

Assignment 2: Word Prediction

Deadline: Sunday, April 18th, by 9pm.

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

In this assignment, we will make a neural network that can predict the next word in a sentence given the previous three.

In doing this prediction task, our neural networks will learn about words and about how to represent words.

We'll explore the **vector representations** of words that our model produces, and analyze these representations.

You may modify the starter code as you see fit, including changing the signatures of functions and adding/removing helper functions. However, please make sure that you properly explain what you are doing and why.

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import collections

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

from collections import Counter

In [2]: # helper class to print nice outputs
class txt:
    BOLD = '\033[1m'
    UNDERLINE = '\033[4m'
    END = '\033[0m'

    def make_bold(text):
        return txt.BOLD+text + txt.END

    def make_underline(text):
        return txt.UNDERLINE+text + txt.END
```

Question 1. Data (15%)

With any machine learning problem, the first thing that we would want to do is to get an intuitive understanding of what our data looks like. Download the file `raw_sentences.txt` from the course page on Moodle and upload it to Google Drive. Then, mount Google Drive from your Google Colab notebook:

```
In [3]: from google.colab import drive
drive.mount('/content/gdrive')
```

Mounted at /content/gdrive

Find the path to `raw_sentences.txt`:

```
In [4]: file_path = '/content/gdrive/My Drive/Intro_to_Deep_Learning/assignment2/raw_sentence
```

The following code reads the sentences in the file, split each sentence into its individual words, and stores the sentences (list of words) in the variable `sentences`.

```
In [5]: sentences = []
for i, sent in enumerate(open(file_path, 'r')):
    words = line.split()
    sentence = [word.lower() for word in words]
    sentences.append(sentence)
```

There are 97,162 sentences in total, and these sentences are composed of 250 distinct words.

```
In [6]: vocab = set([w for s in sentences for w in s])
print(len(vocab))
97162
250
```

We'll separate our data into training, validation, and test. We'll use 10,000 sentences for test, 10,000 for validation, and the rest for training.

```
In [7]: test, valid, train = sentences[:10000], sentences[10000:20000], sentences[20000:]
```

Part (a) -- 2%

Display 10 sentences in the training set. Explain how punctuations are treated in our word representation, and how words with apostrophes are represented.

```
In [8]: for i, sent in enumerate(train[:1020]):
    print(str(i+1)+":", " ".join(sent))
```

```
1. but,for no, now, i, this is it.
2. she's still there for us.
3. it's part of this game, man.
4. i, how, how do we get, treat
5. but they do nt last too long
6. more are like me , she said .
7. who do you think they want to be like ?
8. no , he could not
9. so i left it up to them .
10. we were nt right .
```

We can tell that punctuations such as dot or a comma is treated as a word by itself. When looking at sentence number 2 for example we observe that words with apostrophes are represented as 2 different words: the prefix of the ord and the apostrophes.

Part (b) -- 3%

Print the 10 most common words in the vocabulary and how often does each of these words appear in the train sentences. Express the second quantity as a percentage (i.e. number of occurrences of the word / total number of words in the training set).

These are useful quantities to compute, because one of the first things a machine learning model learn is to predict the most common class. Getting a sense of the distribution of our data will help you understand our model's behaviour.

You can use Python's `collections.Counter` class if you would like to.

```
In [9]: train_words = []
for sentence in train:
    for word in sentence:
        train_words.append(word)

print(make_bold("10 most common words in the vocabulary:"))
c = Counter(train_words)
freq = most_common(10)
def dataFrame(freq, col_names=['word', 'num of occ'], *args):
    df['percentage'] = df['num of occ']/len(train_words)*100
    df
```

```
Out[9]: word num of occ percentage
0 it 23118 3.845648
1 i 17684 2.941710
2 do 16181 2.691688
3 to 15490 2.576741
4 nt 13009 2.164030
5 ? 12881 2.142737
6 the 12583 2.093165
7 's 12552 2.088008
```

```
In [10]: ax = df.plot.bar(x='word', y='percentage', rot=0)
ax.set_title('10 most common words percentage')
```

```
Out[10]: Text(0.5, 1.0, '10 most common words percentage')
```



```
In [11]: ax = df.plot.bar(x='word', y='num of occ', rot=0)
ax.set_title('10 most common words count')
```

```
Out[11]: Text(0.5, 1.0, '10 most common words count')
```



That distribution is reasonable because each sentence is finished with dot or question mark for example.

Let's check words distribution without punctuation

```
In [12]: import string

# A list of all the words in the vocabulary without punctuation:
c = Counter([w for w in train_words if w not in string.punctuation])
freq = most_common(10)
def dataFrame(freq, col_names=['word', 'num of occ'], *args):
    df['percentage'] = df['num of occ']/len(train_words)*100
    df
```

```
Out[12]: word num of occ percentage
0 it 23118 3.845648
1 i 17684 2.941710
2 do 16181 2.691688
3 to 15490 2.576741
4 nt 13009 2.164030
5 the 12583 2.093165
6 's 12552 2.088008
7 he 12192 2.082123
8 he 12192 2.082123
9 y 10166 1.691011
```

```
In [13]: ax = df.plot.bar(x='word', y='percentage', rot=0)
ax.set_title('10 most common words percentage')
```

```
Out[13]: Text(0.5, 1.0, '10 most common words percentage')
```



```
In [14]: ax = df.plot.bar(x='word', y='num of occ', rot=0)
ax.set_title('10 most common words count')
```

```
Out[14]: Text(0.5, 1.0, '10 most common words count')
```



Part (c) -- 10%

Our neural network will take as input three words and predict the next one. Therefore, we need our data set to be comprised of sequences of four consecutive words in a sentence, referred to as *4grams*.

Complete the helper functions `convert_words_to_indices` and `generate_4grams`, so that the function `process_data` will take a list of sentences (i.e. list of list of words), and generate an $N \times 4$ numpy matrix containing indices of 4 words that appear next to each other, where N is the number of 4grams (sequences of 4 words appearing one after the other) that can be found in the complete list of sentences. Examples of how these functions should operate are detailed in the code below.

You can use the defined `vocab`, `vocab_itos`, and `vocab_itos` in your code.

```
In [15]: # A list of all the words in the data set. We will assign a unique
# identifier for each of these words.
vocab = sorted(list(set([w for s in train for w in s])))
# A mapping of index -> word (this vocab)
vocab_itos = dict(enumerate(vocab))
# A mapping of word -> its index
vocab_stoi = {word:index for index, word in vocab_itos.items()}
```

```
def convert_words_to_indices(sents):
    """
    This function takes a list of sentences (list of list of words)
    and returns a new list with the same structure, but where each word
    is replaced by its index in 'vocab_stoi'.
```

```
>>> convert_words_to_indices([[1, 'one', 'in', 'five', 'are', 'over', 'here'], ['other', 'one', 'in', 'five', 'are', 'over', 'here'], ['other', 'one', 'in', 'five', 'are', 'over', 'here']])
[[148, 98, 70, 23, 154, 89], [151, 148, 181, 246], [248]]
```

```
>>> generate_4grams([[[148, 98, 70, 23, 154, 89], [151, 148, 181, 246], [248]]])
[[148, 98, 70, 23], [98, 70, 23, 154], [70, 23, 154, 89], [151, 148, 181, 246]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

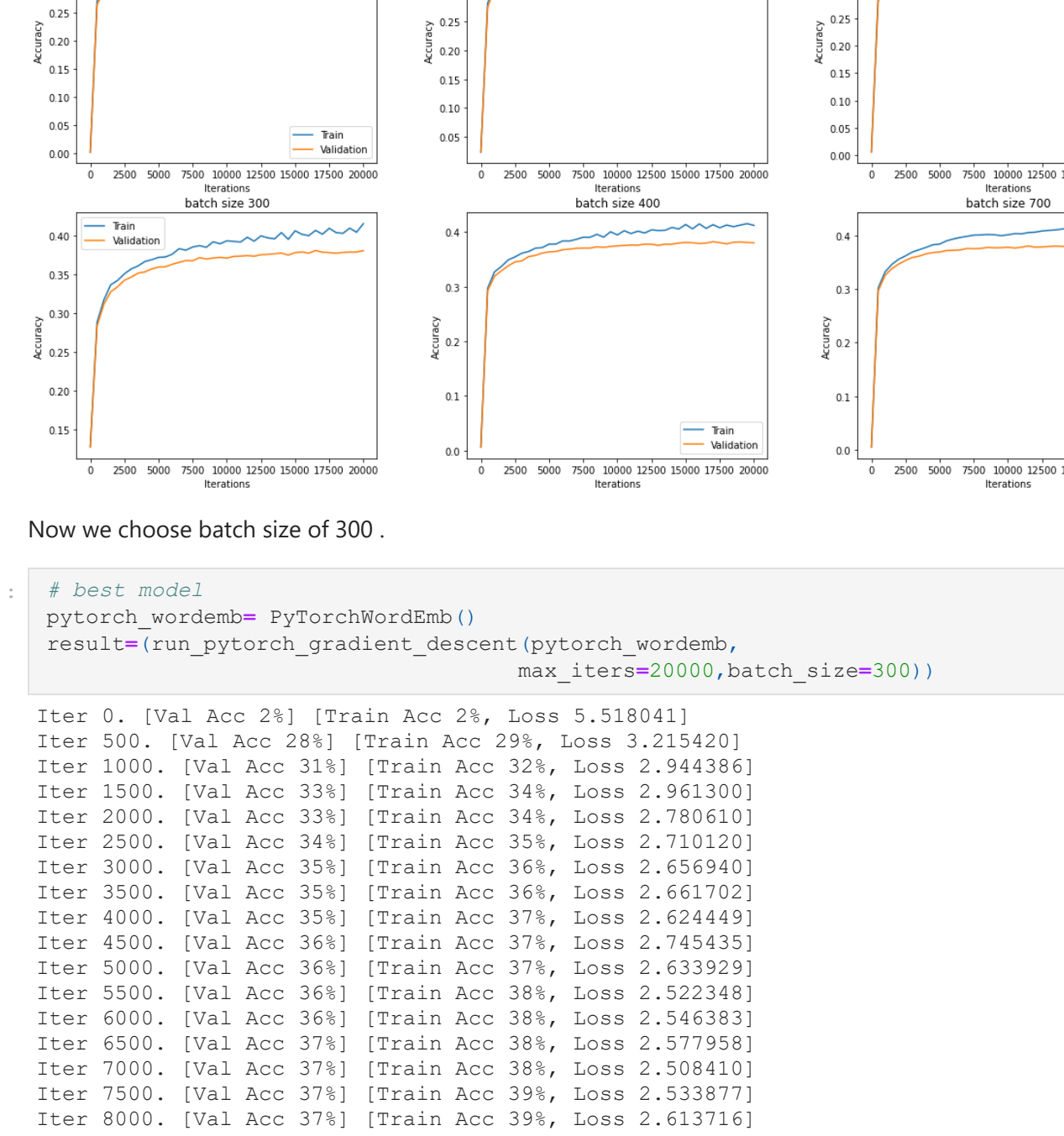
```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1]]
```

```
>>> generate_4grams([[[1, 1, 1, 1], [1, 1, 1, 1]]])
[[1, 1, 1, 1], [1, 1, 1, 1
```


Accuracy per iteration



Now we choose batch size of 300.

```
In [39]: # best model
pytorch_wordemb= PyTorchWordEmb()
result=run_pytorch_gradient_descent(pytorch_wordemb,
max_iters=2000,batch_size=300))

Iter 0. [Val Acc 2%] [Train Acc 2%, Loss 5.518041]
Iter 500. [Val Acc 28%] [Train Acc 29%, Loss 3.215420]
Iter 1000. [Val Acc 31%] [Train Acc 32%, Loss 2.944386]
Iter 1500. [Val Acc 33%] [Train Acc 34%, Loss 2.961300]
Iter 2000. [Val Acc 33%] [Train Acc 34%, Loss 2.780610]
Iter 2500. [Val Acc 34%] [Train Acc 35%, Loss 2.710120]
Iter 3000. [Val Acc 35%] [Train Acc 36%, Loss 2.656940]
Iter 3500. [Val Acc 35%] [Train Acc 36%, Loss 2.661702]
Iter 4000. [Val Acc 35%] [Train Acc 37%, Loss 2.624449]
Iter 4500. [Val Acc 36%] [Train Acc 37%, Loss 2.745435]
Iter 5000. [Val Acc 36%] [Train Acc 37%, Loss 2.63929]
Iter 5500. [Val Acc 36%] [Train Acc 38%, Loss 2.522348]
Iter 6000. [Val Acc 36%] [Train Acc 38%, Loss 2.546383]
Iter 6500. [Val Acc 37%] [Train Acc 38%, Loss 2.577858]
Iter 7000. [Val Acc 37%] [Train Acc 38%, Loss 2.508410]
Iter 7500. [Val Acc 37%] [Train Acc 39%, Loss 2.538777]
Iter 8000. [Val Acc 37%] [Train Acc 39%, Loss 2.613716]
Iter 8500. [Val Acc 37%] [Train Acc 39%, Loss 2.581933]
Iter 9000. [Val Acc 37%] [Train Acc 39%, Loss 2.459998]
Iter 9500. [Val Acc 37%] [Train Acc 39%, Loss 2.522073]
Iter 10000. [Val Acc 37%] [Train Acc 39%, Loss 2.320807]
Iter 10500. [Val Acc 37%] [Train Acc 39%, Loss 2.424054]
Iter 11000. [Val Acc 38%] [Train Acc 39%, Loss 2.315505]
Iter 11500. [Val Acc 37%] [Train Acc 40%, Loss 2.657997]
Iter 12000. [Val Acc 38%] [Train Acc 39%, Loss 2.254560]
Iter 12500. [Val Acc 38%] [Train Acc 40%, Loss 2.381103]
Iter 13000. [Val Acc 38%] [Train Acc 40%, Loss 2.432642]
Iter 13500. [Val Acc 38%] [Train Acc 40%, Loss 2.206470]
Iter 14000. [Val Acc 38%] [Train Acc 40%, Loss 2.385060]
Iter 14500. [Val Acc 38%] [Train Acc 40%, Loss 2.265549]
Iter 15000. [Val Acc 38%] [Train Acc 41%, Loss 2.649423]
Iter 15500. [Val Acc 38%] [Train Acc 40%, Loss 2.392514]
Iter 16000. [Val Acc 38%] [Train Acc 40%, Loss 2.252542]
Iter 16500. [Val Acc 38%] [Train Acc 41%, Loss 2.212328]
Iter 17000. [Val Acc 38%] [Train Acc 40%, Loss 2.485659]
Iter 17500. [Val Acc 38%] [Train Acc 41%, Loss 2.416045]
Iter 18000. [Val Acc 38%] [Train Acc 40%, Loss 2.225331]
Iter 18500. [Val Acc 38%] [Train Acc 40%, Loss 2.386328]
Iter 19000. [Val Acc 38%] [Train Acc 41%, Loss 2.456148]
Iter 19500. [Val Acc 38%] [Train Acc 40%, Loss 2.434269]
Iter 20000. [Val Acc 38%] [Train Acc 42%, Loss 2.480032]
```

Part (b) -- 8%

Use the function `make_prediction` that you wrote earlier to predict what the next word should be in each of the following sentences:

- "You are a"
- "few companies show"
- "There are no"
- "yesterday i was"
- "the game had"
- "yesterday the federal"

How do these predictions compared to the previous model?

Print the output for all of these sentences using the new network and **Write** below how the new results compare to the previous ones.

Just like before, if you encounter overfitting, train your model for more iterations, or change the hyperparameters in your model. You may need to do this even if your training accuracy is >=38%.

```
In [39]: print("you are a",make_bold(make_underline(make_prediction_torch(pytorch_wordemb, ['y
print("few companies shows" , make_bold(make_underline(make_prediction_torch(pytorch_
print("there are no",make_bold(make_underline(make_prediction_torch(pytorch_wordemb,
print("yesterday i was",make_bold(make_underline(make_prediction_torch(pytorch_wordemb,
print("the game had",make_bold(make_underline(make_prediction_torch(pytorch_wordemb,
print("yesterday the federal",make_bold(make_underline(make_prediction_torch(pytorch_

you are a good
few companies shows a
there are no other
yesterday i was at
the game had to
yesterday the federal government
```

We got the same result as in the preview model!

Part (c) -- 4%

Report the test accuracy of your model

```
In [40]: estimate_accuracy_torch(pytorch_wordemb, test4grams)

Out[40]: 0.3844901367289427
```

Question 4. Visualizing Word Embeddings (15%)

While training the `PyTorchMLP`, we trained the `word_emb_layer`, which takes a one-hot representation of a word in our vocabulary, and returns a low-dimensional vector representation of that word. In this question, we will explore these word embeddings, which are a key concept in natural language processing.

Part (a) -- 5%

The code below extracts the **weights** of the word embedding layer, and converts the PyTorch tensor into a numpy array. Explain why each row of `word_emb` contains the vector representing of a word. For example `word_emb[vocab_stoi["any"],:]` contains the vector representation of the word "any".

```
In [45]: word_emb_weights = list(pytorch_wordemb.word_emb_layer.parameters())[0]
word_emb = word_emb_weights.detach().numpy().T
```

Explanation: each row of `word_emb` is the result of multiply and sums all weights with the `i`'th input for the layer, which is a one hot representation of the `i`'th word. But since the one-hot representation will gives us vectors of zeros and only 1 in the `i`'th index we get that each row is the vector represening of the word itself.

For example, `word_emb[vocab_stoi["any"],:]` contains the vector representation of the word "any".

Part (b) -- 5%

One interesting thing about these word embeddings is that distances in these vector representations of words make some sense! To show this, we have provided code below that computes the *cosine similarity* of every pair of words in our vocabulary. This measure of similarity between vector `v` and `w` is defined as

$$d_{cos}(v, w) = \frac{v \cdot w}{||v|| ||w||}$$

We also pre-scale the vectors to have a unit norm, using Numpy's `norm` method.

```
In [46]: norms = np.linalg.norm(word_emb, axis=1)
word_emb_norm = (word_emb.T / norms).T
similarities = np.matmul(word_emb_norm, word_emb_norm.T)
```

Some example distances. The first one should be larger than the second
print(make_bold("similarity between 'any' and 'many'"), similarities[vocab_stoi['any'],
print(make_bold("similarity between 'any' and 'government'"), similarities[vocab_stoi['any'],

similarity between 'any' and 'many': 0.18362503
similarity between 'any' and 'government': 0.15161606

Compute the 5 closest words to the following words:

- "four"
- "go"
- "what"
- "should"
- "school"
- "your"
- "not"

```
In [46]: for word in ["four", "go", "what", "should", "school", "your", "yesterday", "not"]:
```

words_similarities[(w,similarities[vocab_stoi[word], vocab_stoi[w]])] **for** w in vocab
similarities dict = {}
[similarities dict.update((k,v)) **for** k,v in words_similarities]
Counter(similarities dict)
five_most_common=most_common(5)
df=pd.DataFrame(five_most_common,columns=['word', 'similarity'])
ax=plt.plot.bar(x=word, y=similarity, rot=0)
ax.set_title("%s closest words for \"%word\"" % word)
print(make_bold("%s closest words for \"%word\"" % word), five_most_common)
print("=====")

5 closest words for "four": [('three', 0.6514785), ('five', 0.5263721), ('many', 0.37343597), ('two', 0.34956172), ('several', 0.30247074)]

5 closest words for "go": [('come', 0.45551512), ('going', 0.4357491), ('back', 0.36501148), ('down', 0.3545441), ('same', 0.31410465)]

5 closest words for "what": [('where', 0.4293098), ('dr.', 0.3823012), ('who', 0.3806837), ('how', 0.3584918), ('part', 0.34056073)]

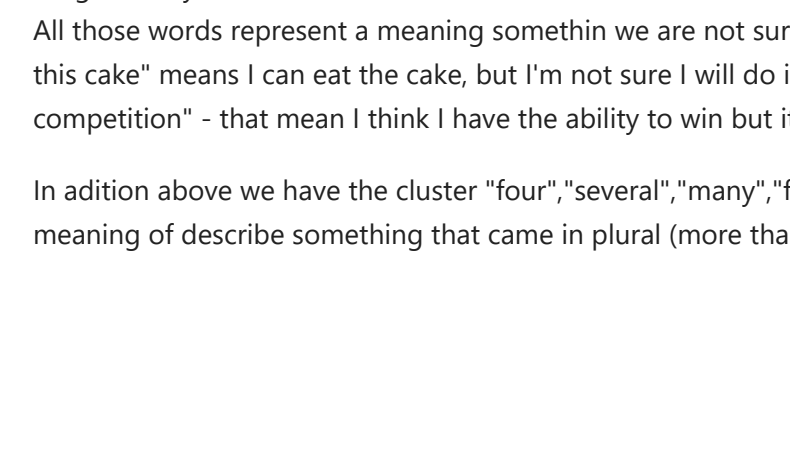
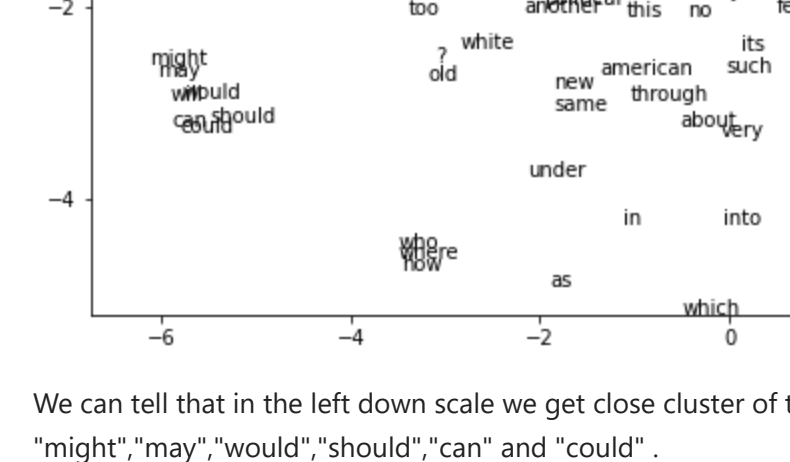
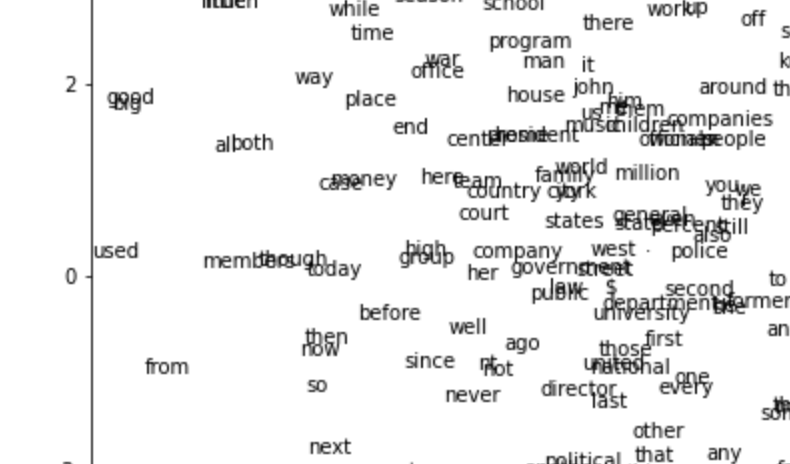
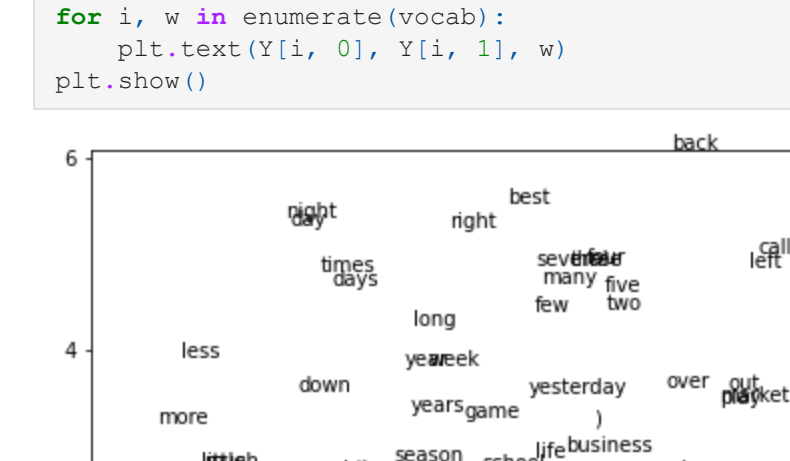
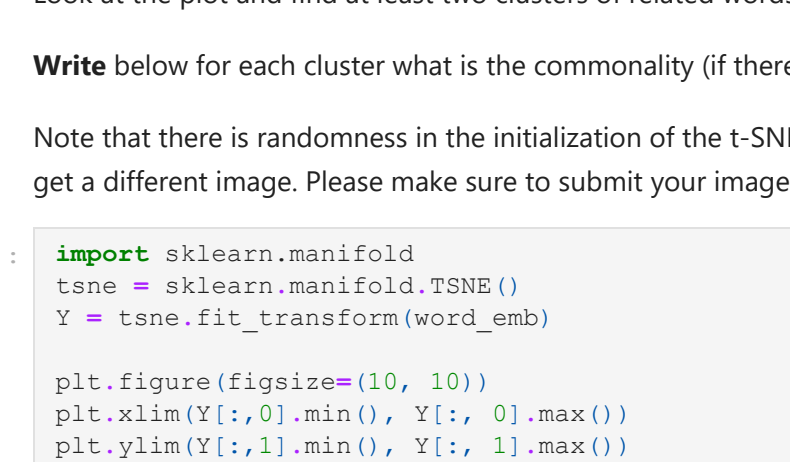
5 closest words for "should": [('could', 0.4523806), ('would', 0.4472415), ('may', 0.4060487), ('can', 0.3821438), ('might', 0.3427131)]

5 closest words for "school": [('program', 0.4482231), ('week', 0.47374997), ('year', 0.4049222), ('other', 0.2969315), ('right', 0.28540703)]

5 closest words for "your": [('my', 0.5222146), ('each', 0.49735525), ('federal', 0.3004222), ('other', 0.2969315), ('right', 0.28540703)]

5 closest words for "yesterday": [('today', 0.4687232), ('I', 0.4189693), ('ago', 0.4180559), ('said', 0.38785744), ('well', 0.3846288)]

5 closest words for "not": [('nt', 0.47753257), ('never', 0.3488891), ('might', 0.34298548), ('also', 0.2864333), ('around', 0.27710342)]



With word embedding words with close meaning will have vectors with close distance between them .

Here for example the 5 closest word to "four" are "three","five","many","several" and two. All of this options can replace the word four in a senetence and the sentence will still be reasonable and with similar meaning.

Part (c) -- 5%

We can visualize the word embeddings by reducing the dimensionality of the word vectors to 2D. There are many dimensionality reduction techniques that we could use, and we will use an algorithm called t-SNE. (You don't need to know what this is for the assignment; we will cover it later in the course.) Nearby points in this 2-D space are meant to correspond to nearby points in the original, high-dimensional space.

The following code runs the t-SNE algorithm and plots the result.

Look at the plot and find at least two clusters of related words.

Write below for each cluster what is the commonality (if there is any) and if they make sense.

Note that there is randomness in the initialization of the t-SNE algorithm. If you re-run this code, you may get a different image. Please make sure to submit your image in the PDF file.

```
In [47]: import sklearn.manifold
t_sne = sklearn.manifold.TSNE()
Y = t_sne.fit_transform(wordemb)

plt.figure(figsize=(10, 10))
plt.xlim(Y[:,0].min(), Y[:,0].max())
plt.ylim(Y[:,1].min(), Y[:,1].max())
for i, w in enumerate(vocab):
    plt.text(Y[i, 0], Y[i, 1], w)
plt.show()
```



We can tell that in the left down scale we get close cluster of the words :
"might","may","would","should","can" and "could" .

All those words represent a meaning something we are not sure about in the present or future: "I might eat this cake" means I can eat the cake, but I'm not sure I will do it. Same as "I could have won the competition" - that mean I think I have the ability to win but it does not mean I won.

In addition above we have the cluster "four","several","many","few","five","two" - which all have the same meaning of describe something that came in plural (more than one object).