

CSC413 Assignment 3. Text Denoising Autoencoder for News Headlines

Deadline: TBD

Submission: Submit a PDF report containing your code, outputs, and your written solutions. You may export the completed notebook, but if you do so **it is your responsibility to make sure that your code and answers do not get cut off.**

Late Submission: Please see the syllabus for the late submission criteria.

Working with a partner: You may work with a partner for this assignment. If you decide to work with a partner, please create your group on Markus by the deadline date, even if you intend to use grace tokens. Markus does not allow you to create groups past the deadline, even if you have grace tokens remaining.

In this assignment, we'll explore a more advanced use of deep learning on a natural language process (NLP) task involving news headlines. In particular, we'll be working with a dataset of Reuters news headlines collected over a span of 15 months, covering some of 2018, 2019, and early 2020. This assignment will combine several of the concepts that we discussed in class, including recurrent neural networks, data augmentation, autoencoders (soon), and working with embeddings.

To be more specific, we'll be building an **autoencoder** of news headlines. We will build an **encoder** model that maps a news headline to a vector embedding, and a **decoder** that reconstructs the news headline. By building a model that learns to reconstruct the news headlines from the vector embedding, the model will learn good embeddings of these headlines.

We'll see a similar idea with image autoencoders and image VAEs, but both our encoder and decoder networks will be Recurrent Neural Networks. You'll have a chance build networks that takes a sequence as an input, and a network that generates a sequence as an output.

This project is organized as follows:

- Question 1. Data exploration
- Question 2. Background Math
- Question 3. Building the autoencoder
- Question 4. Training the autoencoder using *data augmentation*
- Question 5. Analyzing the embeddings (interpolating between headlines)
- Question 6. Work Allocation

Much of the idea behind this assignment is motivated by Shen et al [1]. We'll use the data augmentation rules proposed in that work to improve the robustness of the autoencoder.

[1] Shen et al (2019) "Educating Text Autoencoders: Latent Representation Guidance via Denoising" <https://arxiv.org/pdf/1905.12777.pdf>

```
In [1]:  
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  
  
%matplotlib inline
```

Question 1

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 [2]:  
from google.colab import drive  
drive.mount('/content/gdrive')  
  
train_path = '/content/gdrive/My Drive/CSC413/a3/reuters_train.txt' # Update me  
valid_path = '/content/gdrive/My Drive/CSC413/a3/reuters_valid.txt' # Update me
```

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

We will be using PyTorch's `torchtext` utilities to help us load, process, and batch the data. This package is useful, but takes a bit of time to get used to.

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 [4]:  
import torchtext  
  
# 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()  
    # the .split() will split the sequence/headline into words  
  
# Data field (column) representing our *text*.  
text_field = torchtext.legacy.data.Field(  
    # i changed this from torchtext.data.Field to torchtext.legacy.data.Field  
    sequential=True, # this field consists of a sequence  
    tokenize=tokenize_headline, # how to split sequences into words (use the fun  
    include_lengths=True, # to track the length of sequences, for batching  
    batch_first=True, # similar to batch_first=True in nn.RNN demonstr
```

```
use_vocab=True) # to turn each character into an integer index
train_data = torchtext.legacy.data.TabularDataset(
    # I changed this from torchtext.data.TabularDataset to torchtext.legacy.data
    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) -- 2 points

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 [5]:

```
# Include your histogram and your written explanations

# Here is an example of how to plot a histogram in matplotlib:
# plt.hist(np.random.normal(0, 1, 40), bins=20)

# Here are some sample code that uses the train_data object:
print(len(train_data[5].title) - 2)
print(train_data[5].title)
num_words_per_headline = []
for example in train_data:
    num_words = (len(example.title) - 2) # number of words in this headline. -2 to
    num_words_per_headline.append(num_words)
    # print(example.title)

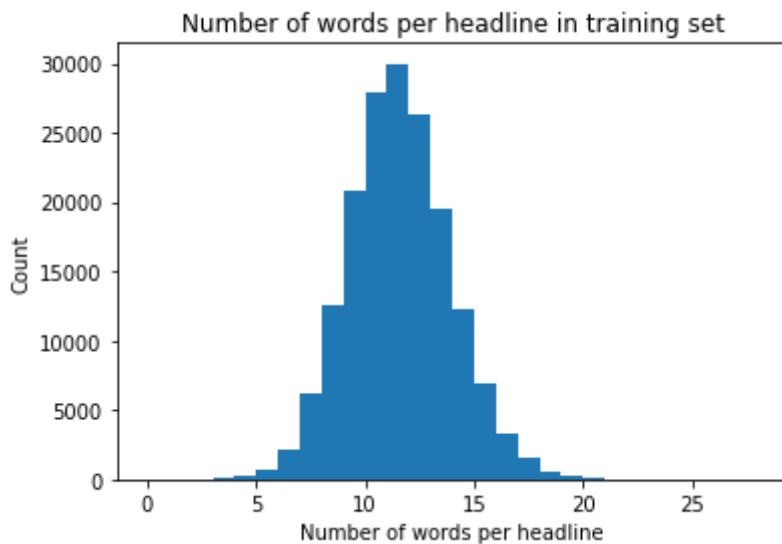
plt.hist(num_words_per_headline, bins = [i for i in range(max(num_words_per_headline)+1)])
plt.title('Number of words per headline in training set')
plt.xlabel('Number of words per headline')
plt.ylabel('Count')

# explanation
# The reason why we are interested in the histogram is that it allows us to see
# the lengths of the headlines. We can then make predictions and make choices on
# the how the model looks like such as when choosing the sizes of the hidden
# layer or embedding size. Such things are affected by sequence length, if we'd
# like to remember a certain headline of length n, we know then that we would
# need an embedding size of at least n. Seeing the lengths of the headlines also
# gives some insight on padding/batching should we want to pad/batch, since we
# would want to pad the shorter sequences to the length of the longest sequence
# in the batch.
```

11

```
['<bos>', 'u.s.', 'navy', 'pursuing', 'block', 'buy', 'of', 'two', 'aircraft',
'carriers', '-', 'senator', '<eos>']
```

Out[5]: Text(0, 0.5, 'Count')



Part (b) -- 2 points

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

```
In [6]: # 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

c = Counter()
for headline in train_data:
    # print(headline.title)
    words = headline.title
    for word in words:
        c[word] += 1

num_distinct_words = str(len(c) - 2)
print("Number of distinct words in training data: " + num_distinct_words)
```

Number of distinct words in training data: 51298

Part (c) -- 2 points

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?

```
In [7]: # Report your values here. Make sure that you report the actual values,
# and not just the code used to get those values

appears_exactly_once = 0 # number of words that appear exactly one time
appears_exactly_twice = 0 # number of words that appear exactly two times
for word in c: # keys
    if c[word] == 1:
        appears_exactly_once += 1
    if c[word] == 2:
```

```

    appears_exactly_twice += 1

print("Number of words that appear exactly once in the training set: " +
      str(appears_exactly_once))
print("Number of words that appear exactly twice in the training set: " +
      str(appears_exactly_twice))

```

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) -- 2 points

Explain why we may wish to replace these infrequent words with an `<unk>` tag, instead of learning embeddings for these rare words. (Hint: Consider words in the validation set that might not appear in training)

In [9]:

```

# Include your explanation here
# The <unk> tag allows the model to pick up on unknown words not in our
# vocabulary
# and therefore embed those unknown/infrequent words. A set such as
# the validation set would contain
# these kind of words that wouldn't be seen in the training set and we would
# need to embed them and we can use the <unk> tag to do so. This will allow us
# to find embeddings for headlines whose words are infrequent or unknown.
# without this <unk> system in place, we would have no embeddings for these
# kind of words and we wouldn't be able to get anywhere,
# and these infrequent/unknown words will be excluded by the model,
# the model will not pick them up which will affect our predictions greatly.

```

Part (e) -- 2 points

We will only model the top 9995 words in the training set, excluding the tags `<bos>`, `<eos>`, and other possible tags we haven't mentioned yet (including those, we will have a vocabulary size of exactly 10000 tokens).

What percentage of word occurrences will be supported? Alternatively, what percentage of word occurrences in the training set will be set to the `<unk>` tag?

In [10]:

```

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

del c["<eos>"]
del c["<bos>"]

# total occurrences of words in the whole training data set
all_occurrences = sum(c.values())

# get the total number of times a word appears in the top 9995
num_occurrences_top_9995 = 0
for pair in c.most_common(9995):
    num_occurrences_top_9995 += pair[1]

percentage_supported_words = num_occurrences_top_9995 / all_occurrences

```

```
percentage_words_not_supported = 1 - percentage_supported_words
print("percentage of word occurrences that will be supported: ",
      100 * percentage_supported_words)
print("percentage of word occurrences that will be set to <unk> tag (not supported): "
      100 * percentage_words_not_supported)
```

```
percentage of word occurrences that will be supported: 93.97857393100142
percentage of word occurrences that will be set to <unk> tag (not supported): 6.021426068998581
```

Our `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. You can interpret these indices as sparse representations of one-hot vectors.

In [11]:

```
# 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
# for instances, we can convert from string to (unique) index
print(vocab.stoi["hello"])
print(vocab.itos[10])      # ... and from word index to string

# The size of our vocabulary is actually 10000
vocab_size = len(text_field.vocab.stoi)
print(vocab_size) # should be 10000

# The reason is that torchtext adds two more tokens 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
10000
<unk>
<pad>
```

Question 2

Choosing the right model architecture is key for any successful deep learning system. In this question, we will compare the learning performance of RNNs and GRUs from the perspective of the vanishing/exploding gradient problem that arises during backpropagation.

Part (a) -- 4 pts

First, we will analyze the recurrent weight matrix of an RNN using Singular Value Decomposition (SVD). SVD says that any real matrix $M \in \mathbb{R}^{m \times n}$ can be written as $M = U \Sigma V^T$ where $U \in \mathbb{R}^{m \times m}$ and $V \in \mathbb{R}^{n \times n}$ are square orthogonal matrices, and $\Sigma \in \mathbb{R}^{m \times n}$ is a rectangular diagonal matrix. Recall that the values of Σ are the eigenvalues of $M^T M$.

(For a quick overview of eigenvalues and eigenvectors, see <https://www.youtube.com/watch?v=PFDu9oVAE-g> 0:44-5:22 and 13:05-end. The last section explains visually why we are

interested in working with eigenvalues when working with an RNN)

Consider a simple simple RNN-like architecture that computes $x_{t+1} = \text{sigmoid}(W x_t)$. You can view this architecture as a deep fully connected network that uses the same weight matrix at each layer. Suppose the largest singular value of the weight matrix is $\sigma_{\max}(W) = \frac{1}{2}$.

Show that the largest singular value of the input-output Jacobian has the following bound: $\sigma_{\max}(\partial \mathbf{x}_n / \partial \mathbf{x}_1) \leq (\frac{1}{2})^n$. $\text{Hint: if } C = AB, \text{ then } \sigma_{\max}(C) \leq \sigma_{\max}(A) \sigma_{\max}(B)$. Also, the input-output Jacobian is the multiplication of layerwise Jacobians).

What does this tell us about the input-output Jacobian $\frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_1}$ as $n \rightarrow \infty$?

2. (a)

SVD: $W = U \Sigma V^T$, a matrix $M \in \mathbb{R}^{m \times n}$ can be rewritten like this.

WT matrix

$$W = \underbrace{U}_{m \times m} \underbrace{\sum}_{m \times n} \underbrace{V^T}_{n \times n}$$

Diagonal matrix

The sum of the eigenvalues of $W^T W$'s and WW^T 's are the singular values in Σ .

backpropagating:

$$\begin{aligned} \frac{\partial L}{\partial W} &= \sum_{t=1}^T \frac{\partial L_t}{\partial W} \\ &= \sum_{t=1}^T \frac{\partial L_t}{\partial z_t} \frac{\partial z_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial W} \end{aligned}$$

* $\frac{\partial h_t}{\partial h_k} = \underbrace{\frac{\partial h_t}{\partial h_{t-1}}}_{\substack{\text{Jacobian matrix,} \\ \text{derivative of the hidden state} \\ \text{at time } t \text{ wrt hidden state} \\ \text{at time } k}}$... $\frac{\partial h_{t+1}}{\partial h_t}$

$$= \prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}}$$

Jacobian matrix

we know $\frac{\partial h_{t+1}}{\partial h_t} = \sigma'(W \cdot h_t) \cdot W$ (chain rule)

because $h_{t+1} = \sigma(W \cdot h_t)$ (given).

$\frac{\partial h_t}{\partial h_k} = \prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=k+1}^t \sigma'(W \cdot h_{i-1}) \cdot W$

We know $\frac{1}{2}$ is the max sing. value of W , so $\frac{1}{2}$ is a eigenvalue of W (the largest eigenvalue of W). Multiplying these Jacobians results in the

1 of 8
2 of 8

eigenvalues of the Jacobian matrices getting multiplied too. Since $(\frac{1}{2})^2$ is an eigenvalue of the Jacobian matrix $\frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_1}$, and since $\frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_1}$ can be

written as n Jacobian matrices getting multiplied with each other, we have that the eigen value $(\frac{1}{2})^2$ will be raised to the n^{th} power $\left[\left(\frac{1}{2}\right)^2\right]^n = \left(\frac{1}{2}\right)^{2n}$. So the largest singular

value of the Jacobian matrix is $\sqrt{\left(\frac{1}{2}\right)^{2n}} = \left(\frac{1}{2}\right)^n$, because the square roots of the eigenvalues of a matrix written in SVD form, are the singular values of that matrix, and the Jacobian matrix is written in SVD because W has been rewritten using SVD, and we've seen that the Jacobian matrices can be written in terms of W .

So we know for sure that the largest

singular value of the Jacobian matrix $\frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_1}$ is $\frac{1}{2}$ due to the fact that $\frac{1}{2}$ is the largest singular value of W and that

a Jacobian matrix can be written in terms of W ie $\frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_{n-1}} = \sigma'(W \cdot h_{t-1}) \cdot W$.

\therefore we have that $\sigma'_{\max}\left(\frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_1}\right) \leq \left(\frac{1}{2}\right)^n$.

We also know that $\sigma'_{\max}\left(\frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_1}\right) \geq 0$

because if $\sigma'_{\max}(W) = \frac{1}{2}$, then the largest eigenvalue of WTW (or WW^T) is $\left(\frac{1}{2}\right)^2$, and the eigenvalues of W is equivalent to the eigenvalues of W^T , so when we multiply both W and W^T , their eigenvalues will multiply as well (the eigenvalues will be squared) so the eigenvalues can never be negative which means

$\sigma'_{\max}\left(\frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_1}\right) \geq 0$

$\therefore 0 \leq \sigma'_{\max}\left(\frac{\partial \mathbf{x}_n}{\partial \mathbf{x}_1}\right) \leq \left(\frac{1}{2}\right)^n$

4 of 8

Explanation:

As n goes infinity, the Jacobian's singular values will be getting closer and closer to 0. This will cause the gradient to vanish because as the number of hidden states n increase, the more times the weight matrix W will get multiplied with each other (in the equation $\frac{\partial h_t}{\partial h_k} = \prod_{i=k+1}^t \frac{\partial h_i}{\partial h_{i-1}} = \sum_{i=k+1}^t W \cdot \dots \cdot W \cdot h_{i-1}$), which means the eigenvalues of W will be raised to the power of n (if we were to have n hidden states), but since n is going to infinity, the eigenvalues will become the dominating term since it's getting raised to a very high power. Since the eigenvalue here is $1/2$, raising $1/2$ to a large power will result in an extremely small number, and when we multiply this small number with the gradients of the hidden states ($W \cdot h_{i-1}$), the gradients of each hidden state at each time step will also become a small number, causing a vanishing gradient. These gradients make up the columns of the Jacobian matrix, so Jacobian matrix will vanish too.

Part (b) -- 4 pts

We will now compare the gradients of a vanilla RNN unit and a Gated Recurrent Unit (GRU). For both parts (b) and (c), assume that all weights are scalars, and that we have an input sequence of length T with $x_1 = 1$ and $x_t = 0$ for all other t . Also, assume that $h_0 = 0$ and that after T timesteps, we calculate the squared loss $L = \frac{1}{2}(y_T - o_T)^2$ where o_T is the target at timestep T .

Consider the vanilla RNN units that compute h_t at each timestep t as follows:

$$\$ \$ m_t = W_x x_t + W_h h_{t-1} \$ \$ \$ h_t = \tanh(m_t) \$ \$ \$ y_t = W_y h_t \$ \$$$

Compute $\frac{\partial L}{\partial W_x}$ using backpropagation. You should obtain an expression in terms of the quantities given (like o_T , y_T , etc...)

Do you see a vanishing gradient problem? What about exploding gradient? Explain.

(b) find $\frac{\partial L}{\partial W_x}$

Computing graph.

$\frac{\partial L_{t+1}}{\partial W_x} = \frac{\partial L_{t+1}}{\partial y_{t+1}} \cdot \frac{\partial y_{t+1}}{\partial h_{t+1}} \cdot \frac{\partial h_{t+1}}{\partial W_x}$ [①] the next hidden state (h_{t+1}) depends on the current hidden state (h_t).
 $+ \left[\frac{\partial L_{t+1}}{\partial y_{t+1}} \frac{\partial y_{t+1}}{\partial h_{t+1}} \frac{\partial h_{t+1}}{\partial h_t} \right] \frac{\partial h_t}{\partial W_x}$ [②] h_t depends on x_t .

$$L_t = \frac{1}{2} (y_t - o_t)^2$$

$$L_{t+1} = y_{t+1} - o_{t+1}$$

$$\frac{\partial y_{t+1}}{\partial h_{t+1}} = W_y$$

$$h_{t+1} = \tanh(m_{t+1})$$

$$m_{t+1} = W_x x_{t+1} + W_h h_t$$

$$\frac{\partial h_{t+1}}{\partial W_x} = \tanh(m_{t+1}) \frac{\partial (W_x x_{t+1} + W_h h_t)}{\partial W_x}$$

$$= \tanh(m_{t+1}) \cdot X_{t+1}$$

$$\frac{\partial h_{t+1}}{\partial h_t} = \frac{\partial \tanh(Wx X_{t+1} + Wh h_t)}{\partial h_t}$$

(C)

$$\frac{\partial h_t}{\partial Wx} = \frac{\partial \tanh(Wx X_t + Wh h_{t-1})}{\partial Wx}$$

$$\begin{aligned} \therefore \frac{\partial o_{t+1}}{\partial Wx} &= (y_{t+1} - o_{t+1}) W_y \cdot \frac{\partial (\tanh(Wx X_{t+1} + Wh h_t))}{\partial Wx} \\ &\quad + (y_{t+1} - o_{t+1}) \cdot W_y \frac{\partial (\tanh(Wx X_{t+1} + Wh h_t))}{\partial Wx} \cdot \frac{\partial \tanh(1)}{\partial x} \end{aligned}$$

Explanation

Yes there is a vanishing gradient problem here. In this gradient as the computed \$y_i\$ gets closer to either 0 or 1, the values of the terms getting multiplied with \$y_{t+1} - o_{t+1}\$ will get

smaller and make the gradient signal vanish. Also, since this is a vanilla RNN, we will constantly be applying the tanh activation function over and over again, which will cause cliffs, which will result in having both vanishing and exploding gradients in itself

Note that in my calculations, h_{t+1} is the next hidden state and h_t is the current hidden state, that's why my final answer is in terms of quantities like o_{t+1} and y_{t+1} as opposed to o_t and y_t .

Part (c) -- 4 pts

Now, let's consider GRU units that uses a gating mechanism, and computes h_t at each timestep t as follows:

$$\begin{aligned} z_t &= \sigma(W_x x_t + U_z h_{t-1}) \\ r_t &= \sigma(W_r x_t + U_r h_{t-1}) \\ h_t &= (1 - z_t) h_{t-1} + z_t \hat{h}_t \\ y_t &= W_y h_t \end{aligned}$$

Where σ is the sigmoid function.

Compute $\frac{\partial L}{\partial W_x}$ using backpropagation.

Can the vanishing gradient problem be prevented? *Hint* : Consider the term $\frac{\partial h_t}{\partial h_{t-1}}$, what role does z_t play in this gradient? Can it help alleviate the vanishing gradient problem?

Explanation

Yeah z_t can help alleviate the vanishing gradient problem.

Since $L = \sum L_t$, then we have that

$$\frac{\partial L}{\partial W_x} = \sum_{t=1}^n \frac{\partial L_t}{\partial W_x}$$

Because we are using squared loss and y_t is the same as the previous question,

$$\frac{\partial L_t}{\partial W_x} = \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial W_x} = (y_t - o_t) W_y \frac{\partial h_t}{\partial W_x}$$

Now we solve for $\frac{\partial h_t}{\partial W_x}$

$$\frac{\partial h_t}{\partial W_x} = \sum_{i=1}^t \frac{\partial h_t}{\partial h_i} \frac{\partial h_i}{\partial W_x} = \sum_{i=1}^t \left(\left(\prod_{j=1}^{t-1} \frac{\partial h_{j+1}}{\partial h_j} \right) \frac{\partial h_i}{\partial W_x} \right)$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial h_t}{\partial \hat{h}_t} \frac{\partial \hat{h}_t}{\partial h_{t-1}} + \frac{\partial h_t}{\partial z_t} \frac{\partial z_t}{\partial h_{t-1}} + \frac{\partial \overline{h_t}}{\partial h_{t-1}} = \frac{\partial h_t}{\partial \hat{h}_t} \left(\frac{\partial \hat{h}_t}{\partial r_t} \frac{\partial r_t}{\partial h_{t-1}} + \frac{\partial \overline{h_t}}{\partial h_{t-1}} \right) + \frac{\partial h_t}{\partial z_t} \frac{\partial z_t}{\partial h_{t-1}} + \frac{\partial \overline{h_t}}{\partial h_{t-1}}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = z_t (W_r ((W_h ((1-\hat{h}) \hat{h}) h_{t-1} r (1-r)) + ((W_h ((1-\hat{h}) \hat{h}) r_t) + W_z ((h_{t-1} \hat{h}_t z_t (1-z_t)) + (1-z)$$

Question 3

Building a text autoencoder is a little more complicated than an image autoencoder, so we'll need to thoroughly understand the model that we want to build before actually building our model. Note that the best and fastest way to complete this assignment is to spend a *lot* of time upfront understanding the architecture. The explanations are quite dense, and you might need to stop every sentence or two to understand what's going on. You won't feel productive for a while since you won't be writing code, but this initial investment will help you become more productive later on. Understanding this architecture will also help you understand other machine learning papers you might come across. So, take a deep breath, and let's do this!

Here is a diagram showing our desired architecture:

{ width=90% }



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 $\{\bf h\}^7$ 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 $\{\text{bf } h\}^{(7)}$) 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 $\{\text{bf } x\}^{(t)}$ to our RNN, and will be used to update the hidden state $\{\text{bf } m\}^{(t)}$. Since each of our sequences begin with the `<bos>` token, we'll set $\{\text{bf } x\}^{(1)}$ to be the `<bos>` token.
- 2) Compute the updates to the hidden state $\{\text{bf } m\}^{(t)}$ based on the previous hidden state $\{\text{bf } m\}^{(t-1)}$ and $\{\text{bf } x\}^{(t)}$. Intuitively, this hidden state vector $\{\text{bf } m\}^{(t)}$ 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 $\{\text{bf } m\}^{(t)}$, 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 $K=10000$ distinct words. As in a classification problem, the fully-connected neural network will compute a *probability distribution* over these 10,000 words. In the diagram, we are using $\{\text{bf } z\}^{(t)}$ 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 $\{\text{bf } z\}^{(t)}$. We can do this in a number of ways, which we'll discuss in question 4. 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 $\{\text{bf } x\}^{(t+1)}$ to our RNN, which is used to update our hidden state $\{\text{bf } m\}^{(t+1)}$ ---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.

Unfortunately, we can't *train* this autoencoder in the way we just described. That is, we can't just compare our generated sequence with our ground-truth sequence, and get gradients. Both sequences are **discrete** entities, so we won't be able to compute gradients at all! In particular, **sampling is a discrete process**, and so we won't be able to back-propagate through any kind of sampling that we do.

You might wonder whether we can get away with computing gradients by comparing the distributions $\{\text{bf } z\}^{(t)}$ with the ground truth words at each time step. Like any multi-class classification problem, we can represent the ground-truth words as a one-hot vector, and use the cross-entropy loss.

In theory, we can do this. In practice, there are a few issues. One is that the generated sequence might be longer or shorter than the actual sequence, meaning that there may be more/fewer $\{\text{bf } z\}^{(t)}$ 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 $\{\text{bf } z\}^{\{(t)\}}$, we will forgo sampling, and use the **ground truth** $\{\text{bf } x\}^{\{(t)\}}$ in the next step.

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

{ width=90% }



We will use the RNN generator to compute the logits $\{\text{bf } z\}^{\{(1)\}}, \{\text{bf } z\}^{\{(2)\}}, \dots, \{\text{bf } z\}^{\{(T)\}}$. 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 . (This is similar to what we did in a pixel-wise prediction problem.)

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, or initialize our parameters to GloVe embeddings. To prevent students who don't have Colab access from having to download a 1GB file, we won't do that. The word embedding component is represented with blue arrows in the diagram.
- The encoder RNN (which will use Gated Recurrent Units) 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 Gated Recurrent Units) 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** (one fully-connected layer) that computes a distribution over the next word to generate, given a decoder RNN hidden state.

Part (a) -- 8 pts

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`. (4 points)
2. Complete the `forward` method, which uses teacher forcing and computes the logits $\{z\}^{\{(t)\}}$ of the reconstruction of the sequence. (4 points)

You should first try to understand the `encode` and `decode` methods, which are written for you. The `encode` method mimics a discriminative RNN (see the sentiment analysis notebook). The `decode` method is a generative RNN and is a bit more complex (see the text generation tutorial notebook). 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 very careful of the input that you feed into either `decode` or to `self.decoder_rnn`. Refer to the teacher-forcing diagram.

In [12]:

```

class AutoEncoder(nn.Module):
    def __init__(self,
                vocab_size, # number of unique words in our vocabulary
                emb_size, # size of word embeddings ( $x^t$ )
                hidden_size # size of the hidden state in the RNNs
                ):
        super(self, object).__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):
        pass

```

```

"""
Computes the encoder output given a sequence of words.
"""

emb = self.embed(inp) # create an embedding of the input
# pass in the embedding into the encoder
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) # create an embedding of the input
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*
"""

# TODO
embedding = self.embed(inp) # need to create an embedding of the input
output_encoding, last_hidden_state_enc = self.encoder_rnn(embedding)
output_decoding, last_hidden_state_dec = self.decoder_rnn(embedding,
last_hidden_state_enc)

# predict distribution over next tokens
output = self.proj(output_decoding)
return output, last_hidden_state_dec

```

Part (b) -- 5 pts

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 [13]:

```

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

torch.Size([1, 9])
tensor([[ 2, 5258,   91, 9117,     6,    25,   637,   118,      3]])
['<bos>', 'zambian', 'president', 'swears', 'in', 'new', 'army', 'chief', '<eos>']

```

We are looking for the way that you set up your loss function corresponding to the figure above.

Be very careful of off-by-ones.

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

In [14]:

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

for it in range(300):

    # TODO
    optimizer.zero_grad() # set gradients to 0 before doing back prop
    # for the input, <EOS> is never an input token. in the first iteration,
    # target is bos, hidden is the one from bos
    output, hidden = model(input_seq[:, :-1])
    target = input_seq[:, 1:] # <BOS> is never a target token
    loss = criterion(output.reshape(-1, vocab_size), # reshape to 2D tensor
                     target.reshape(-1)) # dont include <bos> in input_seq
    loss.backward() # accumulates the gradient
    optimizer.step() # performs a parameter update based on the current gradient

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

```
[Iter 50] Loss 0.094657
[Iter 100] Loss 0.023950
[Iter 150] Loss 0.015026
[Iter 200] Loss 0.010384
[Iter 250] Loss 0.007750
[Iter 300] Loss 0.006060
```

Part (c) -- 2 pt

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 [15]:

```
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

hidden = model.encode(input_seq)

for i in range(5):
    print(sample_sequence(model, hidden))

```

```

['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief']

```

Part (d) -- 3 pt

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**.

```

In [16]: # Include the generated sequences and explanation in your PDF report.
print("temperature = 1.5")
for i in range(5):
    print(sample_sequence(model, hidden, temperature=1.5))

print("temperature = 2.0")
for i in range(5):
    print(sample_sequence(model, hidden, temperature=2.0))

print("temperature = 5.0")
for i in range(5):
    print(sample_sequence(model, hidden, temperature=5.0))

# By having the temperature too high, the words being sampled become too
# randomized. generated sequences will be more varies and will have lower
# quality, Seeing when the temperature is 5.0, the headlines generated
# don't make sense at all. We can also notice there are less <eos> tokens.

```

```

temperature = 1.5
['zambian', 'president', 'ventures', 'in', 'new', 'army', 'chief']
['deployment', 'deepen', 'rules', 'tampering', 'affects', 'seizes', 'itself', 'communications', 'slams', 'army', 'chief', 'potential', 'population', 'sam', 'men d', 'exposure', 'dangers', 'erdogan', 'hess', 'quitting']
['itv', 'army', 'chief', 'headaches', 'hung', '40,000', 'watchdogs', 'conviction', 'equality', 'e15', 'satellites', 'alstom', 'fm', 'army', 'chief', 'port', 'momentum', 'disappointed', 'cypress', 'inks']
['zambian', 'president', 'swears', 'in', 'new', 'army', 'urge', 'legend', 'chief']
['zambian', 'chip', 'stones', 'taste', 'influence', 'vetoes', 'sea', 'clients', 'ebbs', 'suppliers', 'democrat', 'wsj', 'kavanaugh', 'wounds', 'scam', 'garcia', 'slower', 'rep.', 'pricing', 'chief']
temperature = 2.0
['altria', 'knows', 'door', '-wsj', 'burden', 'grows', 'caps', 'flowers', 'mav s', 'attempt', 'coronation', 'westbrook', 'athletics', 'deliberations', 'were', 'aviation', 'toilet', 'entities', 'receives', 'adam']
['game', 'shipwreck', 'hikes', 'accusers', 'slovakia', 'profitability', 'swear s', 'in', 'orioles', 'blaming', 'battling', 'motegi', 'impeachment', 'hails', '_num_singapore', 'berrettini', 'charts', 'osram', 'library', 'colombian']
['responds', 'france-klm', 'parkland', 'transplant', 'bids', 'chest', 'bosnia', 'against', 'steadily', 'happens', 'army', 'broadly', 'khartoum', 'marine', 'kobe', '_num_sri', 'year', 'answer', 'hebei', 'mets']
['running', 'seize', 'withdraws', 'dips', 'global', 'warming', 'nurse', 'rule', 'leaves', 'dissent', 'catastrophe', '_num_s.korean', 'tells', 'forgotten', 'ind exes', 'new', 'movil', 'pan', 'shut', 'condemns']
['divide', 'affordable', 'libya', 'ebay', 'turnout', 'finally', 'wipro', 'uses', 'english', 'freeze', 'halts', 'dutch', 'vodafone', 'reversal', 'kidney', 'explosion', 'accusation', 'uproar', 'averting', 'lagerfeld']
temperature = 5.0
['relieved', 'somalia', 'permission', 'male', 'net', 'miami', 'legitimate', 'opp oses', 'friends', 'resorts', 'emotional', 'leonard', 'brent', 'high-tech', 'faci lity', 'terrible', 'tribune', 'hottest', 'news', 'studying']
['oyo', 'grade', 'arabiya', 'loading', 'forum', 'feared', 'journal', 'time', 'de cry', 'india', 'bruised', '_num_ohlahoma', 'cibc', 'ambulance', 'vans', 'conste llation', 'boc', 'corbyn', 'concert', 'fully']
['flaring', 'coronavirus-hit', 'serious', 'one-year', 'customs', 'concerns', 'tor nado', 'on', 'boards', 'fisher', 'outflow', 'either', 'treasuries', 'virtual', 'lethal', 'advisors', 'finishes', 'focus', 'gbagbo', 'dayton']
['atlantia', 'planemaker', 'mahomes', '_num_coronavirus', 'exercises', 'four-ye ar', 'were', 'input', 'billionaire', 'brain', 'head', '_num_boeing', 'lifescien ces', 'evacuates', 'game', 'airport', 'office', 'the', 'mavs', '_num_-1']
['tosses', 'ousts', 'forum', 'dominant', 'franchisee', 'stocks-tech', 'medal', 'cycling', 'assess', 'recommended', 'siemens', 'mission', 'partisan', 'denied', 'ngos', 'cosmetics', 'debated', 'trauma', 'charm', 'academic']

```

Question 4

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 (but not up/down---why?). Machine learning practitioners can also add Gaussian noise to the image.

When we use data augmentation to train an *autoencoder*, we typically 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) -- 3pt

Give three more examples of data augmentation techniques that we could use if we were training an **image** autoencoder. What are different ways that we can change our input?

```
In [17]: # Include your three answers

# smooth warping
# translation
# Rotation
```

Part (b) -- 2pt

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 [18]: 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
                                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] # removed 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 [19]:

```

# 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[42].title
print(train_data[42].title[1:-1])

for i in range(5):
    new_headline = tokenize_and_randomize(train_data[42].title[1:-1])
    new_headline = [vocab.itos[w] for w in new_headline]
    print(new_headline)

['zambian', 'president', 'swears', 'in', 'new', 'army', 'chief']
['<bos>', 'zambian', 'president', 'in', '<pad>', 'chief', 'army', '<eos>']
['<bos>', 'zambian', 'president', 'swears', '<pad>', 'new', 'army', 'chief', '<eos>']
['<bos>', 'zambian', 'swears', 'president', 'new', 'in', 'unicorn', 'chief', '<eos>']
['<bos>', "m", 'president', 'swears', 'knew', 'recover', 'chief', '<eos>']
['<bos>', 'zambian', 'president', 'in', 'army', 'new', 'chief', '<eos>']

```

Part (c) -- 3 pt

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 for longer if you want, but training tends 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 [20]:

```

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()

    for ep in range(num_epochs):
        # avg_loss = 0
        # We will perform data augmentation by re-reading the input each time
        field = torchtext.legacy.data.Field(sequential=True,
                                            tokenize=tokenize_and_randomize, # <-- data
                                            include_lengths=True,
                                            batch_first=True,
                                            use_vocab=False, # <-- the tokenization fun
                                            pad_token=vocab.stoi['<pad>'])
        dataset = torchtext.legacy.data.TabularDataset(train_path, "tsv", [('tit

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

    for it, ((xs, lengths), _) in enumerate(train_iter):
        # Fill in the training code here

        xs = xs.cuda()

        inp = xs[:, :-1] # <EOS> is never an input token
        # target = xs[:, 1:] # <BOS> is never an target token

        # forward pass
        # output, _ = model(inp) # TODO

        # calculate loss
        # loss = criterion( # TODO
            # output.reshape(-1, vocab_size), # reshape to 2D tensor
            # target.reshape(-1)
            # )

        # backward pass
        # loss.backward()
        # optimizer.step()

        # cleanup
        # optimizer.zero_grad()

    optimizer.zero_grad() # set gradients to 0 before doing back prop
    # for the input, <EOS> is never an input token. in the first
    # iteration, target is bos, hidden is the one from bos
    output, hidden = model(inp)
    target = xs[:, 1:] # <BOS> is never a target token
    loss = criterion(output.reshape(-1, vocab_size), # reshape to 2D ten

```

```

        target.reshape(-1)) # dont include <bos> in input_se
loss.backward() # accumulates the gradient
optimizer.step() # performs a parameter update based on the current
# gradient

# avg_loss += loss

if (it+1) % 100 == 0:
    print("[Iter %d] Loss %f" % (it+1, float(loss)))

# Optional: Compute and track validation loss
#val_loss = 0
#val_n = 0
#for it, ((xs, lengths), _) in enumerate(valid_iter):
#    zs = model(xs)
#    loss = None # TODO
#    val_loss += float(loss)

# Include your training curve or output to show that your training loss is trend
train_autoencoder(model, num_epochs = 1)

# Note: We don't know why our loss is like this, fluctuating so much and not
# going down, we followed the tutorial exactly and we really couldn't find any
# issue with the code we used above.

```

```

[Iter 100] Loss 2.464070
[Iter 200] Loss 2.243648
[Iter 300] Loss 2.175040
[Iter 400] Loss 2.384873
[Iter 500] Loss 2.555398
[Iter 600] Loss 2.588475
[Iter 700] Loss 2.304259
[Iter 800] Loss 2.632010
[Iter 900] Loss 2.755851
[Iter 1000] Loss 2.302110
[Iter 1100] Loss 2.180500
[Iter 1200] Loss 2.505726
[Iter 1300] Loss 2.266066
[Iter 1400] Loss 2.518975
[Iter 1500] Loss 2.337051
[Iter 1600] Loss 2.291297
[Iter 1700] Loss 2.555068
[Iter 1800] Loss 2.094347
[Iter 1900] Loss 2.580861
[Iter 2000] Loss 2.516101
[Iter 2100] Loss 2.108625
[Iter 2200] Loss 2.218971
[Iter 2300] Loss 2.255381
[Iter 2400] Loss 2.304690
[Iter 2500] Loss 2.104535
[Iter 2600] Loss 2.097020

```

Part (d) -- 2 pt

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 <https://www.cs.toronto.edu/~lczhang/321/files/p4model.pk> (11MB).

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

```
In [21]: model = AutoEncoder(10000, 128, 128)
checkpoint_path = '/content/gdrive/My Drive/CSC413/a3/p4model.pk' # Update me
model.load_state_dict(torch.load(checkpoint_path))
```

```
Out[21]: <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 [22]: # 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()

hidden = model.encode(input_seq)
# ...
print("temperature = 0.7")
for i in range(5):
    print(sample_sequence(model, hidden, temperature=0.7))

print("temperature = 0.9")
for i in range(5):
    print(sample_sequence(model, hidden, temperature=0.9))

print("temperature = 1.5")
for i in range(5):
    print(sample_sequence(model, hidden, temperature=1.5))

# With a too small temperature, when we want to generate new headlines,
# we will instead get the same repeated headlines. When the temperature
# is 0.7, we can see that the headlines are more similar than the higher
# temperatures. Low temperatures result in a more skewed distribution,
# the generated sequences won't be as varied which is something unattractive,
# we can see there is a lot of words repeated for temp = 0.7. We would like to
# have some different sentences but low temperatures won't allow for that.
```

```
temperature = 0.7
['wall', 'street', 'rises', ',', 'limps', 'die', 'win', 'at', '$', '<pad>', 'highway', 'investments', 'first']
['wall', 'street', 'rises', ',', 'limps', 'die', 'win', 'at', '$', 'says', 'invements', 'fire', 'signs']
['wall', 'street', 'rises', ',', 'limps', 'die', 'win', 'at', 'of', 'sciences', ':', 'face', 'signs']
['wall', 'street', 'rises', ',', 'limps', 'across', 'the', 'finish', 'line', ',', 'next', 'year', 'prize']
['wall', 'street', 'rises', ',', 'scales', 'vision', 'out', ',', 'update', '12th', '<pad>', 'election', 'sector']
temperature = 0.9
['wall', 'street', 'rises', ',', 'hut', 'out', 'conservatives', 'protesters', 'to', 'highest', 'fraser', 'after', 'strains']
['wall', 'street', 'rises', ',', 'limps', 'die', 'win', 'at', 'of', 'sciences', 'election', 'four']
['wall', 'street', 'rises', ',', 'limps', 'across', 'his', 'home', '<pad>', 'retailers', 'nears', 'after', 'round']
['wall', 'street', 'rises', ',', 'limps', 'die', 'win', 'at', 'of', 'sciences', '<pad>', 'presidential', 'cheers']
```

```
['wall', 'street', 'rises', ',', 'limps', 'across', 'melbourne', 'france', 'an',
 'to', 'peaking', 'jobs', '-']
temperature = 1.5
['wall', 'street', 'rises', ',', 'limps', 'race', 'greener', 'gather', 'to', 'te
ch', 'boko', 'sector', 'london']
['wall', 'street', 'percent', 'blast', 'stronger', 'blow', 'wealthy', 'among',
 'to', 'northeast', 'recession', 'big', 'server']
['wall', 'street', 'rises', ',', 'limps', 'die', 'win', 'estimates', 'unity', 'o
f', 'oil', 'rally', 'interested']
['wall', 'street', 'rises', 'beat', 'march', 'founding', 'violence', 'oregon',
 '<pad>', 'for', 'biggest', 'screens', ':']
['wall', "'s", 'ambev', 'arts', 'to', 'borrowers', 'birth', 'die', ',', 'junio
r', 'big', 'best', 'fire']
```

Question 5

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 [23]: valid_data = torchtext.legacy.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)
```

Part (a) -- 2 pt

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 [24]: # Write your code here
# Show that your resulting PyTorch tensor has shape `[19046, 128]`

emb_size = 128
# tensor that will hold all the embeddings of the validation set
validation_embeddings = torch.Tensor([])

for example in valid_data:
    headline = example.title
    input_seq = torch.Tensor([vocab.stoi[w] for w in headline
                           ]).long().unsqueeze(0)
    emb = model.encode(input_seq).squeeze(0)
    # embed = nn.Embedding(num_embeddings=vocab_size, embedding_dim=emb_size) # th
    # embedding = embed(input_seq)
    validation_embeddings = torch.cat((validation_embeddings, emb), dim=0)

validation_embeddings
```

```
Out[24]: tensor([[ 0.0826,  0.2645, -0.2565, ..., -0.1510, -0.8473, -0.0278],
                  [ 0.3108,  0.0070, -0.4099, ..., -0.5711, -0.7966, -0.8654],
                  [ 0.1921, -0.0109, -0.2046, ..., -0.1550, -0.8173, -0.1998],
                  ...,
                  [-0.2950, -0.0850, -0.0800, ...,  0.2101, -0.8050, -0.0324],
                  [ 0.4920,  0.1950, -0.6876, ...,  0.4365, -0.8328, -0.5217],
```

```
[ 0.3114, -0.0836, -0.1114, ..., -0.0497, -0.8331, -0.3255]],  
grad_fn=<CatBackward>)
```

In [25]: validation_embeddings.shape

Out[25]: torch.Size([19046, 128])

Part (b) -- 2 pt

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

In [26]:

```
# Write your code here. Make sure to include the actual 5 closest headlines.  
valid_data[13].title  
  
val_emb = validation_embeddings.detach().numpy()  
norms = np.linalg.norm(val_emb, axis=1)  
val_emb_norm = (val_emb.T / norms).T  
similarities = np.matmul(val_emb_norm, val_emb_norm.T)  
  
def closefive(word):  
    s = np.argsort(similarities[word])  
    cl5 = s[-6:-1]  
    cl5_list = []  
    for i in cl5:  
        cl5_list.append(valid_data[i].title)  
  
    cl5_list.reverse()  
    return cl5_list  
  
print(valid_data[13].title)  
for headline in closefive(13):  
    print(headline)  
  
['<bos>', 'asia', 'takes', 'heart', 'from', 'new', 'year', 'gains', 'in', 'u.  
s.', 'stock', 'futures', '<eos>']  
['<bos>', 'italy', "'s", 'salvini', 'loses', 'aura', 'of', 'invincibility', 'i  
n', 'emilia', 'setback', '<eos>']  
['<bos>', 'saudi', ',', 'russia', 'look', 'to', 'seal', 'deeper', 'output', 'cut  
s', 'with', 'oil', 'producers', '<eos>']  
['<bos>', 'eu', 'orders', 'quarantine', 'for', 'staff', 'who', 'traveled', 'to',  
'northern', 'italy', '<eos>']  
['<bos>', 'update', '_num_italy', "'s", 'prime', 'minister', 'says', 'new', 'go  
vernment', 'will', 'bicker', 'less', '<eos>']  
['<bos>', 'portugal', "'s", 'moura', 'pays', 'tribute', 'to', 'cod', 'fisherme  
n', 'at', 'milan', 'fashion', 'close', '<eos>']
```

Part (c) -- 2 pt

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

In [27]:

```
# Write your code here.  
# Make sure to include the original headline and the 5 closest headlines.  
print(valid_data[1113].title)  
for headline in closefive(1113):  
    print(headline)
```

```
[ '<bos>', 'emerging', 'markets-emerging', 'market', 'assets', 'jump', 'as', 'dovish', 'fed', 'boosts', 'sentiment', '<eos>' ]
[ '<bos>', 'emerging', 'markets-emerging', 'market', 'stocks', 'snap', 'five-session', 'losing', 'run', ',', 'currencies', 'firm', '<eos>' ]
[ '<bos>', 'emerging', 'markets-emerging', 'market', 'stocks', 'tick', 'up', ',', 'us-china', 'fears', 'cap', 'sentiment', '<eos>' ]
[ '<bos>', 'emerging', 'markets-emerging', 'market', 'stocks', 'dip', ',', 's.african', 'rand', 'soft', 'ahead', 'of', 'rate', 'meet', '<eos>' ]
[ '<bos>', 'us', 'stocks-wall', 'st', 'slips', 'as', 'railroads', 'slide', 'after', 'csx', 'signals', 'trade', 'impact', '<eos>' ]
[ '<bos>', 'us', 'stocks-wall', 'st', 'moves', 'lower', 'on', 'dampened', 'hopes', 'for', 'hefty', 'fed', 'cut', '<eos>' ]
```

Part (d) -- 4 pts

Choose two headlines from the validation set, and find their embeddings. We will **interpolate** between the two embeddings.

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:

$$\begin{aligned} e_1 &= 0.75 e_0 + 0.25 e_4 \\ e_2 &= 0.50 e_0 + 0.50 e_4 \\ e_3 &= 0.25 e_0 + 0.75 e_4 \end{aligned}$$

Decode each of e_1 , e_2 and e_3 five times, with a temperature setting that shows some variation in the generated sequences, while generating sequences that makes sense.

In [28]:

```
# Write your code here. Include your generated sequences.
print("Original headlines:\n")

print(valid_data[3000].title)
print(valid_data[2388].title)

print("\nInterpolated headlines:\n")

input1 = torch.Tensor([vocab.stoi[w] for w in valid_data[3000].title
                     ]).long().unsqueeze(0)
input2 = torch.Tensor([vocab.stoi[w] for w in valid_data[2388].title
                     ]).long().unsqueeze(0)

enc1 = model.encode(input1)
enc2 = model.encode(input2)

e1 = 0.75 * enc1 + 0.25 * enc2
e2 = 0.50 * enc1 + 0.50 * enc2
e3 = 0.25 * enc1 + 0.75 * enc2

temp = 0.9

print("Headline e1 \n")
for i in range(5):
    print(sample_sequence(model, e1, temperature=temp))

print("\nHeadline e2 \n")
for i in range(5):
    print(sample_sequence(model, e2, temperature=temp))
```

```
print("\nHeadline e3 \n")
for i in range(5):
    print(sample_sequence(model, e3, temperature=temp))
```

Original headlines:

```
['<bos>', 'golden', 'cross', 'for', 'stocks', 'does', "n't", 'always', 'glitter', '<eos>']
['<bos>', 'tornado', 'victim', 'recalls', 'total', 'devastation', '<eos>']
```

Interpolated headlines:

Headline e1

```
['genetic', '_num_-s.africa', 'florida', 'salesforce', '_num_-france', 'now', 'grande']
['cave', 'in', 'buyout', 'but', 'vs.', 'forward', ':']
['activists', 'magna', 'four-year', 'must', 'zambian', 'host', 'playoffs']
['nissan', 'mississippi', 'solar', 'nvidia', "n't", 'adviser', 'sukuk']
['rebuke', 'prizes', 'man', 'ends', 'researchers', 'appeal', '300,000']
```

Headline e2

```
['requiring', 'lead', 'policeman', 'misses', 'reopen', 'walker']
['requiring', 'africa', 'be', 'therapy', 'easyjet', 'bath']
['cross', 'international', 'croatian', 'march', 'probed', 'warriors']
['catalan', 'permian', "'s", 'tracking', 'fined', 'mcdonnell']
['blaze', 'endgame', 'migrants', 'lme', 'refund', 'completes']
```

Headline e3

```
['portuguese', 'australian', 'skills', 'debts', 'quarterly']
['portuguese', 'bomb', 'gaga', 'checks', 'angels']
['catalan', 'freeport', 'monitoring', 'rhp', 'estimates']
['portuguese', 'australian', 'responds', 'globally', 'magazine']
['chicago', 'update', 'computer', 'recalls', 'holding']
```

Question 6. Work Allocation -- 2 pts

This question is to make sure that if you are working with a partner, that you and your partner contributed equally to the assignment.

Please have each team member write down the times that you worked on the assignment, and your contribution to the assignment.

In [29]:

```
# Your answer goes here
# Team Members: Wafiqah Raisa and David Pham
# On March 30th, Wafiqah did the coding portions of Q1,
# and David wrote the explanation portions
# On March 31st, Wafiqah finished all of Q3
# On April 9th, Wafiqah finished Q2a and Q2b.
# On April 4th, David finished Q5 and all of Q4.
# On April 10th, David checked Wafiqah's work
# On April 11th, David finished Q2c
# On April 12th, Wafiqah checked David's work
```