Write your code here. Make sure to include the actual 5 closest headlines.*

```
Emb_valid_np = Emb_valid.detach().numpy()
norms = np.linalg.norm(Emb_valid_np, axis=1)
word_emb_norm = (Emb_valid_np.T / norms).T
similarities = np.matmul(word_emb_norm, word_emb_norm.T)
```

In [ ]:

```
five_smallest_idx = np.argpartition(similarities[13], -6)[-6:-1]
```

```
print("The given headline is:", valid_data[13].title)
print("\nThe 5 closest headlines to that are:")
for idx in five_smallest_idx:
    print(valid_data[idx].title)
```

The given headline is: ['<bos>', 'asia', 'takes', 'heart', 'from', 'new', 'year', 'gains', 'in', 'u.s.', 'stock', 'futures', '<eos>']

The 5 closest headlines to that are:

```
['<bos>', 'us', 'stocks-wall', 'street', 'slightly', 'lower', 'as', 'citi', 'results', 'weigh', 'on', 'bank', 'stocks', '<eos>']
['<bos>', 'venezuela', 'lawmaker', 'seeks', 'refuge', 'in', 'argentine', 'embassy', 'after', 'colleague', "'s", 'arrest', '<eos>']
['<bos>', 'fbi', 'joins', 'criminal', 'investigation', 'into', 'boeing', '_num_', 'max', 'certification', '-', 'report', '<eos>']
['<bos>', 'bahrain', 'closer', 'to', 'extradition', 'of', 'footballer', 'held', 'in', 'thailand', '<eos>']
['<bos>', 'congo', 'president', 'says', 'ebola', 'outbreak', 'should', 'be', 'over', 'this', 'year', '<eos>']
```

## Part (c) -- 4%

Find the 5 closest headlines to another headline of your choice.

In [ ]:

```
# Write your code here.
# Make sure to include the original headline and the 5 closest headlines.
five_smallest_idx = np.argpartition(similarities[1], -6)[-6:-1]
```

```
print("The given headline is:", valid_data[1].title)
print("\nThe 5 closest headlines to that are:")
for idx in five_smallest_idx:
    print(valid_data[idx].title)
```

The given headline is: ['<bos>', 'indonesia', 'landslides', 'kill', 'at', 'least', 'two', ',', 'leave', 'dozens', 'missing', '<eos>']

The 5 closest headlines to that are:

```
['<bos>', 'ralph', 'lauren', 'profit', 'gets', 'boost', 'from', 'chinese', 'demand', ';', 'shares', 'rise', '_num_', '%', '<eos>']
['<bos>', 'indonesia', 'landslides', 'kill', 'at', 'least', 'two', ',', 'leave', 'dozens', 'missing', '<eos>']
['<bos>', 'singapore', 'central', 'bank', 'to', 'reveal', 'fx', 'intervention', ',', 'improving', 'transparency', '<eos>']
['<bos>', 'wall', 'street', 'trims', 'gains', 'after', 'kudlow', "'s", 'trade', 'comments', '<eos>']
['<bos>', 'precious-palladium', 'hits', 'record', 'high', ';', 'gold', 'gains', 'on', 'dovish', 'fed', 'talk', '<eos>']
```

## Part (d) -- 8%

Choose two headlines from the validation set, and find their embeddings. We will interpolate between the two embeddings like we did in the example presented in class for training autoencoders on MNIST.

Find 3 points, equally spaced between the embeddings of your headlines. If we let  $e_0$  be the embedding of your first headline and  $e_4$  be the embedding of your second headline, your three points should be:

$$e_1 = 0.75e_0 + 0.25e_4$$

$$e_2 = 0.50e_0$$

$$+ 0.50e_4$$

$$e_3 = 0.25e_0$$

$$+ 0.75e_4$$

Decode each of  $e_1$ ,  $e_2$  and  $e_3$  five times, with a temperature setting that shows some variation in the generated sequences. Try to get a logical and cool sentence (this might be hard).

In [ ]:

```
# Write your code here. Include your generated sequences.
e = [0]*5
e[0] = Emb_valid[60]
e[4] = Emb_valid[61]
for i in range(1,4):
    e[i]=(1-0.25*i)*e[0]+0.25*i*e[4]

for j in range(5):
    print("\nround", j+1, ":")
    for i in range(0,5):
        print("e",i," = ", sample_sequence(model, e[i].unsqueeze(0).unsqueeze(0), max_le
n=20, temperature=1.5))

round 1 :
e 0 = ['update', 'we', 'firmer', 'america', 'new', 'chemical', 'markets-asian', 'sales'
, 'call', 'ends', 'and', 'kill']
e 1 = ['uk', 'says', 'wall', 'void', 'update', 'semi-final', 'seeks', 'on', 'hits', 'ta
lks', '_num_-week']
e 2 = ['u.s.-iran', 'china', 'cup', 'santos', 'delay', 'find', 'as', 'pct', 'decline',
'extension', 'nov']
e 3 = ['dollar', 'to', 'islands', 'party', 'nearly', '_num_-congo', 'top', 'economic',
'year', '_num_-week']
e 4 = ['dollar', 'falls', 'for', 'a', 'since', 'suzuki', 'after', 'slowing', 'but', 'ec
onomic', 'lagarde', '<pad>', 'farage', 'parliament', 'fake']

round 2 :
e 0 = ['update', '_num_-huawei', 'says', 'european', 'smart', 'phones', 'sales', 'up',
'in', 'rail', ':', 'twist']
e 1 = ['update', 'u.s.-iran', 'ups', '_num_-dutch', 'u.s.', 'disorder', '_num_-delta',
',', 'and', '-baker', 'optimism', 'lebanese']
e 2 = ['investors', "'s", 'paso', 'for', 'depend', 'new', 'revenue', 'pound', 'standoff
', 'rate', 'outperforms']
e 3 = ['to', 'breakingviews', 'opec', 'cup', 'of', 'ftse', 'leads', 'weak', 'four', 'sa
mples', 'ecb']
e 4 = ['dollar', 'falls', 'for', 'a', 'third', 'day', 'on', 'rate', 'pause', 'bets']

round 3 :
e 0 = ['update', '_num_-huawei', 'to', 'maryland', '_num_-pompeo', 'testing', 'looking'
, 'to', 'gain', 'hits', 'coast', 'mozambique']
e 1 = ['uk', 'survey', 'toward', 'wing', 'for', 'santos', 'fans', 'nearly', 'continue',
'retreats', ':']
e 2 = ['dollar', 'update', 'wider', 'to', 'repsol', 'big', 'business', 'flooding', 'fad
es', 'of', 'boon']
e 3 = ['dollar', 'to', 'supervisory', ',', 'third', '100m', 'of', 'slow', 'trump', 'eye
ing', 'recovery']
e 4 = ['dollar', 'falls', 'for', 'a', 'third', 'day', 'on', 'rate', 'pause', 'bets']

round 4 :
e 0 = ['update', 'erdogan', 'weigh', 'wall', 'style', 'ask', 'sales', 'in', 'tumble', '
through', 'olympics', 'came']
e 1 = ['too', 'update', 'golden', 'corruption', 'trade', 'pickup', 'sales', 'from', 'in
', 'mlb', 'january']
e 2 = ['investors', "'s", 'russian', 'responses', 'to', 'output', 'ups', 'blackouts', '
as', 'decade', 'wave']
e 3 = ['dollar', 'to', 'islands', 'party', 'in', 'slump', 'since', 'cuts', 'job', 'tsit
sipas', 'u.s.-china']
e 4 = ['dollar', 'falls', 'for', 'a', 'after', 'third', 'proposal', 'year', 'rate', 'ta
lk']

round 5 :
e 0 = ['uk', 'survey', 'toward', 'wing', 'for', 'santos', 'fans', 'nearly', 'continue',
'retreats', ':']
e 1 = ['update', 'u.s.-iran', 'ups', '_num_-dutch', 'u.s.', 'disorder', '_num_-delta',
',', 'and', '-baker', 'optimism', 'lebanese']
e 2 = ['investors', "'s", 'paso', 'for', 'depend', 'new', 'revenue', 'pound', 'standoff
', 'rate', 'outperforms']
e 3 = ['to', 'breakingviews', 'opec', 'cup', 'of', 'ftse', 'leads', 'weak', 'four', 'sa
mples', 'ecb']
e 4 = ['dollar', 'falls', 'for', 'a', 'third', 'day', 'on', 'rate', 'pause', 'bets']
```

```
e 0 = ['uk', 'reverses', 'update', 'pleads', '-study', 'or', 'science', 'in', 'sink', 'sales', 'mulls', 'partial']
e 1 = ['update', 'unlikely', 'hopes', 'pills', 'new', 'disorder', 'nasdaq', ';', 'march', 'imo', 'standoff']
e 2 = ['erdogan', ',', 'offsets', '_num_', 'nexon', 'semis', 'in', 'prices', 'key', 'and', 'texans']
e 3 = ['dollar', 'rising', 'to', 'damascus', 'gdp', 'after', 'steady', 'as', 'dubai', 'january']
e 4 = ['dollar', 'falls', 'for', 'a', 'third', 'day', 'on', 'rate', 'pause', 'bets']
```