# UDACITY

DISCUSS ON STUDENT HUB

# Sentiment Analysis with Neural Networks

| REVIEW |
| --- |
| HISTORY |

## Requires Changes

1 specification requires changes

## Almost passing . . .

**Overall, you are doing Excellent!! You are almost there.**

- You have forgotten to calculate and display **validation accuracy** in your training implementation.

## Please read through the comments above, correct the issues and resubmit.

**Good Luck with your re-submission AND with your remaining 2 projects.**

## It has been a pleasure reviewing your project.

# As for your question posted in the Student Notes . . .

I think this answer at **Udacity's Knowledge Base -
https://knowledge.udacity.com/questions/30112** answers things much better than I
ever could:

Hi **Hahnsang,**

This error is being caused by your LSTM layer. This project is slightly different from the
Sentiment RNN you saw in the lessons, one difference is that in this project the inputs
come in shape (seq_length, batch_size) instead of (batch_size, seq_length) as in the lesson.
Because of this, your LSTM layer should have the parameter:

```
1    batch_first=False
```

Another difference is the forward pass. As **Ajit** said you shouldn't be reshaping the
lstm_out instead you should obtain the **last output of the sequence,** so if the outputs
from your LSTM are shape (seq_length, batch_size, lstm_size) you should do:

```
1    lstm_out[-1,:,:]
```

instead of:

```
1    lstm_out.contiguous().view(-1, self.lstm_size)
```

After you make this change you just have to pass the outputs through the fully connected
layer and logsoftmax without any reshaping, and then returning the logps and the hidden
state from the LSTM.

(edited)

**Ricardo D** almost 3 years ago

↩ **SHOW 3 COMMENTS**                                          **REPORT**

## Importing Twits

Print the number of twits in the dataset.

# Exactly correct!

- You were asked to print out the number of twits encountered (which required you to properly read in the provided **twits.json** file and then understand how to **count** the number of entries) and your code does that properly.

- Nice job – you got the correct number - there are expected to be: **1548010**.

**Length of Data**

Now let's look at the number of twits in dataset. Print the number of twits below.

```
1  """print out the number of twits"""
2
3  # TODO Implement
4  twits_total = len(twits['data'])
5  twits_total
```

: 1548010

## Preprocessing the Data

The function  `preprocess`  correctly lowercases, removes URLs, removes ticker symbols, removes punctuation, tokenizes, and removes any single character tokens.

# All requested preprocessing of the tweets has been implemented correctly

You are correctly:

- Lowercasing
- Removing urls
- Removing ticker symbols
- Removing usernames
- Removing non-letter words
- Removing any tokens that are not at least 2 characters long.

Excellent work!

```
1   nltk.download('wordnet')
2
3   def preprocess(message):
4       """
5       This function takes a string as input, then performs these operations:
6           - lowercase
7           - remove URLs
8           - remove ticker symbols
9           - removes punctuation
10          - tokenize by splitting the string on whitespace
11          - removes any single character tokens
12
13      Parameters
14      ----------
15          message : The text message to be preprocessed.
16
17      Returns
18      -------
19          tokens: The preprocessed text into tokens.
20      """
21      #TODO: Implement
22
23      # Lowercase the twit message
24      text = message.lower()
25
26      # Replace URLs with a space in the message
27      # https://ihateregex.io/expr/url/
28      text = re.sub(r'https?:\/\/(www\.)?[-a-z0-9@:%._\+~#=]{1,256}\.[a-z0-9()]{1,6}\b([-a-z0-9()!@:%_\+.~#?&\/\/=]*)', ' ',
29
30      # Replace ticker symbols with a space. The ticker symbols are any stock symbol that starts with $.
31      text = re.sub(r'\$[a-z0-9]{1,6}', ' ', text)
32
33      # Replace StockTwits usernames with a space. The usernames are any word that starts with @.
34      text = re.sub(r'@[a-z0-9@._]+', ' ', text)
35
36      # Replace everything not a letter with a space
37      text = re.sub(r'[^a-z]+', ' ', text)
38
39      # Tokenize by splitting the string on whitespace into a list of words
40      tokens = text.split(' ')
41
42      # Lemmatize words using the WordNetLemmatizer. You can ignore any word that is not longer than one character.
43      wnl = nltk.stem.WordNetLemmatizer()
44      tokens = [wnl.lemmatize(token, pos='v') for token in tokens if len(token) > 1]  # root verbs
45      tokens = [wnl.lemmatize(token, pos='n') for token in tokens]  # root nouns
46
47      assert type(tokens) == list, 'Tokens should be list'
48      return tokens
```
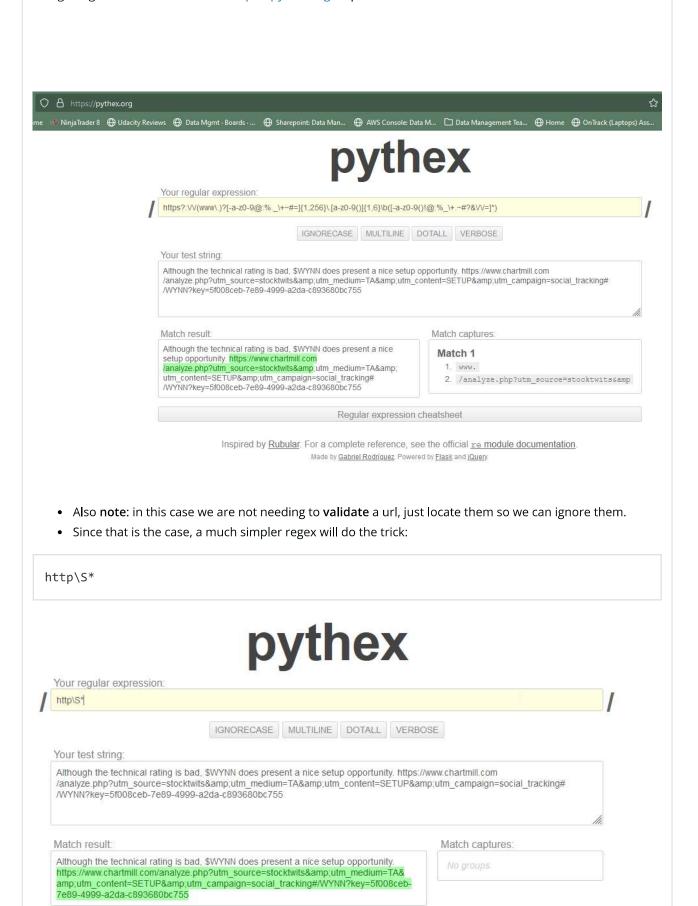
# Please Note:

- Your regular expression for urls is quite complex and in fact does not accurately capture ever part of every url.
- While this is not a reason for you to not pass this section, I figured you'd want to be aware of it . . .
- Shown below is one message for which your regex doesn't function properly. The tokens highlighted in red should have been removed:

```
Message: [69]
Although the technical rating is bad, $WYNN does present a nice setup opportunity. https://www.chartmill.com/analyze.php?
utm_source=stocktwits&amp;utm_medium=TA&amp;utm_content=SETUP&amp;utm_campaign=social_tracking#/WYNN?key=5f008ceb-7e89-49
99-a2da-c893680bc755
-----------------------------------------------------------------
['although', 'the', 'technical', 'rat', 'be', 'bad', 'do', 'present', 'nice', 'setup', 'opportunity', 'utm', 'medium', 't
a', 'amp', 'utm', 'content', 'setup', 'amp', 'utm', 'campaign', 'social', 'track', 'wynn', 'key', 'ceb', 'da', 'bc']
```

Using a regex tester like the one at https://pythex.org helps us see:



- Also **note**: in this case we are not needing to **validate** a url, just locate them so we can ignore them.
- Since that is the case, a much simpler regex will do the trick:

```
http\S*
```

Inspired by Rubular. For a complete reference, see the official re module documentation.

Made by Gabriel Rodríguez. Powered by Flask and jQuery.

---

Preprocess all the twits into the `tokenized` variable.

---

# Implemented correctly and the first few tokenized elements look appropriate!

- You are using your **preprocess** function and applying it to each message.
- Nice use of a list comprehension!

## Preprocess All the Twits

Now we can preprocess each of the twits in our dataset. Apply the function `preprocess` to all the twit messages.

```
1  # TODO Implement
2  tokenized = [preprocess(message) for message in messages]
```

---

Create a bag of words using the tokenized data.

---

# You have properly created a bag of words as you've been taught!

- You have used the **Counter** class properly to create a bag of words from all of the tokens!

## Bag of Words

Now with all of our messages tokenized, we want to create a vocabulary and count up how often each word appears in our entire corpus. Use the `Counter` function to count up all the tokens.

```
1  from collections import Counter
2
3
4  """
5  Create a vocabulary by using Bag of words
6  """
7
8  # TODO: Implement
9  bow = Counter([word for message in tokenized for word in message])
```

```
1  bow.most_common(10)
```

```
[('be', 582331),
 ('the', 398861),
 ('to', 379769),
 ('amp', 295583),
 ('for', 273575),
 ('utm', 270812),
 ('on', 241828),
 ('of', 211358),
 ('and', 208609),
 ('in', 205362)]
```

Remove most common and rare words by defining the following variables: `freqs`, `low_cotoff`, `high_cutoff`, `K_most_common`.

---

# You have removed the most common and rare words properly and left a good amount of data for your model to train with!!

- Good Choice for **high_cutoff**
- Good Choice for **low_cutoff**
- You are properly creating **K_most_common**
- You are properly creating **freqs**
- You are removing the rare and most common words in your filtered words list

## Well done!

```
1  """
2  Set the following variables:
3      freqs
4      low_cutoff
5      high_cutoff
6      K_most_common
7  """
8
9  # TODO Implement
10
11 # Dictionary that contains the Frequency of words appearing in messages.
12 # The key is the token and the value is the frequency of that word in the corpus.
13 freqs = {word: count/twits_total for word, count in bow.items()}
14 print(f'vocabulary: {len(freqs)}')
15
16 # Float that is the frequency cutoff. Drop words with a frequency that is lower or equal to this number.
17 low_cutoff = 0.000002
18
19 # Integer that is the cut off for most common words. Drop words that are the `high_cutoff` most common words.
20 high_cutoff = 20
21
22 # The k most common words in the corpus. Use `high_cutoff` as the k.
23 K_most_common = [word for word, _ in bow.most_common(high_cutoff)]
24 print(K_most_common)
25
26 filtered_words = [word for word in freqs if (freqs[word] > low_cutoff and word not in K_most_common)]
27 print(f'filtered words: {len(filtered_words)}')
```

```
vocabulary: 92727
['be', 'the', 'to', 'amp', 'for', 'utm', 'on', 'of', 'and', 'in', 'this', 'it', 'at', 'will', 'up', 'buy', 'report', 'go', 'have', 'metric']
filtered words: 27474
```

---

```
1  freqs = {word: count / twits_total for word, count in bow.items()}
2  max_freq = 0.0
3  min_freq = 1.0
4  target_low_cutoff=0.000002
5  max_word = ''
6  min_word = ''
7  remove_word_count = 0
8  for word, freq in freqs.items():
9      if freq <= target_low_cutoff:
10         remove_word_count += 1
11     if freq > max_freq:
12         max_freq = freq
13         max_word = word
14     if freq < min_freq:
15         min_freq = freq
16         min_word = word
17 print(f"The word with Maximum Frequency is: '{max_word}' - Frequency: {max_freq}")
18 print(f"The word with Minimum Frequency is: '{min_word}' - Frequency: {min_freq}")
19 print(f" With a low cutoff of {target_low_cutoff}, we would remove {remove_word_count} words out of {len(freqs)}")
20 print(f" That is {100 * (remove_word_count / (len(freqs)))}% of our data . . .")
21
```

```
The word with Maximum Frequency is: 'be' - Frequency: 0.3761803864316122
The word with Minimum Frequency is: 'raisinf' - Frequency: 6.459906589750712e-07
With a low cutoff of 2e-06, we would remove 65233 words out of 92727
That is 70.349520635845011% of our data . . .
```

Defining the variables : 'vacab', 'id2vocab' and 'filtered' correctly.

# vocab, id2vocab, and filtered are all properly created – Bravo!!

- you have properly created **vocab** as a dictionary where the key is the **word** and the value is the **id** of the word.
- you have properly created **id2vocab** as a dictionary where the key is the **id** of the word and the value is the **word**.
- you have properly created **filtered** as a list of **only** the words in vocab from the original **tokenized list**.

## Additionally, if you take any word and call id2vocab[vocab[word]] you get back the original word.

- Well done!

**Updating Vocabulary by Removing Filtered Words**

Let's creat three variables that will help with our vocabulary.

```
1  """
2  Set the following variables:
3      vocab
4      id2vocab
5      filtered
6  """
7
8  #TODO Implement
9
10 # A dictionary for the `filtered_words`. The key is the word and value is an id that represents the word.
11 vocab = {word: index for index, word in enumerate(filtered_words)}
12
13 # Reverse of the `vocab` dictionary. The key is word id and value is the word.
14 id2vocab = {index: word for word, index in vocab.items()}
15
16 # tokenized with the words not in `filtered_words` removed.
17 filtered = [[word for word in message if word in vocab] for message in tokenized]
18
19 assert set(vocab.keys()) == set(id2vocab.values()), 'Check vocab and id2vocab dictionaries'
```

```
1  id2vocab[vocab['great']]
```

```
'great'
```

# Neural Network

The init function correctly initializes the following parameters: `self.vocab_size`, `self.embed_size`, `self.lstm_size`, `self.lstm_layers`, `self.dropout`, `self.embedding`, `self.lstm`, and `self.fc`.

## Nicely done:

- you are correctly initializing **self.vocab_size**
- you are correctly initializing **self.embed_size**
- you are correctly initializing **self.lstm_size**

- you are correctly initializing **self.lstm_layers**
- you are correctly initializing **self.dropout**
- you are correctly initializing **self.embedding**

```python
3    def __init__(self, vocab_size, embed_size, lstm_size, output_size, lstm_layers=1, dropout=0.1):
4        """
5        Initialize the model by setting up the layers.
6
7        Parameters
8        ----------
9            vocab_size : The vocabulary size.
10           embed_size : The embedding layer size.
11           lstm_size : The LSTM layer size.
12           output_size : The output size.
13           lstm_layers : The number of LSTM layers.
14           dropout : The dropout probability.
15       """
16
17       super().__init__()
18
19       self.vocab_size = vocab_size
20       self.embed_size = embed_size
21       self.lstm_size = lstm_size
22       self.output_size = output_size
23       self.lstm_layers = lstm_layers
24
25       # TODO Implement
26
27       # Setup embedding layer
28       self.embedding = nn.Embedding(vocab_size, embed_size)
29
30       # Setup LSTM layer
31       self.lstm = nn.LSTM(embed_size, lstm_size, lstm_layers, dropout=dropout, batch_first=False)
32
33       # Setup dropout layer
34       self.dropout = nn.Dropout(dropout)
35
36       # Setup fully connected linear (dense) layer
37       self.fc = nn.Linear(lstm_size, output_size)
38
39       # Setup softmax layer
40       self.lsoftmax = nn.LogSoftmax(dim=1)
41
```

The 'init_hidden' function generates a hidden state

## Yes! - you are properly generating a hidden_state using the incoming batch_size.

```
42      def init_hidden(self, batch_size):
43          """
44          Initializes hidden state
45
46          Parameters
47          ----------
48              batch_size : The size of batches.
49
50          Returns
51          -------
52              hidden_state
53
54          """
55
56          # TODO Implement
57
58          # Create two new tensors with sizes lstm_layers x batch_size x lstm_size,
59          # initialized to zero, for hidden state and cell state of LSTM
60          weight = next(self.parameters()).data
61
62          hidden_state = (weight.new(self.lstm_layers, batch_size, self.lstm_size).zero_(),
63                          weight.new(self.lstm_layers, batch_size, self.lstm_size).zero_())
64
65          return hidden_state
66
```

The 'forward' function performs a forward pass of the model the parameter input using the hidden state.

# Excellent! - you are performing a forward pass of the model on the incoming nn_input and using the incoming hidden_state.

- Using an embedding layer helps because there are a large number of words in the vocabulary!

```
67      def forward(self, nn_input, hidden_state):
68          """
69          Perform a forward pass of our model on nn_input.
70
71          Parameters
72          ----------
73              nn_input : The batch of input to the NN.
74              hidden_state : The LSTM hidden state.
75
76          Returns
77          -------
78              logps: log softmax output
79              hidden_state: The new hidden state.
80
81          """
82
83          # TODO Implement
84          #batch_size = nn_input.size(0)   # since batch_first param is set to True
85
86          # Pass thru embedding layer
87          embeds = self.embedding(nn_input)
88          #print(f'embeds shape: {embeds.shape}')
89
90          # Pass thru LSTM layer
91          lstm_out, hidden_state = self.lstm(embeds, hidden_state)
92          #print(f'lstm_out shape: {lstm_out.shape}')
93          #print(f'hidden_state [0] shape: {hidden_state[0].shape}')
94
95          # Stack up LSTM outputs
96          #lstm_stack = lstm_out.contiguous().view(-1, self.lstm_size)
97          lstm_stack = lstm_out[-1, :, :]
98          #print(f'lstm_stack shape: {lstm_stack.shape}')
99
100         # Pass thru dropout layer
101         out = self.dropout(lstm_stack)
102
103         # Pass thru dense layer
104         out = self.fc(out)
105         #print(f'out shape: {out.shape}')
106
107         # Pass thru log softmax layer
108         log_ps = self.lsoftmax(out)
109         #print(f'log ps shape: {log_ps.shape}')
109         #print(f'log_ps shape: {log_ps.shape}')
110
111         # Reshape to be batch_size first
112         #log_ps = log_ps.view(batch_size, -1)
113
114         # Get last batch of output
115         #log_ps = log_ps[:, -1]
116
117         return log_ps, hidden_state
```

# Training

Correctly split the data into `train_features` , `valid_features` , `train_labels` , and `valid_labels` .

## Excellent!

- You have created a nice **80 / 20** split of training and validation data!

```
1  """
2  Split data into training and validation datasets. Use an appropriate split size.
3  The features are the `token_ids` and the labels are the `sentiments`.
4  """
5
6  # TODO Implement
7
8  split_frac = 0.8
9  split_pos = int(len(token_ids) * split_frac)
10
11 train_features, valid_features = token_ids[:split_pos], token_ids[split_pos:]
12 train_labels, valid_labels = sentiments[:split_pos], sentiments[split_pos:]
```

```
1  print("\t\t\tFeature Sizes:")
2  print("Train size: \t\t{}".format(len(train_features)),
3        "\nValidation size: \t{}".format(len(valid_features)))
4
```

```
                      Feature Sizes:
Train size:           823349
Validation size:      205838
```

Train your model with dropout and clip the gradient. Print out the training progress with the loss and accuracy.

## Almost there!

- You are calling **optimizer.zerograd()** properly.
- You are getting output from the model.
- You are using **criterion** to calculate the loss and calling **backward()** to perform backpropagation.
- You are using **clip_grad_norm** to prevent exploding gradient problem.
- You are calculating and printing **training loss and validation loss for every 100 steps**

Excellent job so far . . .

## However:

- You were asked to calculate and display validation accuracy as well for every 100 steps

You may find this post at Udacity's Knowledge Base
https://knowledge.udacity.com/questions/30981 quite helpful!!

```python
1    """
2    Train your model with dropout. Make sure to clip your gradients.
3    Print the training loss, validation loss, and validation accuracy for every 100 steps.
4    """
5
6    # Loss and optimazation
7    learning_rate = 0.001
8
9    criterion = nn.NLLLoss()
10   optimizer = optim.Adam(model.parameters(), lr=learning_rate)
11
12   # Training parameters
13   batch_size = 1024
14   sequence_length = 100
15
16   epochs = 3
17   clip = 5 # gradient clipping
18
19   # Printing parameters
20   print_every = 100
21
22   # Turn on training mode
23   model.train()
24
25   # Epoch loop
26   for e in range(epochs):
27       steps = 0
28       print('Starting epoch {}/{}'.format(e + 1, epochs))
29
30       # Batch loop
31       for text_batch, labels in dataloader(train_features, train_labels, sequence_length=sequence_length, batch_size=batch_s
32           # TODO Implement: Train Model
33           steps += 1
34
35           # Initialize hidden state
36           hidden = model.init_hidden(labels.shape[0])  # at the last iteration of the epoch, rows of the batch will be left-
37
38           # Set Device
39           text_batch, labels = text_batch.to(device), labels.to(device)
40           for each in hidden:
41               each.to(device)
42
43           # Zero accumulated gradients
44           model.zero_grad()
45
46           # Get the output from the model
47           logps, hidden = model(text_batch, hidden)
48
49           # calculate the loss and perform backprop
50           loss = criterion(logps.squeeze(), labels)
51           loss.backward()
52
53           # `clip_grad_norm` helps prevent the exploding gradient problem in RNNs / LSTMs.
54           nn.utils.clip_grad_norm_(model.parameters(), clip)
55           optimizer.step()
56
57           # Loss stats
58           if steps % print_every == 0:
59               model.eval()
60
61               # TODO Implement: Print metrics
62               valid_losses = []
63
64               for text_batch, labels in dataloader(valid_features, valid_labels, sequence_length=sequence_length, batch_size
65                   valid_hidden = model.init_hidden(labels.shape[0])  # at the last iteration of the epoch, rows of the batch
66
67                   text_batch, labels = text_batch.to(device), labels.to(device)
68                   for each in valid_hidden:
69                       each.to(device)
70
71                   logps, valid_hidden = model(text_batch, valid_hidden)
72                   valid_loss = criterion(logps.squeeze(), labels)
73
74                   valid_losses.append(valid_loss.item())
75
76               print("Epoch: {}/{}...".format(e + 1, epochs),
77                   "Step: {}...".format(steps),
78                   "Loss: {:.6f}...".format(loss.item()),
79                   "Valid Loss: {:.6f}".format(np.mean(valid_losses)))
80
81               model.train()
```

```
Starting epoch 1/3
Epoch: 1/3... Step: 100... Loss: 0.942924... Valid Loss: 1.016169
Epoch: 1/3... Step: 200... Loss: 0.886049... Valid Loss: 0.934194
Epoch: 1/3... Step: 300... Loss: 0.867592... Valid Loss: 0.902369
Epoch: 1/3... Step: 400... Loss: 0.838153... Valid Loss: 0.883949
Epoch: 1/3... Step: 500... Loss: 0.853582... Valid Loss: 0.865797
Epoch: 1/3... Step: 600... Loss: 0.813141... Valid Loss: 0.856329
Epoch: 1/3... Step: 700... Loss: 0.825543... Valid Loss: 0.845845
Epoch: 1/3... Step: 800... Loss: 0.754917... Valid Loss: 0.838022
Starting epoch 2/3
Epoch: 2/3... Step: 100... Loss: 0.753580... Valid Loss: 0.839457
Epoch: 2/3... Step: 200... Loss: 0.776417... Valid Loss: 0.836485
Epoch: 2/3... Step: 300... Loss: 0.731900... Valid Loss: 0.833061
Epoch: 2/3... Step: 400... Loss: 0.791676... Valid Loss: 0.826225
Epoch: 2/3... Step: 500... Loss: 0.812037... Valid Loss: 0.827216
Epoch: 2/3... Step: 600... Loss: 0.756848... Valid Loss: 0.818247
Epoch: 2/3... Step: 700... Loss: 0.709338... Valid Loss: 0.816720
Epoch: 2/3... Step: 800... Loss: 0.674014... Valid Loss: 0.816762
Starting epoch 3/3
Epoch: 3/3... Step: 100... Loss: 0.686234... Valid Loss: 0.833953
Epoch: 3/3... Step: 200... Loss: 0.727408... Valid Loss: 0.838917
Epoch: 3/3... Step: 300... Loss: 0.739165... Valid Loss: 0.836471
Epoch: 3/3... Step: 500... Loss: 0.718194... Valid Loss: 0.830025
Epoch: 3/3... Step: 600... Loss: 0.711602... Valid Loss: 0.833537
Epoch: 3/3... Step: 700... Loss: 0.768566... Valid Loss: 0.831728
Epoch: 3/3... Step: 800... Loss: 0.742700... Valid Loss: 0.827539
```

## Making Predictions

The `predict` function correctly prints out the prediction vector from the trained model.

### You have properly displayed an accurately created prediction vector from the trained model.

- Well done!

```
1  def predict(text, model, vocab):
2      """
3      Make a prediction on a single sentence.
4
5      Parameters
6      ----------
7          text : The string to make a prediction on.
8          model : The model to use for making the prediction.
9          vocab : Dictionary for word to word ids. The key is the word and the value is the word id.
10
11     Returns
12     -------
13         pred : Prediction vector
14     """
15
16     # TODO Implement
17
18     tokens = preprocess(text)
19
20     # Filter non-vocab words
21     tokens = [token for token in tokens if token in vocab]
22
23     # Convert words to ids
24     tokens =  [vocab[token] for token in tokens]
25
26     # Adding a batch dimension
27     text_input = torch.from_numpy(np.asarray(torch.LongTensor(tokens).view(-1, 1)))
28
29     # Get the NN output
30     hidden = model.init_hidden(1)
31     logps, _ = model.forward(text_input, hidden)
32
33     # Take the exponent of the NN output to get a range of 0 to 1 for each label.
34     pred = torch.exp(logps)
35
36     return pred
```

```
1  text = "Google is working on self driving cars, I'm bullish on $goog"
2
3  model.eval()
4  model.to("cpu")
5
6  predict(text, model, vocab)
```

tensor([[ 0.0000,  0.0071,  0.0065,  0.8282,  0.1582]])

Answer what the prediction of the model is and the uncertainty of the prediction.

# Excellent Answer based on your displayed prediction vector!

- You have indicated which label the model is most confident with
- You have indicated the **probability** of that label
- You have indicated the **uncertaintly**

## Well Done

```
1  text = "Google is working on self driving cars, I'm bullish on $goog"
2
3  model.eval()
4  model.to("cpu")
5
6  predict(text, model, vocab)
```

```
tensor([[ 0.0000,   0.0071,   0.0065,   0.8282,   0.1582]])
```

**Questions: What is the prediction of the model? What is the uncertainty of the prediction?**

**TODO: Answer Question**

The prediction of the model gave highest positive sentiment on the twit 0.8282 (83%). Models always have uncertainty, in this case, this model has 17% uncertainty.

☑ RESUBMIT

⤓ DOWNLOAD PROJECT

## Best practices for your project resubmission

Ben shares 5 helpful tips to get you through revising and resubmitting your project.

▶ Watch Video (3:01)

RETURN TO PATH