

▼ Assignment 2: Word Prediction

Deadline: Sunday, December 11th, by 8pm.

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.

```
import pandas
import numpy as np
import matplotlib.pyplot as plt
import collections

import torch
import torch.nn as nn
import torch.optim as optim
```

▼ Question 1. Data (18%)

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:

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

Mounted at /content/gdrive
```

Find the path to `raw_sentences.txt`:

```
file_path = '/content/gdrive/MyDrive/Colab Notebooks/raw_sentences.txt'
```

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

```
sentences = []
for line in open(file_path):
    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.

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

```
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.

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

▼ Part (a) -- 3%

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

```
for i in range(10,20):
    print(f"Sentence number {i-9}: \n {train[i]}")
```

```

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

```

Punctuations are treated as words in this word representation, since it gets its own place on the list. Also, words with apostrophes are treated as two different words. We can see that the words with apostrophes are cutted to two different word, after the punctuation mark (').

▼ Part (b) – 4%

Print the 10 most common words in the vocabulary and how often does each of these words appear in the training 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 will 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.

```

from collections import Counter
words_list = [item for sublist in sentences for item in sublist]

most_common_words= [word for word, word_count in Counter(words_list).most_common(10)]
print(most_common_words)

train_words_list = [item for sublist in train for item in sublist]
count_train = Counter(train_words_list)
for word in most_common_words:
    print(f'The word: "{word}" appears {count_train[word]} times in the train data.')
    print(f'That is {(count_train[word]/len(train_words_list)*100):.2f}% of the train words. \n')

['.', 'it', ',', 'i', 'do', 'to', 'nt', '?', 'the', 'that']
The word: "." appears 64297 times in the train data.
That is 10.70% of the train words.

The word: "it" appears 23118 times in the train data.
That is 3.85% of the train words.

The word: ",", appears 19537 times in the train data.
That is 3.25% of the train words.

The word: "i" appears 17684 times in the train data.
That is 2.94% of the train words.

The word: "do" appears 16181 times in the train data.
That is 2.69% of the train words.

The word: "to" appears 15490 times in the train data.
That is 2.58% of the train words.

The word: "nt" appears 13009 times in the train data.
That is 2.16% of the train words.

The word: "?" appears 12881 times in the train data.
That is 2.14% of the train words.

The word: "the" appears 12583 times in the train data.
That is 2.09% of the train words.

The word: "that" appears 12535 times in the train data.
That is 2.09% of the train words.

```

▼ Part (c) – 11%

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_stoi` in your code.

```
from parsing.core import WordEnd
# 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 (string)
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`.

    Example:
    >>> convert_words_to_indices([[ 'one', 'in', 'five', 'are', 'over', 'here'], [ 'other', 'one', 'since', 'yesterday'], [ 'you']])
    [[148, 98, 70, 23, 154, 89], [151, 148, 181, 246], [248]]
    """

    # Write your code here
    indices_list = sents.copy()
    for i, sen in enumerate(sents):
        for j, word in enumerate(sen):
            indices_list[i][j] = vocab_stoi[word]
    return indices_list

def generate_4grams(seqs):
    """
    This function takes a list of sentences (list of lists) and returns
    a new list containing the 4-grams (four consequentially occurring words)
    that appear in the sentences. Note that a unique 4-gram can appear multiple
    times, one per each time that the 4-gram appears in the data parameter `seqs`.

    Example:

    >>> 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]]
    """

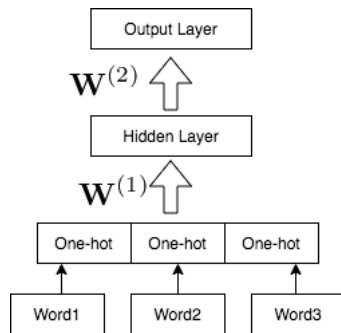
    # Write your code here
    _4gram_list = []
    for seq in seqs:
        if len(seq) < 4:
            continue
        else:
            pointer = 4
            while pointer <= len(seq):
                _4gram_list += [seq[pointer-4:pointer]]
                pointer += 1
    return _4gram_list

def process_data(sents):
    """
    This function takes a list of sentences (list of lists), and generates an
    numpy matrix with shape [N, 4] containing indices of words in 4-grams.
    """
    indices = convert_words_to_indices(sents)
    fourgrams = generate_4grams(indices)
    return np.array(fourgrams)

# We can now generate our data which will be used to train and test the network
train4grams = process_data(train)
valid4grams = process_data(valid)
test4grams = process_data(test)
```

▼ Question 2. A Multi-Layer Perceptron (44%)

In this section, we will build a two-layer multi-layer perceptron. Our model will look like this:



Since the sentences in the data are comprised of 250 distinct words, our task boils down to classification where the label space \mathcal{S} is of cardinality $|\mathcal{S}| = 250$ while our input, which is comprised of a combination of three words, is treated as a vector of size 750×1 (i.e., the concatenation of three one-hot 250×1 vectors).

The following function `get_batch` will take as input the whole dataset and output a single batch for the training. The output size of the batch is explained below.

Implement yourself a function `make_onehot` which takes the data in index notation and output it in a onehot notation.

Start by reviewing the helper function, which is given to you:

```
def make_onehot(data):
    """
    Convert one batch of data in the index notation into its corresponding onehot
    notation. Remember, the function should work for both xt and st.

    input - vector with shape D (1D or 2D)
    output - vector with shape (D,250)
    """

    # Write your code here
    base = np.eye(250) # create an 250*250 identity matrix. we will generate onehot from here
    out = np.stack([base[word] for word in data])
    return out

def get_batch(data, range_min, range_max, onehot=True):
    """
    Convert one batch of data in the form of 4-grams into input and output
    data and return the training data (xt, st) where:
    - `xt` is a numpy array of one-hot vectors of shape [batch_size, 3, 250]
    - `st` is either
        - a numpy array of shape [batch_size, 250] if onehot is True,
        - a numpy array of shape [batch_size] containing indices otherwise

    Preconditions:
    - `data` is a numpy array of shape [N, 4] produced by a call
      to `process_data`
    - range_max > range_min
    """
    xt = data[range_min:range_max, :3]
    xt = make_onehot(xt)
    st = data[range_min:range_max, 3]
    if onehot:
        st = make_onehot(st).reshape(-1, 250)
    return xt, st
```

▼ Part (a) -- 8%

We build the model in PyTorch. Since PyTorch uses automatic differentiation, we only need to write the *forward pass* of our model.

Complete the `forward` function below:

```
class PyTorchMLP(nn.Module):
    def __init__(self, num_hidden=400):
        super(PyTorchMLP, self).__init__()
        self.layer1 = nn.Linear(750, num_hidden)
        self.layer2 = nn.Linear(num_hidden, 250)
        self.num_hidden = num_hidden
    def forward(self, inp):
```

```

inp = inp.reshape([-1, 750])
# TODO: complete this function
# Note that we will be using the nn.CrossEntropyLoss(), which computes the softmax operation internally, as loss criterion

first_layer = self.layer1(inp)
second_layer = self.layer2(first_layer)

return second_layer

```

▼ Part (b) -- 10%

We next train the PyTorch model using the Adam optimizer and the cross entropy loss.

Complete the function `run_pytorch_gradient_descent`, and use it to train your PyTorch MLP model.

Obtain a training accuracy of at least 35% while changing only the hyperparameters of the train function.

Plot the learning curve using the `plot_learning_curve` function provided to you, and include your plot in your PDF submission.

```

def estimate_accuracy_torch(model, data, batch_size=5000, max_N=100000):
    """
    Estimate the accuracy of the model on the data. To reduce
    computation time, use at most `max_N` elements of `data` to
    produce the estimate.
    """
    correct = 0
    N = 0
    for i in range(0, data.shape[0], batch_size):
        # get a batch of data
        xt, st = get_batch(data, i, i + batch_size, onehot=False)

        # forward pass prediction
        y = model(torch.Tensor(xt))
        y = y.detach().numpy() # convert the PyTorch tensor => numpy array
        pred = np.argmax(y, axis=1)
        correct += np.sum(pred == st)
        N += st.shape[0]

    if N > max_N:
        break
    return correct / N

def run_pytorch_gradient_descent(model,
                                train_data=train4grams,
                                validation_data=valid4grams,
                                batch_size=100,
                                learning_rate=0.001,
                                weight_decay=0,
                                max_iters=1000,
                                checkpoint_path=None):
    """
    Train the PyTorch model on the dataset `train_data`, reporting
    the validation accuracy on `validation_data`, for `max_iters`
    iteration.

    If you want to **checkpoint** your model weights (i.e. save the
    model weights to Google Drive), then the parameter
    `checkpoint_path` should be a string path with `{}` to be replaced
    by the iteration count:

    For example, calling

    >>> run_pytorch_gradient_descent(model, ...,
                                     checkpoint_path = '/content/gdrive/My Drive/Intro_to_Deep_Learning/mlp/ckpt-{}.pk')

    will save the model parameters in Google Drive every 500 iterations.
    You will have to make sure that the path exists (i.e. you'll need to create
    the folder Intro_to_Deep_Learning, mlp, etc...). Your Google Drive will be populated with files:

    - /content/gdrive/My Drive/Intro_to_Deep_Learning/mlp/ckpt-500.pk
    - /content/gdrive/My Drive/Intro_to_Deep_Learning/mlp/ckpt-1000.pk
    - ...

    To load the weights at a later time, you can run:

    >>> model.load_state_dict(torch.load('/content/gdrive/My Drive/Intro_to_Deep_Learning/mlp/ckpt-500.pk'))

    This function returns the training loss, and the training/validation accuracy,
    which we can use to plot the learning curve.
    """
    criterion = nn.CrossEntropyLoss()
    optimizer = optim.Adam(model.parameters()),

```

```

        lr=learning_rate,
        weight_decay=weight_decay)

iters, losses = [], []
iters_sub, train_accs, val_accs = [], [], []

n = 0 # the number of iterations
while True:
    for i in range(0, train_data.shape[0], batch_size):
        if (i + batch_size) > train_data.shape[0]:
            break

        # get the input and targets of a minibatch
        xt, st = get_batch(train_data, i, i + batch_size, onehot=False)

        # convert from numpy arrays to PyTorch tensors
        xt = torch.Tensor(xt)
        st = torch.Tensor(st).long()

        zs = model(xt)          # compute prediction logit
        loss = criterion(zs, st) # compute the total loss
        optimizer.zero_grad()   # compute updates for each parameter
        loss.backward()          # make the updates for each parameter
        optimizer.step()         # a clean up step for PyTorch

        # save the current training information
        iters.append(n)
        losses.append(float(loss)/batch_size) # compute *average* loss

    if n % 500 == 0:
        iters_sub.append(n)
        train_cost = float(loss.detach().numpy())
        train_acc = estimate_accuracy_torch(model, train_data)
        train_accs.append(train_acc)
        val_acc = estimate_accuracy_torch(model, validation_data)
        val_accs.append(val_acc)
        print("Iter %d. [Val Acc %.0f%%] [Train Acc %.0f%%, Loss %f]" % (
            n, val_acc * 100, train_acc * 100, train_cost))

        if (checkpoint_path is not None) and n > 0:
            torch.save(model.state_dict(), checkpoint_path.format(n))

    # increment the iteration number
    n += 1

    if n > max_iters:
        return iters, losses, iters_sub, train_accs, val_accs

def plot_learning_curve(iters, losses, iters_sub, train_accs, val_accs):
    """
    Plot the learning curve.
    """
    plt.title("Learning Curve: Loss per Iteration")
    plt.plot(iters, losses, label="Train")
    plt.xlabel("Iterations")
    plt.ylabel("Loss")
    plt.show()

    plt.title("Learning Curve: Accuracy per Iteration")
    plt.plot(iters_sub, train_accs, label="Train")
    plt.plot(iters_sub, val_accs, label="Validation")
    plt.xlabel("Iterations")
    plt.ylabel("Accuracy")
    plt.legend(loc='best')
    plt.show()

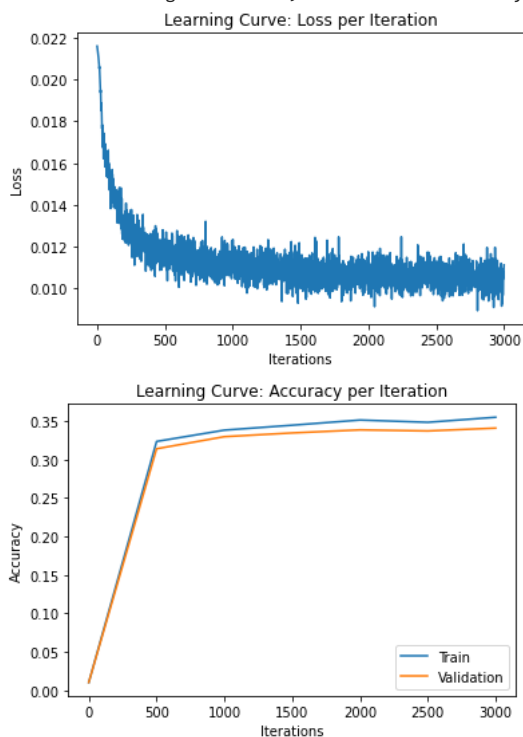
# see which rate is best between [0.0001, 0.001, 0.01, 0.1], and use this rate
max_acc = 0
learning_rates = [0.0001, 0.001, 0.01, 0.1]
for rate in learning_rates:
    print(f"This is the information for learning rate {rate}:")
    pytorch_mlp_temp = PyTorchMLP()
    iters, losses, iters_sub, train_accs, val_accs = run_pytorch_gradient_descent(pytorch_mlp_temp,
                                                                                  batch_size=256, learning_rate=rate,
                                                                                  weight_decay=0, max_iters=3000)

    if train_accs[-1] > max_acc:
        best_rate = rate
        max_acc = train_accs[-1]
        pytorch_mlp = pytorch_mlp_temp
        learning_curve_info = iters, losses, iters_sub, train_accs, val_accs

```

```
print(f"Plot for learning rate {best_rate}, with maximus accuracy of {max_acc:.2f}.")
plot_learning_curve(*learning_curve_info)
```

```
This is the information for learning rate 0.0001:
Iter 0. [Val Acc 2%] [Train Acc 2%, Loss 5.511062]
Iter 500. [Val Acc 18%] [Train Acc 19%, Loss 4.208350]
Iter 1000. [Val Acc 23%] [Train Acc 24%, Loss 4.011002]
Iter 1500. [Val Acc 26%] [Train Acc 27%, Loss 3.672452]
Iter 2000. [Val Acc 28%] [Train Acc 28%, Loss 3.264605]
Iter 2500. [Val Acc 29%] [Train Acc 30%, Loss 3.163743]
Iter 3000. [Val Acc 30%] [Train Acc 30%, Loss 3.227946]
This is the information for learning rate 0.001:
Iter 0. [Val Acc 1%] [Train Acc 1%, Loss 5.531754]
Iter 500. [Val Acc 31%] [Train Acc 32%, Loss 2.968134]
Iter 1000. [Val Acc 33%] [Train Acc 34%, Loss 3.099179]
Iter 1500. [Val Acc 33%] [Train Acc 34%, Loss 2.760988]
Iter 2000. [Val Acc 34%] [Train Acc 35%, Loss 2.723919]
Iter 2500. [Val Acc 34%] [Train Acc 35%, Loss 2.680916]
Iter 3000. [Val Acc 34%] [Train Acc 36%, Loss 2.846323]
This is the information for learning rate 0.01:
Iter 0. [Val Acc 17%] [Train Acc 18%, Loss 5.518847]
Iter 500. [Val Acc 30%] [Train Acc 31%, Loss 3.164366]
Iter 1000. [Val Acc 30%] [Train Acc 31%, Loss 3.300217]
Iter 1500. [Val Acc 30%] [Train Acc 32%, Loss 2.971757]
Iter 2000. [Val Acc 31%] [Train Acc 32%, Loss 2.868073]
Iter 2500. [Val Acc 31%] [Train Acc 31%, Loss 2.901245]
Iter 3000. [Val Acc 31%] [Train Acc 33%, Loss 3.054958]
This is the information for learning rate 0.1:
Iter 0. [Val Acc 19%] [Train Acc 19%, Loss 5.524472]
Iter 500. [Val Acc 17%] [Train Acc 17%, Loss 15.559907]
Iter 1000. [Val Acc 16%] [Train Acc 17%, Loss 17.489374]
Iter 1500. [Val Acc 16%] [Train Acc 17%, Loss 17.116291]
Iter 2000. [Val Acc 17%] [Train Acc 17%, Loss 16.755619]
Iter 2500. [Val Acc 17%] [Train Acc 17%, Loss 13.363822]
Iter 3000. [Val Acc 15%] [Train Acc 16%, Loss 14.480177]
Plot for learning rate 0.001, with maximus accuracy of 0.35.
```



▼ Part (c) – 10%

Write a function `make_prediction` that takes as parameters a PyTorchMLP model and sentence (a list of words), and produces a prediction for the next word in the sentence.

```
def make_prediction_torch(model, sentence):
    """
    Use the model to make a prediction for the next word in the
    sentence using the last 3 words (sentence[:-3]). You may assume
    that len(sentence) >= 3 and that `model` is an instance of
    PYTorchMLP.
```

This function should return the next word, represented as a string.

Example call:

```
>>> make_prediction_torch(pytorch_mlp, ['you', 'are', 'a'])
```

```

"""
global vocab_stoi, vocab_itos

# Write your code here
sentence = [sentence[-3:] + ['?']]
sentence_input = make_onehot((process_data(sentence))[:, :3])

output = model(torch.Tensor(sentence_input))
output = output.detach().numpy()

return vocab_itos[np.argmax(output, axis=1)[0]]

```

▼ Part (d) -- 10%

Use your code 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"

Do your predictions make sense?

In many cases where you overfit the model can either output the same results for all inputs or just memorize the dataset.

Print the output for all of these sentences and **Write** below if you encounter these effects or something else which indicates overfitting, if you do train again with better hyperparameters.

```

# Write your code here

def next_word(model, sentences):
    print("Our predictions are:")
    for sentence in sentences:
        words = sentence.split()
        words = [word.lower() for word in words]
        next = make_prediction_torch(model, words)
        print(f"'{sentence}' -> '{next}'")

sentences = ["You are a", "few companies show", "There are no", "yesterday i was", "the game had", "yesterday the federal"]
next_word(pytorch_mlp, sentences)

Our predictions are:
'You are a' -> 'good'
'few companies show' -> '.'
'There are no' -> 'more'
'yesterday i was' -> 'nt'
'the game had' -> 'to'
'yesterday the federal' -> 'government'

```

Write your answers here:

Most of our predictions make sense. In fact, the only sentence that shows bad results is "Few companies show.". The dot resembles sentences from the dataset like "This was a show.", which is of course a different meaning of the word "show".

Some bad results were expected, since we have accuracy of 36% when training the model.

In addition, we went through the raw data and found only "You are a good" appear in the actual sentences. This alone could have ment we're overfitting, but combined with the other results (which do not match the dataset) and noting that we output different results for each sentence, we do not see signs of overfitting.

Despite the above, we should consider that in the English language some words appear more than others. Therefore, the dataset used is not a balanced one and some word may have a higher chance of being predicted.

▼ Part (e) -- 6%

Report the test accuracy of your model

```

# Write your code here
test_accuracy = estimate_accuracy_torch(pytorch_mlp, test4grams)

print(f"The test accuracy is: {(test_accuracy*100):.2f} %")

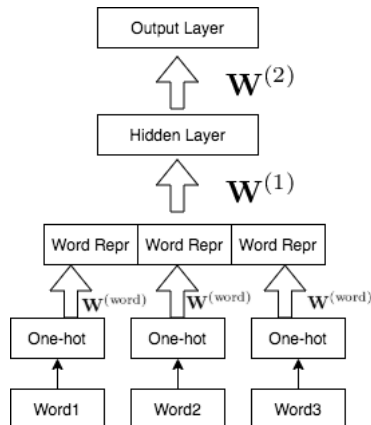
The test accuracy is: 34.22 %

```


▼ Question 3. Learning Word Embeddings (24 %)

In this section, we will build a slightly different model with a different architecture. In particular, we will first compute a lower-dimensional *representation* of the three words, before using a multi-layer perceptron.

Our model will look like this:



This model has 3 layers instead of 2, but the first layer of the network is **not** fully-connected. Instead, we compute the representations of each of the three words **separately**. In addition, the first layer of the network will not use any biases. The reason for this will be clear in question 4.

▼ Part (a) -- 10%

The PyTorch model is implemented for you. Use `run_pytorch_gradient_descent` to train your PyTorch MLP model to obtain a training accuracy of at least 38%. Plot the learning curve using the `plot_learning_curve` function provided to you, and include your plot in your PDF submission.

```
class PyTorchWordEmb(nn.Module):
    def __init__(self, emb_size=100, num_hidden=300, vocab_size=250):
        super(PyTorchWordEmb, self).__init__()
        self.word_emb_layer = nn.Linear(vocab_size, emb_size, bias=False)
        self.fc_layer1 = nn.Linear(emb_size * 3, num_hidden)
        self.fc_layer2 = nn.Linear(num_hidden, 250)
        self.num_hidden = num_hidden
        self.emb_size = emb_size

    def forward(self, inp):
        embeddings = torch.relu(self.word_emb_layer(inp))
        embeddings = embeddings.reshape([-1, self.emb_size * 3])
        hidden = torch.relu(self.fc_layer1(embeddings))
        return self.fc_layer2(hidden)

pytorch_wordemb = PyTorchWordEmb()

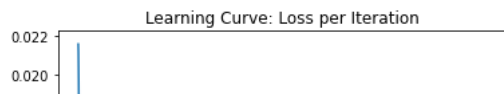
result = run_pytorch_gradient_descent(pytorch_wordemb, max_iters=20000, batch_size=256)

# plot_learning_curve(*result)
```

```

Iter 0. [Val Acc 2%] [Train Acc 2%, Loss 5.528602]
Iter 500. [Val Acc 28%] [Train Acc 28%, Loss 3.228059]
Iter 1000. [Val Acc 31%] [Train Acc 32%, Loss 3.256534]
Iter 1500. [Val Acc 32%] [Train Acc 33%, Loss 2.976914]
Iter 2000. [Val Acc 33%] [Train Acc 34%, Loss 2.820710]
Iter 2500. [Val Acc 34%] [Train Acc 35%, Loss 2.752198]
Iter 3000. [Val Acc 34%] [Train Acc 35%, Loss 2.872323]
Iter 3500. [Val Acc 35%] [Train Acc 36%, Loss 2.520480]
Iter 4000. [Val Acc 35%] [Train Acc 36%, Loss 2.780762]
Iter 4500. [Val Acc 35%] [Train Acc 37%, Loss 2.548871]
Iter 5000. [Val Acc 36%] [Train Acc 37%, Loss 2.505401]
Iter 5500. [Val Acc 36%] [Train Acc 37%, Loss 2.575140]
Iter 6000. [Val Acc 36%] [Train Acc 38%, Loss 2.685529]
Iter 6500. [Val Acc 36%] [Train Acc 38%, Loss 2.431973]
Iter 7000. [Val Acc 37%] [Train Acc 38%, Loss 2.533333]
Iter 7500. [Val Acc 37%] [Train Acc 39%, Loss 2.708552]
Iter 8000. [Val Acc 37%] [Train Acc 39%, Loss 2.379351]
Iter 8500. [Val Acc 37%] [Train Acc 38%, Loss 2.443505]
Iter 9000. [Val Acc 37%] [Train Acc 39%, Loss 2.511275]
Iter 9500. [Val Acc 37%] [Train Acc 39%, Loss 2.291260]
Iter 10000. [Val Acc 37%] [Train Acc 39%, Loss 2.370033]
Iter 10500. [Val Acc 37%] [Train Acc 40%, Loss 2.270216]
Iter 11000. [Val Acc 37%] [Train Acc 39%, Loss 2.330172]
Iter 11500. [Val Acc 37%] [Train Acc 39%, Loss 2.416735]
Iter 12000. [Val Acc 37%] [Train Acc 40%, Loss 2.494140]
Iter 12500. [Val Acc 37%] [Train Acc 39%, Loss 2.495757]
Iter 13000. [Val Acc 37%] [Train Acc 39%, Loss 2.259490]
Iter 13500. [Val Acc 37%] [Train Acc 40%, Loss 2.292390]
Iter 14000. [Val Acc 38%] [Train Acc 40%, Loss 2.438367]
Iter 14500. [Val Acc 37%] [Train Acc 40%, Loss 2.324364]
Iter 15000. [Val Acc 38%] [Train Acc 40%, Loss 2.291252]
Iter 15500. [Val Acc 38%] [Train Acc 40%, Loss 2.448252]
Iter 16000. [Val Acc 38%] [Train Acc 41%, Loss 2.403137]
Iter 16500. [Val Acc 38%] [Train Acc 40%, Loss 2.368816]
Iter 17000. [Val Acc 38%] [Train Acc 40%, Loss 2.302421]
Iter 17500. [Val Acc 38%] [Train Acc 41%, Loss 2.197944]
Iter 18000. [Val Acc 38%] [Train Acc 41%, Loss 2.362378]
Iter 18500. [Val Acc 38%] [Train Acc 40%, Loss 2.371853]
Iter 19000. [Val Acc 38%] [Train Acc 41%, Loss 2.481299]
Iter 19500. [Val Acc 38%] [Train Acc 41%, Loss 2.430783]
Iter 20000. [Val Acc 38%] [Train Acc 40%, Loss 2.484284]

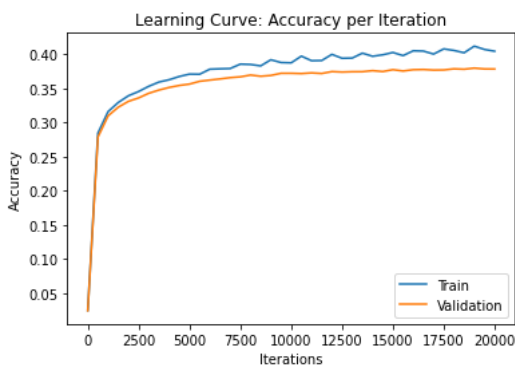
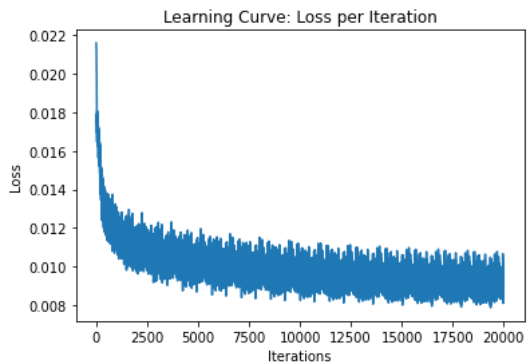
```



```

# adding plot for PDF saving
plot_learning_curve(*result)

```



▼ Part (b) -- 10%

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 $\geq 38\%$.

```
# Your code goes here
next_word(pytorch_wordemb, sentences)

Our predictions are:
'You are a' -> 'good'
'few companies show' -> '.'
'There are no' -> 'people'
'yesterday i was' -> 'nt'
'the game had' -> 'to'
'yesterday the federal' -> 'government'
```

Write your explanation here:

We see that the results produced with this model are the same as the MLP results.

We will note that with the MLP model we reached a 36% accuracy, and with Word Embeddings we reached 41%, so while it is better on paper, both models are imperfect and are expected to have mistakes. Once again our training does not seem to overfit the dataset.

▼ Part (c) -- 4%

Report the test accuracy of your model

```
# Write your code here
test_accuracy_wordemb = estimate_accuracy_torch(pytorch_wordemb, test4grams)

print(f"The test accuracy is: {(test_accuracy_wordemb*100):.2f} %")

The test accuracy is: 38.27 %
```

▼ Question 4. Visualizing Word Embeddings (14%)

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) -- 4%

The code below extracts the **weights** of the word embedding layer, and converts the PyTorch tensor into an 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".

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

Write your explanation here:

As represented in the graph for question 3, $wordRepr = weights \cdot oneHot$, Where $wordRepr$ and $oneHot$ are column vectors.

Since $oneHot$ is 1 at the requested word and 0 elsewhere, we expect each *column* of weights to represent a word.

In the code line above, we see that $wordEmb = weights^T$, therefor each *row* in `word_emb` represents a word.

▼ 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 \mathbf{v} and \mathbf{w} is defined as

$$d_{\cos}(\mathbf{v}, \mathbf{w}) = \frac{\mathbf{v}^T \mathbf{w}}{\|\mathbf{v}\| \|\mathbf{w}\|}.$$

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

```
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(similarities[vocab_stoi['any'], vocab_stoi['many']])
print(similarities[vocab_stoi['any'], vocab_stoi['government']])

0.24776883
0.075421795
```

Compute the 5 closest words to the following words:

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

Write your code here

```
words_list = ["four", "go", "what", "should", "school", "your", "yesterday", "not"]

for word in words_list:
    # word_ind = vocab_stoi[word]
    distances = similarities[vocab_stoi[word], :]
    closest_words = distances.argsort()[-6:-1][::-1]
    print("\nThe 5 closest words to \"{0}\" are:\n".format(word))
    for i, closest_word_index in enumerate(closest_words):
        print("{0}) \"{1}\" with similarity of {2:.4f}".format(str(i+1), vocab_itos[closest_word_index], distances[closest_word_index]))

    4) "few" with similarity of 0.3274
    5) "office" with similarity of 0.2861

    The 5 closest words to "go" are:

    1) "going" with similarity of 0.5233
    2) "back" with similarity of 0.4410
    3) "come" with similarity of 0.4326
    4) "right" with similarity of 0.3775
    5) "down" with similarity of 0.3315

    The 5 closest words to "what" are:

    1) "how" with similarity of 0.4273
    2) "where" with similarity of 0.4130
    3) "who" with similarity of 0.3691
    4) "which" with similarity of 0.3621
    5) "when" with similarity of 0.3538

    The 5 closest words to "should" are:

    1) "could" with similarity of 0.5471
    2) "would" with similarity of 0.5056
    3) "can" with similarity of 0.4666
    4) "will" with similarity of 0.4468
    5) "ago" with similarity of 0.3456

    The 5 closest words to "school" are:

    1) "week" with similarity of 0.3908
    2) "office" with similarity of 0.3897
    3) "market" with similarity of 0.3820
```

```

2) our with similarity of 0.3577
3) "mr." with similarity of 0.4017
4) "her" with similarity of 0.3917
5) "the" with similarity of 0.3372

```

The 5 closest words to "yesterday" are:

```

1) "times" with similarity of 0.5218
2) "today" with similarity of 0.4905
3) "street" with similarity of 0.3774
4) "ago" with similarity of 0.3629
5) "night" with similarity of 0.3591

```

The 5 closest words to "not" are:

```

1) "nt" with similarity of 0.4575
2) "never" with similarity of 0.4336
3) "national" with similarity of 0.3725
4) "even" with similarity of 0.3528
5) "office" with similarity of 0.3442

```

▼ Part (c) -- 5%

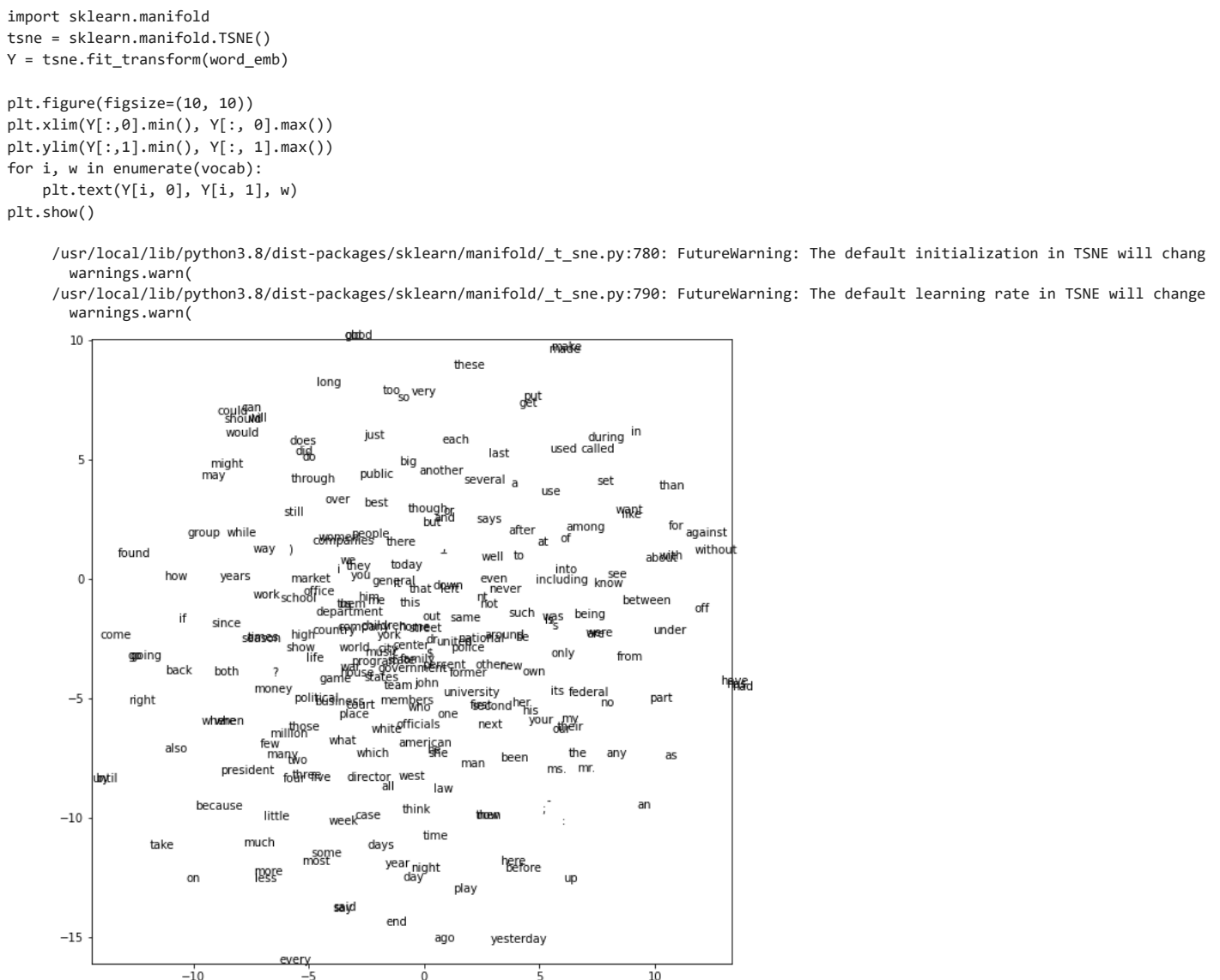
We can visualize the word embeddings by reducing the dimensionality of the word vectors to 2D. There are many dimensionality reduction techniques that we could use, and we will use an algorithm called t-SNE. (You don't need to know what this is for the assignment; we will cover it later in the course.)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



Explain and discuss your results here:

We can see clusters such as:

- Might, May
- More, less
- Said, Say
- Two, Three, Four, Five
- Does, Do, Did
- Go, Going
- When, Where
- Could, Should, Can

These seem great, as some clusters have semantic similarities and some are of the same root.

✓ 0s completed at 2:20 PM

● ×