

▼ 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.

```
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/My Drive/deep_learning/raw_sentences.txt' # TODO - UPDATE ME!
```

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.

```
# Your code goes here
for sentence in train[10:20]:
    print(sentence)

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

Write your answers here:

From the above sentences we can see that:

- Punctuations are treated as words in our word representation.
- Words with apostrophes are separated into the parts before and after the apostrophes, like in the word "she's" on the second sentence printed.

▼ **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.

```
print(collections.Counter([item for sublist in train for item in sublist]).most_common(10))

[(('.', 64297), ('it', 23118), (',', 19537), ('i', 17684), ('do', 16181), ('to', 15490),
```


▼ **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.

```
# 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', 'or
    [[148, 98, 70, 23, 154, 89], [151, 148, 181, 246], [248]]
    """
    return [[vocab_stoi[j.lower()] for j in i] for i in sents]

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]]
    """
    new_list = []
    for sequence in seqs:
        for i in range(len(sequence)-3):
            if len(sequence[i:i+4]) == 4 :
                new_list.append(sequence[i:i+4])
    return new_list

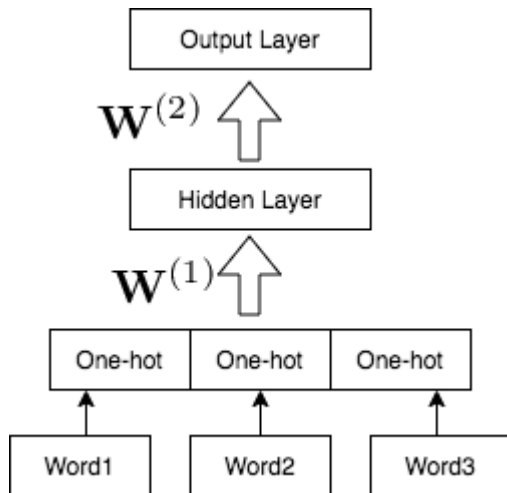
def process_data(sents):
    """
    This function takes a list of sentences (list of lists), and generates an
    numpy matrix with shape [N, 4] containing indices of words in 4-grams.
    """
    indices = convert_words_to_indices(sents)
    fourgrams = generate_4grams(indices)
    return np.array(fourgrams)

# 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)
    """

    if data.ndim==1:
        shape = (data.size, 250)
        one_hot = np.zeros(shape)
        rows = np.arange(data.size)
        one_hot[rows, data] = 1
        return one_hot
    else :
        shape = (data.shape[0],data.shape[1], 250)

```

```

one_hot = np.zeros(shape)
rows = np.arange(data.shape[0])
one_hot[rows,0, data[:,0]] = 1
one_hot[rows,1, data[:,1]] = 1
one_hot[rows,2, data[:,2]] = 1
return one_hot

```

```

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])
        hidden = torch.relu(self.layer1(inp))
        return self.layer2(hidden)

```

▼ 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=5000,
                                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

- /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 = ...                # compute prediction logit
        # loss =                  # compute the total loss
        # ...                     # compute updates for each parameter
        # ...                     # make the updates for each parameter
        # ...                     # a clean up step for PyTorch

        zs = model(xt)
        loss = criterion(zs,st)
        optimizer.zero_grad()
        loss.backward()
```



```

optimizer.step()

# 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()

pytorch_mlp = PyTorchMLP()
# learning_curve_info = run_pytorch_gradient_descent(pytorch_mlp, ...)
learning_curve_info = run_pytorch_gradient_descent(pytorch_mlp)#, ...,checkpoint_path = '/cor

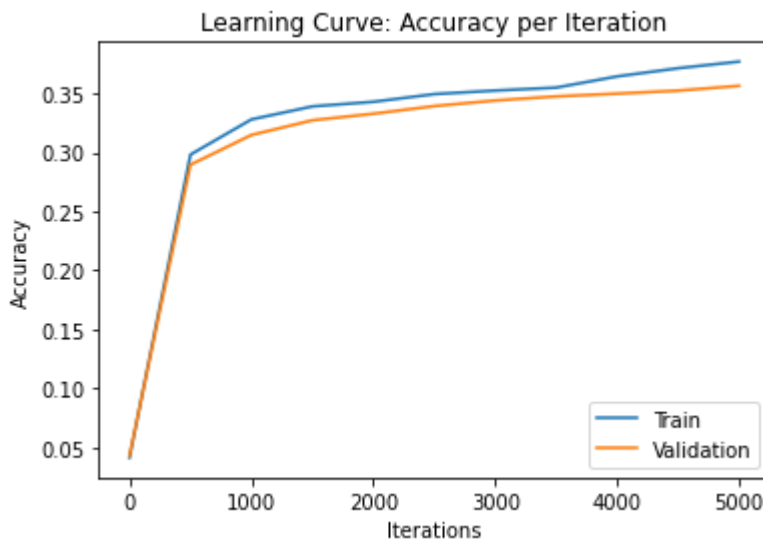
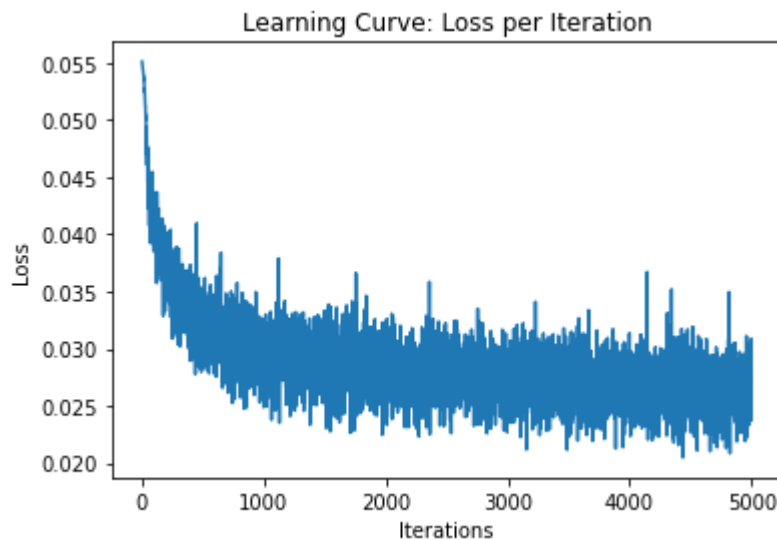
# plot_learning_curve(*learning_curve_info)
plot_learning_curve(*learning_curve_info)

```

```

Iter 0. [Val Acc 4%] [Train Acc 4%, Loss 5.511876]
Iter 500. [Val Acc 29%] [Train Acc 30%, Loss 2.926858]
Iter 1000. [Val Acc 31%] [Train Acc 33%, Loss 2.719736]
Iter 1500. [Val Acc 33%] [Train Acc 34%, Loss 2.686385]
Iter 2000. [Val Acc 33%] [Train Acc 34%, Loss 2.636719]
Iter 2500. [Val Acc 34%] [Train Acc 35%, Loss 2.623256]
Iter 3000. [Val Acc 34%] [Train Acc 35%, Loss 2.459397]
Iter 3500. [Val Acc 35%] [Train Acc 35%, Loss 2.417365]
Iter 4000. [Val Acc 35%] [Train Acc 36%, Loss 2.605546]
Iter 4500. [Val Acc 35%] [Train Acc 37%, Loss 2.611154]
Iter 5000. [Val Acc 36%] [Train Acc 38%, Loss 2.647678]

```



▼ 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
```

```
#vocab_itos = dict(enumerate(vocab))
```

```
# A mapping of word => its index
```

```
#vocab_stoi = {word:index for index, word in vocab_itos.items()}
```

```
#print(sentence[:3])
```

```
#print(convert_words_to_indices(sentence[:3]))
```

```
data = np.array(convert_words_to_indices([sentence]))
```

```
a = make_onehot(data)
```

```
b = model(torch.Tensor(a))
```

```
c = np.argmax(b.detach().numpy())
```

```
return vocab_itos[c]
```

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

```
good
```

▼ 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
```

```

predic1 = make_prediction_torch(pytorch_mlp,['you','are','a'])
predic2 = make_prediction_torch(pytorch_mlp,['few','companies','show'])
predic3 = make_prediction_torch(pytorch_mlp,['There','are','no'])
predic4 = make_prediction_torch(pytorch_mlp,['yesterday','i','was'])
predic5 = make_prediction_torch(pytorch_mlp,['the','game','had'])
predic6 = make_prediction_torch(pytorch_mlp,['yesterday','the','federal'])

```

```

print(predic1)
print(predic2)
print(predic3)
print(predic4)
print(predic5)
print(predic6)

```

```

good
.
other
nt
been
government

```

Write your answers here:

Luckily it seems our output more or less makes sense.

The only exception is the input "few companies show" which probably should have resulted in something more like "success", "promise" or so, as sentences do not usually end with a verb. This may be due to partial memorization of the dataset.

▼ Part (e) -- 6%

Report the test accuracy of your model

```

test_acc = estimate_accuracy_torch(pytorch_mlp, test4grams)
print("Test accuracy is %.0f%%." %( test_acc * 100))

```

```

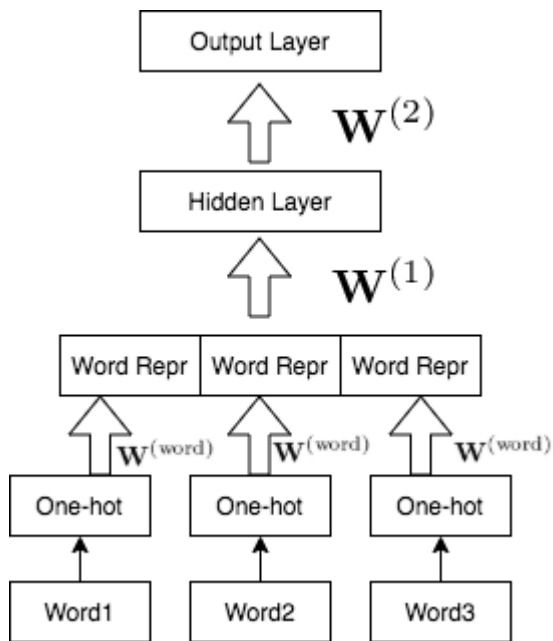
Test accuracy is 36%.

```

▼ 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.

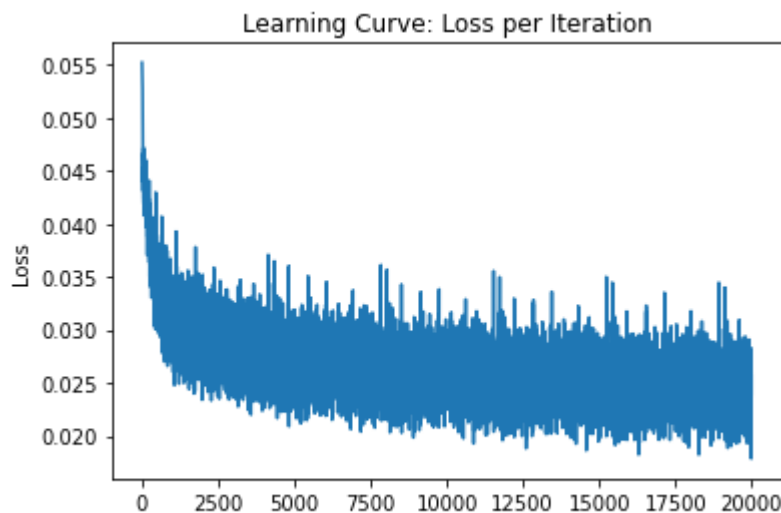
```
class PyTorchWordEmb(nn.Module):
    def __init__(self, emb_size=100, num_hidden=300, vocab_size=250):
        super(PyTorchWordEmb, self).__init__()
        self.word_emb_layer = nn.Linear(vocab_size, emb_size, bias=False)
        self.fc_layer1 = nn.Linear(emb_size * 3, num_hidden)
        self.fc_layer2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

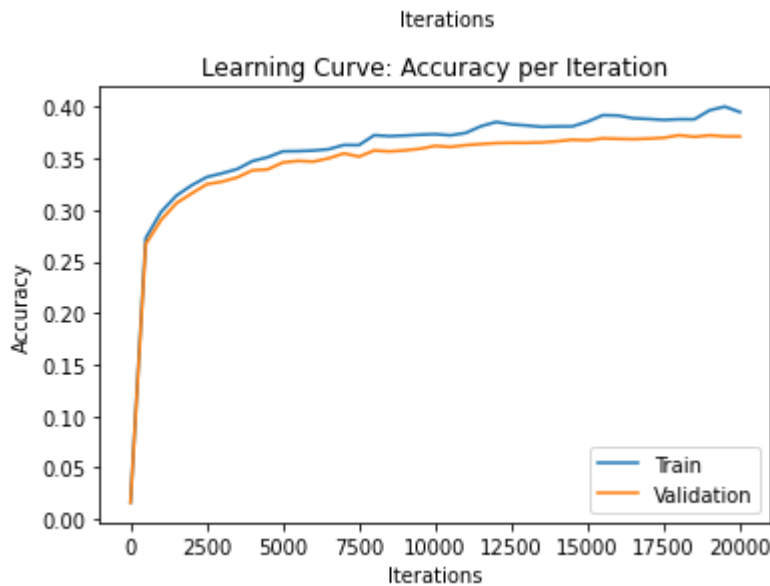
```
pytorch_wordemb= PyTorchWordEmb()
```

```
result = run_pytorch_gradient_descent(pytorch_wordemb,max_iters=20000)
```

```
plot_learning_curve(*result)
```

```
↳ Iter 0. [Val Acc 2%] [Train Acc 2%, Loss 5.521066]
Iter 500. [Val Acc 27%] [Train Acc 27%, Loss 3.116298]
Iter 1000. [Val Acc 29%] [Train Acc 30%, Loss 2.881423]
Iter 1500. [Val Acc 31%] [Train Acc 31%, Loss 2.800914]
Iter 2000. [Val Acc 32%] [Train Acc 32%, Loss 2.757010]
Iter 2500. [Val Acc 32%] [Train Acc 33%, Loss 2.849174]
Iter 3000. [Val Acc 33%] [Train Acc 34%, Loss 2.554440]
Iter 3500. [Val Acc 33%] [Train Acc 34%, Loss 2.464811]
Iter 4000. [Val Acc 34%] [Train Acc 35%, Loss 2.748654]
Iter 4500. [Val Acc 34%] [Train Acc 35%, Loss 2.722729]
Iter 5000. [Val Acc 35%] [Train Acc 36%, Loss 2.797253]
Iter 5500. [Val Acc 35%] [Train Acc 36%, Loss 2.530044]
Iter 6000. [Val Acc 35%] [Train Acc 36%, Loss 2.724496]
Iter 6500. [Val Acc 35%] [Train Acc 36%, Loss 2.570428]
Iter 7000. [Val Acc 35%] [Train Acc 36%, Loss 2.967023]
Iter 7500. [Val Acc 35%] [Train Acc 36%, Loss 2.720179]
Iter 8000. [Val Acc 36%] [Train Acc 37%, Loss 2.203897]
Iter 8500. [Val Acc 36%] [Train Acc 37%, Loss 2.833435]
Iter 9000. [Val Acc 36%] [Train Acc 37%, Loss 2.437212]
Iter 9500. [Val Acc 36%] [Train Acc 37%, Loss 2.672007]
Iter 10000. [Val Acc 36%] [Train Acc 37%, Loss 2.296863]
Iter 10500. [Val Acc 36%] [Train Acc 37%, Loss 2.585612]
Iter 11000. [Val Acc 36%] [Train Acc 37%, Loss 2.622194]
Iter 11500. [Val Acc 36%] [Train Acc 38%, Loss 2.640025]
Iter 12000. [Val Acc 36%] [Train Acc 39%, Loss 2.445100]
Iter 12500. [Val Acc 37%] [Train Acc 38%, Loss 2.510798]
Iter 13000. [Val Acc 37%] [Train Acc 38%, Loss 2.209301]
Iter 13500. [Val Acc 37%] [Train Acc 38%, Loss 2.809682]
Iter 14000. [Val Acc 37%] [Train Acc 38%, Loss 2.580041]
Iter 14500. [Val Acc 37%] [Train Acc 38%, Loss 2.752034]
Iter 15000. [Val Acc 37%] [Train Acc 39%, Loss 2.256427]
Iter 15500. [Val Acc 37%] [Train Acc 39%, Loss 2.495715]
Iter 16000. [Val Acc 37%] [Train Acc 39%, Loss 2.210522]
Iter 16500. [Val Acc 37%] [Train Acc 39%, Loss 2.379016]
Iter 17000. [Val Acc 37%] [Train Acc 39%, Loss 2.468591]
Iter 17500. [Val Acc 37%] [Train Acc 39%, Loss 2.527767]
Iter 18000. [Val Acc 37%] [Train Acc 39%, Loss 2.730582]
Iter 18500. [Val Acc 37%] [Train Acc 39%, Loss 2.302812]
Iter 19000. [Val Acc 37%] [Train Acc 40%, Loss 2.409322]
Iter 19500. [Val Acc 37%] [Train Acc 40%, Loss 2.537691]
Iter 20000. [Val Acc 37%] [Train Acc 39%, Loss 2.396756]
```





▼ 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
predic1 = make_prediction_torch(pytorch_wordemb,['you','are','a'])
predic2 = make_prediction_torch(pytorch_wordemb,['few','companies','show'])
predic3 = make_prediction_torch(pytorch_wordemb,['There','are','no'])
predic4 = make_prediction_torch(pytorch_wordemb,['yesterday','i','was'])
predic5 = make_prediction_torch(pytorch_wordemb,['the','game','had'])
predic6 = make_prediction_torch(pytorch_wordemb,['yesterday','the','federal'])

print(predic1)
print(predic2)
print(predic3)
```

```
print(predic4)
print(predic5)
print(predic6)
```

```
good
up
other
nt
to
government
```

Write your explanation here:

Using this model we got better results than the previous one. the input "few companies show" resulted in the output "up" which makes much more sense.

Also another difference is in the prediction from the input "the game had". The previous model's output was "been" and this model's output is "to". Both seem to be reasonable in our opinion and may be caused by the different algorithms.

▼ Part (c) -- 4%

Report the test accuracy of your model

```
test_acc = estimate_accuracy_torch(pytorch_wordemb, test4grams)
print("Test accuracy is %.0f%%." %( test_acc * 100))
```

```
Test accuracy is 38%.
```

▼ 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 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".


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

Write your explanation here:

`word_emb` is a matrix that when multiplied by a one-hot representation of a word, gives as a result a new vector with lower dimensions (to reduce the amount of computations needed in the actual model).

We remember that one-hot representation is a vector which consists of zeros in all indices except for one.

Therefore in order for the multiplication to produce 1 vector per word, that vector needs to be the row of the matrix on the specific index that is 1 in the one-hot representation. (Matrix multiplication rules)

▼ 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']])

0.2023307
-0.110096626
```

Compute the 5 closest words to the following words:

- "four"
- "go"
- "what"

- "should"
- "school"
- "your"
- "yesterday"
- "not"

```
def get_close_word(word):
    closest_word = "."
    closest_word_similarity = -50
    for index, similarity in enumerate(similarities[vocab_stoi[word]]):
        if(closest_word_similarity < similarity and vocab_itos[index] != word):
            closest_word_similarity = similarity
            closest_word = vocab_itos[index]
    return closest_word

print("closest word to four is", get_close_word('four'))
print("closest word to go is", get_close_word('go'))
print("closest word to what is", get_close_word('what'))
print("closest word to should is", get_close_word('should'))
print("closest word to school is", get_close_word('school'))
print("closest word to your is", get_close_word('your'))
print("closest word to yesterday is", get_close_word('yesterday'))
print("closest word to not is", get_close_word('not'))
```

```
closest word to four is five
closest word to go is come
closest word to what is who
closest word to should is could
closest word to school is money
closest word to your is my
closest word to yesterday is home
closest word to not is nt
```

▼ 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,

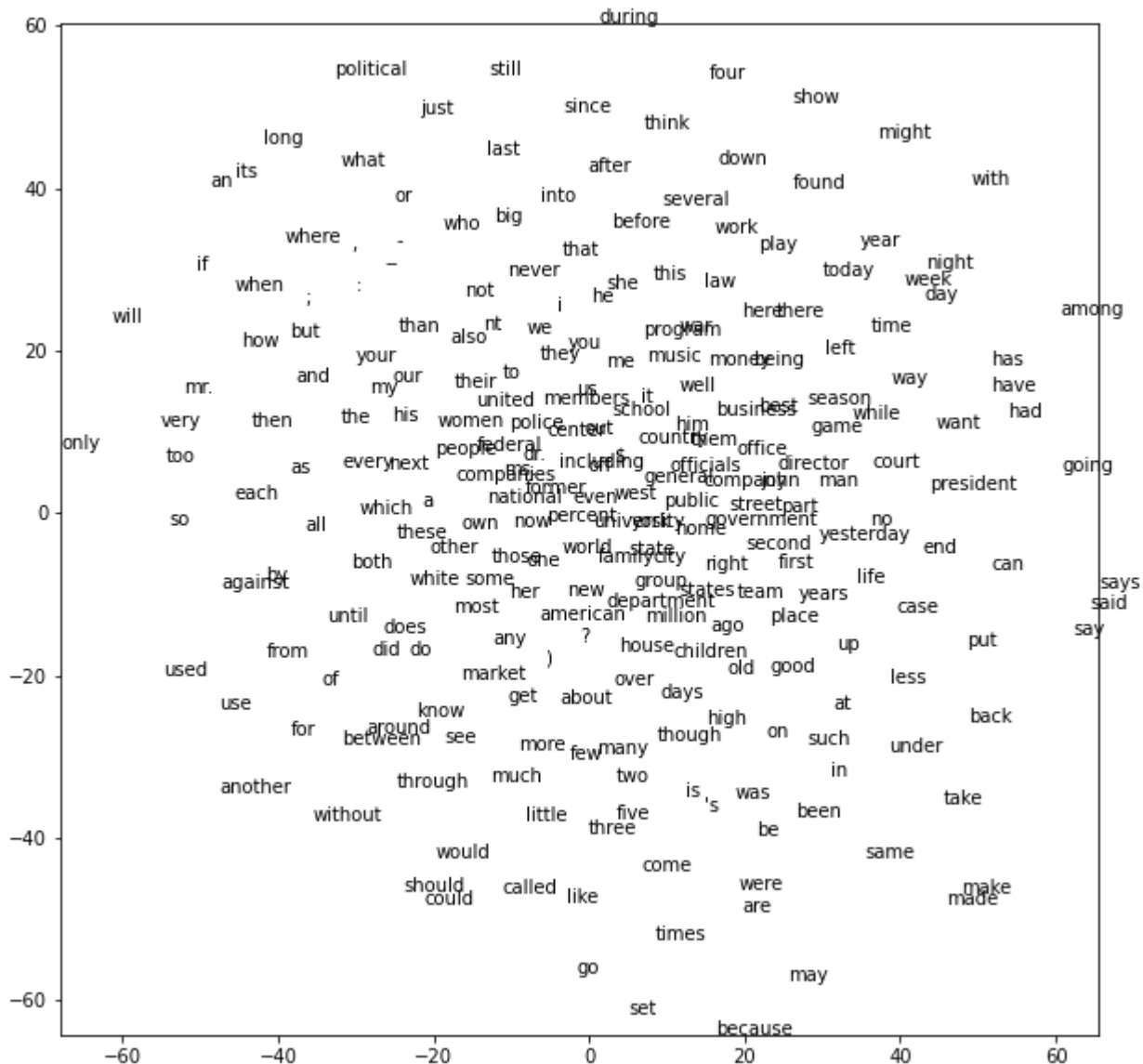
```
import sklearn.manifold
tsne = sklearn.manifold.TSNE()
Y = tsne.fit_transform(word_emb)
```

```
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()
```

```

/usr/local/lib/python3.7/dist-packages/sklearn/manifold/_t_sne.py:783: FutureWarning: T
  FutureWarning,
/usr/local/lib/python3.7/dist-packages/sklearn/manifold/_t_sne.py:793: FutureWarning: T
  FutureWarning,

```



Explain and discuss your results here:

We can see near the coordinates (40, 30) a cluster which contains a few words which are associated with time: "day", "night", "week", "year", "time".

Also near the coordinates (-30, 30) we can see a cluster of punctuation signs such as ; : ' _ - . This seems appropriate as they all belong to the group of punctuation signs and not actual words.

✓ 1s completed at 2:05 PM

