

# HWK4 - N-gram and neural language models

March 20, 2023

CSCI 4140: Natural Language Processing

CSCI/DASC 6040: Computational Analysis of Natural Languages

Spring 2023

Homework 4 - N-gram and neural language models

Due Sunday, March 26, at 11:59 PM

Do not redistribute without the instructor's written permission.

## 1 Setup

Notes:

- You must run the code for Q2 on a computer with GPU (running it on CPU will take much, much longer). [Google Colab](#) is a good choice.
  - If you're using Colab, make sure you upload the `wiki` files.
  - If you're using other computer, update the path to `wiki` files (`fname = "...`).
- The neural language model may take up to 10 minutes to train, so **start early!**
- The rest of the cells are designed so that you can run them in a few minutes of computation time. If it is taking longer than that, you probably have made a mistake in your code.

```
[1]: import torch, pickle, os, sys, random, time
import torch.nn.functional as F
import matplotlib.pyplot as plt
from torch import nn, optim
from collections import *
import numpy as np
```

```
[2]: import torch

if torch.cuda.is_available():
    device = torch.device("cuda")          # a CUDA device object
    print('Using GPU:', torch.cuda.get_device_name())
else:
    device = torch.device("cpu")           # a CPU device object
    print('Using CPU')
```

Using GPU: NVIDIA GeForce RTX 3070 Ti

We'll start by loading the data. The WikiText language modeling dataset is a collection of tokens extracted from the set of verified Good and Featured articles on Wikipedia.

```
[3]: data = {'test': '', 'train': '', 'valid': ''}

for data_split in data:
    fname = "wiki.{}.tokens".format(data_split)
    with open(fname, 'r', encoding='utf8') as f_wiki:
        data[data_split] = f_wiki.read().lower().split()

vocab = list(set(data['train']))
```

Now have a look at the data by running this cell.

```
[4]: print('train : %s ...' % data['train'][:10])
print('dev : %s ...' % data['valid'][:10])
print('test : %s ...' % data['test'][:10])
print('first 10 words in vocab: %s' % vocab[:10])
```

```
train : ['=', 'valkyria', 'chronicles', 'iii', '=', 'senjō', 'no', 'valkyria',
'3', ':'] ...
dev : ['=', 'homarus', 'gammarus', '=', 'homarus', 'gammarus', ',', 'known',
'as', 'the'] ...
test : ['=', 'robert', '<unk>', '=', 'robert', '<unk>', 'is', 'an', 'english',
'film'] ...
first 10 words in vocab: ['stegosauria', 'diêm', 'frakes', 'onça', 'probes',
'rightful', 'tasked', 'ethiopia', 'apted', '970']
```

## 2 Q1. N-gram Language model (40pts)

### 2.1 Q1.1: Train N-gram language model (15pts)

Complete the following `train_ngram_lm` function based on the following input/output specifications. If you've done it right, you should pass the tests in the cell below.

*Input:* + **data**: the data object created in the cell above that holds the tokenized Wikitext data  
+ **order**: the order of the model (i.e., the “n” in “n-gram” model). If `order=3`, we compute  $p(w_2|w_0, w_1)$ .

*Output:* + **lm**: A dictionary where the key is the history and the value is a probability distribution over the next word computed using the maximum likelihood estimate from the training data. Importantly, this dictionary should include *backoff* probabilities as well; e.g., for `order=4`, we want to store  $p(w_3|w_0, w_1, w_2)$  as well as  $p(w_3|w_1, w_2)$  and  $p(w_3|w_2)$ .

Each key should be a single string where the words that form the history have been concatenated using spaces. Given a key, its corresponding value should be a dictionary where each word type in the vocabulary is associated with its probability of appearing after the key. For example, the entry for the history ‘w1 w2’ should look like:

```
lm['w1 w2'] = {'w0': 0.001, 'w1' : 1e-6, 'w2' : 1e-6, 'w3': 0.003, ...}
```

In this example, we also want to store `lm['w2']` and `lm['']`, which contain the bigram and unigram distributions respectively.

*Hint:* You might find the `defaultdict` and `Counter` classes in the `collections` module to be helpful.

```
[5]: def train_ngram_lm(data, order=3):
    """
        Train n-gram language model
    """

    # pad (order-1) special tokens to the left
    # for the first token in the text
    order -= 1
    data = ['<S>'] * order + data #
    lm = defaultdict(Counter)
    for o in range(order + 1):
        for i in range(len(data) - o):
            history = ' '.join(data[i:i+o])
            next_word = data[i+o]
            lm[history][next_word] += 1

    # Normalize counts to probabilities
    for history, next_word_counts in lm.items():
        total_count = sum(next_word_counts.values())
        for next_word, count in next_word_counts.items():
            lm[history][next_word] = count / total_count

    return lm
```

```
[6]: def test_ngram_lm():

    print('checking empty history ...')
    lm1 = train_ngram_lm(data['train'], order=1)
    assert '' in lm1, "empty history should be in the language model!"

    print('checking probability distributions ...')
    lm2 = train_ngram_lm(data['train'], order=2)
    sample = [sum(lm2[k].values()) for k in random.sample(list(lm2), 10)]
    assert all([a > 0.999 and a < 1.001 for a in sample]), "lm[history][word]_
↳should sum to 1!"

    print('checking lengths of histories ...')
    lm3 = train_ngram_lm(data['train'], order=3)
    assert len(set([len(k.split()) for k in list(lm3)])) == 3, "lm object_
↳should store histories of all sizes!"

    print('checking word distribution values ...')
```

```

    assert lm1['']['the'] < 0.064 and lm1['']['the'] > 0.062 and \
           lm2['the']['first'] < 0.017 and lm2['the']['first'] > 0.016 and \
           lm3['the first']['time'] < 0.106 and lm3['the first']['time'] > 0.
↪105, \
           "values do not match!"

    print("Congratulations, you passed the ngram check!")

test_ngram_lm()

```

```

checking empty history ...
checking probability distributions ...
checking lengths of histories ...
checking word distribution values ...
Congratulations, you passed the ngram check!

```

## 2.2 Q1.2: Generate text from n-gram language model (10pts)

Complete the following `generate_text` function based on these input/output requirements:

*Input:*

- **lm**: the lm object is the dictionary you return from the `train_ngram_lm` function
- **vocab**: vocab is a list of unique word types in the training set, already computed for you during data loading.
- **context**: the input context string that you want to condition your language model on, should be a space-separated string of tokens
- **order**: order of your language model (i.e., “n” in the “n-gram” model)
- **num\_tok**: number of tokens to be generated following the input context

*Output:*

- generated text, should be a space-separated string

*Hint:*

After getting the next-word distribution given history, try using `numpy.random.choice` to sample the next word from the distribution.

```

[7]: # generate text
def generate_text(lm, vocab, context="he is the", order=3, num_tok=25):

    # The goal is to generate new words following the context
    # If context has more tokens than the order of lm,
    # generate text that follows the last (order-1) tokens of the context
    # and store it in the variable `history`
    order -= 1
    history = context.split()[-order:]
    # `out` is the list of tokens of context
    # you need to append the generated tokens to this list
    out = context.split()

```

```

for i in range(num_tok):
    # Get the probability distribution over the next word given the history
    if tuple(history) in lm:
        dist = lm[tuple(history)]
    else:
        # If the history is not in the lm, we choose a random word from the
        ↪ vocabulary
        dist = np.ones(len(vocab))/len(vocab)
        # Sample the next word from the distribution
        next_word = np.random.choice(vocab, p=dist)
        # Append the next word to the output
        out.append(next_word)
        # Update the history by removing the first word and adding the next word
        history = history[1:] + [next_word]

    # Concatenate the tokens in the `out` list into a single string and return
    ↪ it
return ' '.join(out)

```

Now try to generate some texts! Read the texts generated by ngram language model with different orders

```

[8]: order = 1
generate_text(train_ngram_lm(data['train'], order=order), vocab, context='he is
    ↪ the', order=order)

```

```

[8]: 'he is the convoys llosa frightened 1215 temporal afrodisiac insects flour asia
    falkland jowell celebrities winter slept coal wasn youth veteran umbo training
    copulation mentioning honeymoon saturn lee'

```

```

[9]: order = 2
generate_text(train_ngram_lm(data['train'], order=order), vocab, context='he is
    ↪ the', order=order)

```

```

[9]: 'he is the emigration patriarchal across atenism egitto interpretive balloon
    thorns scofield data sexes berengaria proving switzerland antics rightful jai
    disguised arcade reliably worthless abe pursuit sivaji war'

```

```

[10]: order = 3
generate_text(train_ngram_lm(data['train'], order=order), vocab, context='he is
    ↪ the', order=order)

```

```

[10]: 'he is the technology helix privy maniac spectrograph nectar denouncing
    superfluous idiosyncratic polite mycena pressured deterioration valentin posse
    1619 cosmetic aralt mayfair opener sacred this mearns subset unemployed'

```

```
[11]: order = 4
generate_text(train_ngram_lm(data['train'], order=order), vocab, context='he is_
↳the', order=order)
```

```
[11]: 'he is the istanbul 176 graduates heart stereotype supergiants affray glasgow
cruel scaffidi reason hubbardton desserts ridot extreme peronism diva expresses
squadron inisfallen liga lakota generosity grow francetić'
```

### 2.3 Q1.3 : Evaluate the models (15pts)

Now let's evaluate the models quantitatively using the intrinsic metric **perplexity**.

Recall perplexity is the inverse probability of the test text

$$PP(w_1, \dots, w_t) = P(w_1, \dots, w_t)^{-\frac{1}{T}}$$

For an n-gram model, perplexity is computed by

$$PP(w_1, \dots, w_t) = \left[ \prod_{t=1}^T P(w_t | w_{t-1}, \dots, w_{t-n+1}) \right]^{-\frac{1}{T}}$$

To address the numerical issue (underflow), we usually compute

$$PP(w_1, \dots, w_t) = \exp \left( -\frac{1}{T} \sum_i \log P(w_t | w_{t-1}, \dots, w_{t-n+1}) \right)$$

*Input:*

- **lm**: the language model you trained (the object you returned from the `train_ngram_lm` function)
- **data**: test data
- **vocab**: the list of unique word types in the training set
- **order**: order of the lm

*Output:*

- the perplexity of test data

*Hint:*

- If the history is not in the **lm** object, back-off to (n-1) order history to check if it is in **lm**. If no history can be found, just use  $1/|V|$  where  $|V|$  is the size of vocabulary.

```
[12]: import math
def compute_perplexity(lm, data, vocab, order=3):
    # pad according to order
    order -= 1
    data = ['<S>'] * order + data
    log_sum = 0
    N = len(data) - order
```

```

for i in range(N):
    h, w = ' '.join(data[i: i+order]), data[i+order]
    # if h not in lm, back-off to n-1 gram and look up again
    while order > 0 and h not in lm:
        order -= 1
        h = ' '.join(data[i: i+order])
    if h in lm:
        p = lm[h].get(w, 1/len(vocab))
    else:
        p = 1/len(vocab)
    log_sum += math.log(p)
# compute perplexity
perplexity = math.exp(-log_sum/N)
return perplexity

```

Let's evaluate the language model with different orders. You should see a decrease in perplexity as the order increases. As a reference, the perplexity of the unigram, bigram, trigram, and 4-gram language models should be around 795, 203, 141, and 130 respectively.

```

[13]: for o in [1, 2, 3, 4]:
        lm = train_ngram_lm(data['train'], order=o)
        print('order {} ppl {}'.format(o, compute_perplexity(lm, data['test'],
↪vocab, order=o)))

```

```

order 1 ppl 794.5377104541699
order 2 ppl 260.5186891747848
order 3 ppl 260.54821138804743
order 4 ppl 260.55785556146014

```

### 3 Q2. Neural language models (70pts)

In this part of the homework, we'll be using PyTorch to play around with neural language models. First, a quick warm up by implementing backpropagation within a *scalar* neural network. Then, you'll implement a neural language model using PyTorch's built-in modules.

Firstly, run the cell below to import pytorch and set up the gradient checking functionality.

```

[14]: import torch
import torch.nn as nn
device = torch.device('cpu')

# checks equality between your gradients and those from autograd
def gradient_check(params, your_gradient):
    all_good = True
    for key in params.keys():
        if params[key].grad.size() != your_gradient[key].size():
            print('GRADIENT ERROR for parameter %s, SIZE ERROR\
↪your size: \
↪%s\nactual size: %s\n' \

```

```

        % (key, your_gradient[key].size(),
            params[key].grad.size()))
    all_good = False
    elif not torch.allclose(params[key].grad, your_gradient[key],
↪atol=1e-6):
        print('GRADIENT ERROR for parameter %s, VALUE ERROR\
nyours:↪
↪%s\nactual: %s\n'\
            % (key, your_gradient[key].detach(),
                params[key].grad))
    all_good = False

return all_good

```

### 3.1 Q2.1 Warm up with single neuron (10 pts)

The following code cell trains a network with scalars (i.e., single neurons) in each layer on a small dataset of ten examples. All you have to do is translate the partial derivatives we computed into code. The network is defined as:

$$h = \tanh(w_1 \cdot \text{input})$$

$$\text{pred} = \tanh(w_2 \cdot h)$$

$$L = 0.5 \cdot (\text{target} - \text{pred})^2$$

If you run the cell below, you should see “GRADIENT ERRORS”. Once you implement the partial derivatives  $\frac{\partial L}{\partial w_1}$  and  $\frac{\partial L}{\partial w_2}$  correctly, you will instead see a “SUCCESS” message. **Do NOT modify any code outside of the block marked “IMPLEMENT BACKPROP HERE”!**

```

[15]: # initialize model parameters
params = {}
params['w1'] = torch.randn(1, 1, requires_grad=True) # input > hidden with↪
↪scalar weight w1
params['w2'] = torch.randn(1, 1, requires_grad=True) # hidden > output with↪
↪scalar weight w2

# set up some training data
inputs = torch.randn(20, 1)
targets = inputs / 2

# training loop
all_good = True
for i in range(len(inputs)):

    ## forward prop, then compute loss.
    a = params['w1'] * inputs[i] # intermediate variable, following lecture↪
↪notes
    hidden = torch.tanh(a)
    b = params['w2'] * hidden

```



```

pred = torch.tanh(b)
loss = 0.5 * (targets[i] - pred) ** 2 # compute square loss
loss.backward() # runs autograd

#####
# TODO: IMPLEMENT BACKPROP HERE
# DO NOT MODIFY ANY CODE OUTSIDE OF THIS BLOCK!!!!
your_gradient = {}
your_gradient['w1'] = torch.zeros(params['w1'].size()) # implement dL/dw1
your_gradient['w2'] = torch.zeros(params['w2'].size()) # implement dL/dw2

# compute gradients
a = params['w1'] * inputs[i] # intermediate variable, following lecture
↳notes
hidden = torch.tanh(a)
b = params['w2'] * hidden
pred = torch.tanh(b)

dL_db = (pred - targets[i]) * (1 - torch.tanh(b) ** 2)
dL_dh = dL_db * params['w2'].item() * (1 - torch.tanh(a) ** 2)

your_gradient['w2'] += dL_db * hidden
your_gradient['w1'] += dL_dh * inputs[i] * (1 - torch.tanh(a) ** 2)

# END
#####

if not gradient_check(params, your_gradient):
    all_good = False
    break

# zero gradients after each training example
params['w1'].grad.zero_()
params['w2'].grad.zero_()

if all_good:
    print('SUCCESS! you passed the gradient check.')

```

GRADIENT ERROR for parameter w1, VALUE ERROR

yours: tensor([[ -0.0013]])

actual: tensor([[ -0.0050]])

### 3.2 Q2.2 RNN language model (20 pts)

For this part of the homework, we will use **PyTorch** to build our model. The following code cell preprocesses the raw text so you can load it directly. The input to your model is a *minibatch* of

sequences which takes the form of a  $N \times L$  matrix where  $N$  is the batch size and  $L$  is the maximum sequence length. For each minibatch, your models should produce an  $N \times L \times V$  tensor where  $V$  is the size of the vocabulary. This tensor stores the predicted probability distribution of the next word for every position of every sequence in the batch. Note that each batch is padded to dimensionality  $L = 40$  using the special padding token `<pad>`; similarly, each sequence begins with the `<bos>` token and ends with the `<eos>` token. Please look at the [PyTorch RNN documentation](#) if you're having problems getting started.

First, run the following code cell to download the data.

Please change your Colab runtime to the GPU backend by going to “Runtime > Change runtime type > Hardware accelerator > GPU”.

```
[16]: import torch, pickle, os, sys, random, time
      from torch import nn, optim

      device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
      print('device: ', device)

      # Load id2word from wikitext pickle
      with open('wikitext.pkl', 'rb') as f_in:
          wikitext = pickle.load(f_in)

      wikitext['train'] = torch.LongTensor(wikitext['train']).to(device)
      wikitext['dev'] = torch.LongTensor(wikitext['valid']).to(device)
      wikitext['test'] = torch.LongTensor(wikitext['test']).to(device)
      idx_to_word = wikitext['id2word']
      word_to_idx = {idx_to_word[k]: k for k in idx_to_word}

      print("Wikitext data loaded!")
      # Demonstrate id2word
      print('There are ' + str(len(idx_to_word)) + ' words in vocabulary')
      for id in range(8):
          print('Word id ' + str(id) + " stands for '" + str(idx_to_word[id]) + "'")
      print('...')
      print((wikitext['train'] > 0).sum())

      print('Set up finished')
```

```
device:  cuda
Wikitext data loaded!
There are 28654 words in vocabulary
Word id 0 stands for '<pad>'
Word id 1 stands for '<unk>'
Word id 2 stands for '<bos>'
Word id 3 stands for '<eos>'
Word id 4 stands for 'the'
Word id 5 stands for ','
Word id 6 stands for '.'
```

Word id 7 stands for 'of'

...

tensor(1622368, device='cuda:0')

Set up finished

The following cell contains code for computing perplexity and training the neural language model. Run the cell, and please make sure you (at least roughly) understand what is happening, but **do not modify any part of it**.

```
[17]: # function to evaluate LM perplexity on some input data, DO NOT MODIFY
def compute_perplexity(dataset, net, bsz=64):
    criterion = nn.CrossEntropyLoss(ignore_index=0, reduction='sum')
    num_examples, seq_len = dataset.size()

    # we'll still use batches because we can't fit the whole
    # validation set into GPU memory
    batches = [(start, start + bsz) for start in range(0, num_examples, bsz)]

    total_unmasked_tokens = 0. # count how many unpadded tokens there are
    nll = 0.
    for b_idx, (start, end) in enumerate(batches):
        batch = dataset[start:end]
        ut = torch.nonzero(batch).size(0)
        preds = net(batch)
        targets = batch[:, 1:].contiguous().view(-1)
        preds = preds[:, :-1, :].contiguous().view(-1, net.vocab_size)
        loss = criterion(preds, targets)
        nll += loss.detach()
        total_unmasked_tokens += ut

    perplexity = torch.exp(nll / total_unmasked_tokens).cpu()
    return perplexity.data

# training loop for language models, DO NOT MODIFY!
def train_lm(dataset, params, net):

    # since the first index corresponds to the PAD token, we just ignore it
    # when computing the loss
    criterion = nn.CrossEntropyLoss(ignore_index=0)

    optimizer = optim.Adam(net.parameters(), lr=params['learning_rate'])
    num_examples, seq_len = dataset.size()
    batches = [(start, start + params['batch_size']) for start in
                range(0, num_examples, params['batch_size'])]

    for epoch in range(params['epochs']):
        ep_loss = 0.
```

```

start_time = time.time()
random.shuffle(batches)
net.train()
# for each batch, calculate loss and optimize model parameters

↪
for b_idx, (start, end) in enumerate(batches):
    batch = dataset[start:end]
    preds = net(batch)

    preds = preds[:, :-1, :].contiguous().view(-1, net.vocab_size)
    targets = batch[:, 1:].contiguous().view(-1)
    loss = criterion(preds, targets)

    loss.backward()
    torch.nn.utils.clip_grad_norm_(net.parameters(), 3)
    optimizer.step()
    optimizer.zero_grad()
    ep_loss += loss

net.eval()
print('epoch: %d, loss: %0.2f, time: %0.2f sec, dev perplexity: %0.2f'↪
↪%\
        (epoch, ep_loss, time.time()-start_time,↪
↪compute_perplexity(wikitext['dev'], net)))

```

Now implement the following class, which defines a recurrent neural language model, by implementing the forward function.

```

[18]: class RNNLM(nn.Module):
    def __init__(self, params):
        super(RNNLM, self).__init__()
        self.vocab_size = params['vocab_size']
        self.d_emb = params['d_emb'] # size of word-embedding vector
        self.d_hid = params['d_hid'] # vector size of the hidden layer
        self.n_layer = 1
        self.batch_size = params['batch_size']

        self.encoder = nn.Embedding(self.vocab_size, self.d_emb)
        self.rnn = nn.RNN(self.d_emb, self.d_hid, self.n_layer,↪
↪batch_first=True)
        self.decoder = nn.Linear(self.d_hid, self.vocab_size)

    def forward(self, batch):
        """
        IMPLEMENT ME!
        Encode the data using the embedding layer you initialized.
        Pass the encoded data and hidden states to your RNN.

```

*Return unnormalized logits for each token's prediction.*

*Why just logits? Check the document of `torch.nn.CrossEntropyLoss`, since it combines `nn.LogSoftmax()` and `nn.NLLLoss()`, you don't need to explicitly use the softmax function!*

```
"""
batch_size, seq_len = batch.shape
hidden = (torch.zeros(self.n_layer, batch_size, self.d_hid).to(device))
x = self.encoder(batch) # (batch_size, seq_len, d_emb)
x, hidden = self.rnn(x, hidden)
x = x.reshape(-1, self.d_hid) # (batch_size*seq_len, d_hid)
x = self.decoder(x) # (batch_size*seq_len, vocab_size)
logits = x.reshape(batch_size, seq_len, self.vocab_size) # (batch_size,
↪seq_len, vocab_size)
return logits
```

Run the following cell to test that your implementation is at least returning tensors of the proper dimensionality. Note that this is just a sanity check. Your RNNLM might still be implemented incorrectly even if it passes. You will have to obtain a reasonable perplexity after training on WikiText to be certain that you've done it right.

```
[19]: def test_RNNLM():
    test_batch = torch.LongTensor(5, 4).random_(0, 10).to(device)
    params = {}
    params['vocab_size'] = len(idx_to_word)
    params['d_emb'] = 8
    params['d_hid'] = 8
    params['batch_size'] = 5
    testnet = RNNLM(params)
    testnet.to(device)
    test_output = testnet(test_batch)
    assert test_output.shape[0] == params['batch_size'], "size of dimension 0_
↪is incorrect, expect %i but got %i" % \
                                                                    5
    ↪(params['batch_size'], test_output.shape[0])
    assert test_output.shape[1] == test_batch.shape[1], "size of dimension 1 is_
↪incorrect, expect %i but got %i" % \
                                                                    (test_batch.shape[1], 4)
    ↪test_output.shape[1])
    assert test_output.shape[2] == params['vocab_size'], "size of dimension 2_
↪is incorrect, expect %i but got %i" % \
                                                                    5
    ↪(params['vocab_size'], test_output.shape[2])
    print("Congratulations, you passed the RNNLM test!")
test_RNNLM()
```

Congratulations, you passed the RNNLM test!

Once you pass the above test, train your RNNLM model on WikiText by running the cell below. It should take a couple minutes per epoch.

```
[20]: # DO NOT CHANGE THESE HYPERPARAMETERS, WE WILL CHECK!
```

```
params = {}
params['vocab_size'] = len(idx_to_word)
params['d_emb'] = 512
params['d_hid'] = 256
params['batch_size'] = 64
params['epochs'] = 5
params['learning_rate'] = 0.001

RNNnet = RNNLM(params)
RNNnet.to(device)
train_lm(wikitext['train'], params, RNNnet)
```

```
epoch: 0, loss: 6398.67, time: 29.49 sec, dev perplexity: 182.50
epoch: 1, loss: 5661.17, time: 29.16 sec, dev perplexity: 155.22
epoch: 2, loss: 5319.08, time: 28.99 sec, dev perplexity: 148.12
epoch: 3, loss: 5062.62, time: 29.03 sec, dev perplexity: 143.56
epoch: 4, loss: 4854.10, time: 29.03 sec, dev perplexity: 146.26
```

After training is finished, run the cell below to get the perplexity on the test set. If you did it right, your perplexity should be around 135-140.

```
[21]: RNNnet.eval() # we're no longer training the network
print('%s perplexity: %0.2f' % ('test', compute_perplexity(wikitext['test'],
↪RNNnet)))
```

```
test perplexity: 137.01
```

### 3.3 Q2.3 Neural Language Model with attention (30 pts)

Only start working at this after you've correctly implemented the RNNLM in the previous problem, as you'll want to copy over some code here. Complete the forward function of both the ATTNLM and Attention modules by following the instructions in the comment block. **Each epoch may take 3-5 minutes to run, so start early!**

```
[22]: # An RNN language model with attention, you implement this!
```

```
class ATTNLM(nn.Module):
    def __init__(self, params):
        super(ATTNLM, self).__init__()

        self.vocab_size = params['vocab_size']
        self.d_emb = params['d_emb']
        self.d_hid = params['d_hid']
        self.n_layer = 1
        self.btz = params['batch_size']
```

```

        self.encoder = nn.Embedding(self.vocab_size, self.d_emb)
        self.attn = Attention(self.d_hid)
        self.rnn = nn.RNN(self.d_emb, self.d_hid, self.n_layer,
↪batch_first=True)
        # the combined_W maps the combined hidden states and context vectors to
↪d_hid
        self.combined_W = nn.Linear(self.d_hid * 2, self.d_hid)
        self.decoder = nn.Linear(self.d_hid, self.vocab_size)

    def forward(self, batch, return_attn_weights=False):
        batch_size, seq_len = batch.shape
        hidden = torch.zeros(self.n_layer, batch_size, self.d_hid).to(device)
        encoded = self.encoder(batch)

        rnn_out, hidden = self.rnn(encoded, hidden)

        context_vectors, attn_weights = self.attn(rnn_out)

        combined = torch.cat((rnn_out, context_vectors), dim=2)
        combined = self.combined_W(combined)
        logits = self.decoder(combined.view(-1, self.d_hid))
        logits = logits.view(batch_size, seq_len, -1)

        if return_attn_weights:
            return attn_weights

        return logits

class Attention(nn.Module):
    def __init__(self, d_hidden):
        super(Attention, self).__init__()
        self.linear_w1 = nn.Linear(d_hidden, d_hidden)
        self.linear_w2 = nn.Linear(d_hidden, 1)

    def forward(self, x):
        batch_size, seq_len, _ = x.shape
        attn_weights = torch.zeros(batch_size, seq_len, seq_len).to(device)

        for t in range(seq_len):
            attn_input = x[:, t, :].unsqueeze(1)
            h_t = torch.tanh(self.linear_w1(x[:, :t+1, :]))
            attn_scores = self.linear_w2(h_t).squeeze(2)
            attn_dist = nn.functional.softmax(attn_scores, dim=1)
            attn_weights[:, t, :t+1] = attn_dist

```

```

        context_vector = torch.bmm(attn_dist.unsqueeze(1), x[:, :t+1, :]).
↪squeeze(1)
        if t == 0:
            context_vectors = context_vector.unsqueeze(1)
        else:
            context_vectors = torch.cat((context_vectors, context_vector.
↪unsqueeze(1)), dim=1)

        return context_vectors, attn_weights.tril() # lower triangular

```

Run the following cell to sanity check your implementation; do not continue until you pass all of the tests!

```

[23]: def test_ATT_NLM():
    test_batch = torch.LongTensor(5, 4).random_(0, 10).to(device)
    params = {}
    params['vocab_size'] = len(idx_to_word)
    params['d_emb'] = 8
    params['d_hid'] = 8
    params['batch_size'] = 5
    testnet = ATT_NLM(params)
    testnet.to(device)
    test_output = testnet(test_batch)
    assert test_output.shape[0] == params['batch_size'], "size of dimension 0_
↪is incorrect, expect %i but got %i" % \
                                                                    (
↪(params['batch_size'], test_output.shape[0])
        assert test_output.shape[1] == test_batch.shape[1], "size of dimension 1 is_
↪incorrect, expect %i but got %i" % \
                                                                    (test_batch.shape[1],
↪test_output.shape[1])
        assert test_output.shape[2] == params['vocab_size'], "size of dimension 2_
↪is incorrect, expect %i but got %i" % \
                                                                    (
↪(params['vocab_size'], test_output.shape[2])
            testnet = ATT_NLM(params)
            testnet.to(device)
            test_output = testnet(test_batch, return_attn_weights=True)
            assert test_output.shape[0] == params['batch_size'], "size of dimension 0_
↪is incorrect, expect %i but got %i" % \
                                                                    (
↪(params['batch_size'], test_output.shape[0])
                assert test_output.shape[1] == test_batch.shape[1], "size of dimension 1 is_
↪incorrect, expect %i but got %i" % \
                                                                    (test_batch.shape[1],
↪test_output.shape[1])

```



```

    assert test_output.shape[2] == test_batch.shape[1], "size of dimension 2 is
↳incorrect, expect %i but got %i" % \
                                                    (test_batch.shape[1],
↳test_output.shape[2])
    prob_dist = torch.sum(test_output, dim=2)[: , 1:]
    assert all([x > 0.99 and x < 1.01 for x in prob_dist.reshape(-1)]),
↳"attention weights not properly normalized, got {}".format(prob_dist)
    print("Congratulations, you passed the ATTNLM test!")

test_ATTNLM()

```

Congratulations, you passed the ATTNLM test!

Now, train your ATTNLM model on WikiText by running the following code cell. If the perplexity on dev set is nan or inf, it is likely the model is corrupted due to gradient exploding/vanishing or other numerical instability issue; stop this cell and run it again.

[24]: *# DO NOT CHANGE THESE HYPERPARAMETERS, WE WILL CHECK!*

```

params = {}
params['vocab_size'] = len(idx_to_word)
params['d_emb'] = 512
params['d_hid'] = 256
params['n_layer'] = 1
params['batch_size'] = 64
params['epochs'] = 6
params['learning_rate'] = 0.0005

ATTNnet = ATTNLM(params)
ATTNnet.cuda()
train_lm(wikitext['train'], params, ATTNnet)

```

```

epoch: 0, loss: 6458.55, time: 67.59 sec, dev perplexity: 195.64
epoch: 1, loss: 5833.66, time: 67.23 sec, dev perplexity: 169.06
epoch: 2, loss: 5552.28, time: 66.82 sec, dev perplexity: 156.35
epoch: 3, loss: 5330.04, time: 67.21 sec, dev perplexity: 148.09
epoch: 4, loss: 5141.65, time: 66.84 sec, dev perplexity: 147.16
epoch: 5, loss: 4978.49, time: 67.66 sec, dev perplexity: 144.76

```

Finally, compute the perplexity on the test set. If you implemented it correctly, you should get a perplexity of around 145-150. Due to random effects, it is possible to get perplexity slightly lower than 145. Make sure you didn't add any additional nonlinearity operation which can lead to lower perplexity.

[25]: *ATTNnet.eval() # we're no longer training the network*

```

print('%s perplexity: %0.2f' % ('test', compute_perplexity(wikitext['test'],
↳ATTNnet)))

```

test perplexity: 135.30

### 3.4 Q2.4 Generate text from the neural LMs (5 pts)

Run the following code cell to generate some text from your RNNLM and ATTNLM.

```
[26]: def sample_from_lm(net, context, max_words=50):

    with torch.no_grad():
        for i in range(max_words):
            data = torch.LongTensor([context]).to(device)
            decoded = net(data)
            decoded = decoded[0, -1].exp().cpu()
            w_i = torch.multinomial(decoded, 1)[0].item()
            if w_i in [1, 2, 3]:
                continue
            context.append(w_i)

        return context

word_to_idx = dict((v,k) for (k,v) in idx_to_word.items())
context = [word_to_idx[w] for w in 'he is the '.split()]

rnn_completion = sample_from_lm(RNNnet, context)
print('rnn completion: ', ' '.join([idx_to_word[w] for w in rnn_completion]))
```

rnn completion: he is the virginia image of 24 .

```
[27]: word_to_idx = dict((v,k) for (k,v) in idx_to_word.items())
context = [word_to_idx[w] for w in 'he is the '.split()]

rnn_completion = sample_from_lm(ATTNnet, context)
print('attention rnn completion: ', ' '.join([idx_to_word[w] for w in
↪rnn_completion]))
```

attention rnn completion: he is the inconsistencies telecommunications koopa  
ruwan bone marrow transplant ahk ibarra .

Do you notice any differences in coherence or grammaticality compared to the n-gram models? What about any differences between the RNNLM and the ATTNLM? If you observed any distinct differences, explain why you think they exist; if not, explain why all of the outputs appear to be of similar quality.

#### 3.4.1 Answer in two to four sentences here.

As compared to the n-gram models, the neural language models (RNNLM and ATTNLM) have demonstrated better coherence and grammaticality. In addition, ATTNLM has outperformed RNNLM in terms of generating more coherent and grammatically correct sentences, as it utilizes attention to focus on different parts of the input sequence. This makes the model more effective in modeling long-range dependencies and capturing context-specific information.

### 3.5 Q2.5 Interpreting attention (5 pts)

Finally, let's visualize some attention heatmaps by running the following two code cells.

```
[28]: def plot_attn_heatmap(sent):

    sent_in_id = [word_to_idx[w] for w in sent.split()]

    with torch.no_grad():
        data = torch.LongTensor([sent_in_id]).to(device)
        weights = ATTNnet(data, return_attn_weights=True)

    fig, ax = plt.subplots()

    sent_sp = sent.split()
    ax.set_xticks(np.arange(len(sent_sp)))
    ax.set_yticks(np.arange(len(sent_sp)))
    ax.set_xticklabels(sent_sp)
    ax.set_yticklabels(sent_sp)
    plt.setp(ax.get_xticklabels(), rotation=45, ha='right',
    ↪rotation_mode="anchor")

    plt.imshow(weights[0, :].cpu())

sent = "top warning signs earth is warming , according to experts"
plot_attn_heatmap(sent)
```

Canceled future for execute\_request message before replies were done

The Kernel crashed while executing code in the the current cell or a previous  
↪cell. Please review the code in the cell(s) to identify a possible cause of  
↪the failure. Click ↪for more info. View Jupyter ↪further details.

```
[ ]: sent = "us cities lose 36 million trees each year . here is why it matters "
plot_attn_heatmap(sent)
```

```
-----
NameError                                Traceback (most recent call last)
~\AppData\Local\Temp\ipykernel_5884\1065687754.py in <module>
      1 sent = "us cities lose 36 million trees each year . here is why it
    ↪matters "
----> 2 plot_attn_heatmap(sent)

~\AppData\Local\Temp\ipykernel_5884\1895593329.py in plot_attn_heatmap(sent)
      1 def plot_attn_heatmap(sent):
```

```

2
----> 3     sent_in_id = [word_to_idx[w] for w in sent.split()]
4
5     with torch.no_grad():

~\AppData\Local\Temp\ipykernel_5884\1895593329.py in <listcomp>(.0)
1 def plot_attn_heatmap(sent):
2
----> 3     sent_in_id = [word_to_idx[w] for w in sent.split()]
4
5     with torch.no_grad():

NameError: name 'word_to_idx' is not defined

```

Each row of these plots represents the attention weights on the history tokens when the model is trying to predict the next word. For example, the third row of the first plot can be interpreted as the attention weights over “top” and “warning” when predicting “signs”; you’ll note that the rest of the row is black (i.e., zero attention on future words). Are these attention maps interpretable? If you (as a human) were solving the same word prediction problem, would you focus on the same words as the ATTNLM does?

### 3.5.1 *Answer in two to four sentences here.*

The attention maps can be somewhat interpretable as they show which previous words the model is focusing on while predicting the next word. However, the maps can be complex and difficult to interpret, especially as the sequences become longer. As a human, I might not focus on the same words as the ATTNLM does since the model might be using certain patterns or associations in the training data that are not obvious or intuitive to humans.