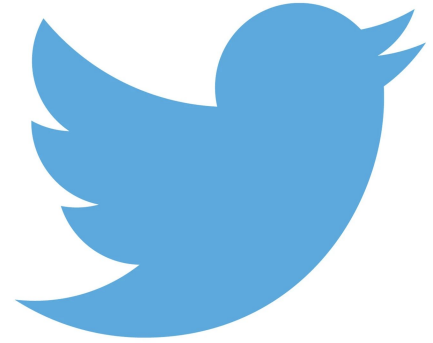


Twitter Sentiment Analysis Using an LSTM RNN

By: Zack Godwin

Overview

Brand management is a key factor for a successful business. Understanding how customers feel about a brand, company, or even a specific topic can help predict how they'll respond to future decisions. Twitter gives us the perfect data stream to investigate this.





Project Goals

- Stream at least 200,000 tweets from Twitter using the Twitter API.
 - ◆ Store data in a Postgres database.
- Use the TextBlob sentiment analysis module to **label** the tweets.
- Using the labelled data set, build a classification model to predict sentiment.
- Train a RNN classifier and compare performance.
 - ◆ Investigate the impact of different degrees of text preprocessing on the performance of the RNN.

Data Capture Methodology

- Using the [Twitter API](#) I wrote a python script to collect Tweets and deposit the data into my Postgres database.
- My Tweet filter was programmed to grab anything with the word “[trump](#)”.
- I collected over 400,000 tweets.

```
1 import tweepy
2 from tweepy import Stream
3 from tweepy import OAuthHandler
4 from tweepy.streaming import StreamListener
5 import time
6 import json
```

```
1 import psycopg2
2
3 # Connect using psycopg2
4 conn = psycopg2.connect(dbname= "Final Project", user= "postgres",
5                          password= "", host="127.0.0.1",
6                          port= "5433")
7
8 curs = conn.cursor()
```

```
1 ckey = ''
2 csecret = ''
3 atoken = ''
4 asecret = ''
```

```
1 class listener(StreamListener):
2
3     def on_data(self, data):
4         try:
5             all_data = json.loads(data)
6             tweet = all_data["text"]
7             username = all_data["user"]["screen_name"]
8             location = all_data["user"]["location"]
9
10            curs.execute("INSERT INTO trump (time, location, username, tweet)
11                          (time.time(), location, username, tweet)")
12
13            conn.commit()
14
15            print((username,tweet))
16
17            return True
18        except KeyError:
19            pass
20
21     def on_error(self, status):
22         if status_code == 420:
23             #returning False in on_data disconnects the stream
24             return False
25
26
27 auth = tweepy.OAuthHandler(ckey, csecret)
28 auth.set_access_token(atoken, asecret)
29 twitterStream = Stream(auth, listener())
30 twitterStream.filter(track=['trump'], languages=['en'])
```

The Data

- Dataframe Dimension:
 - **430228** Rows
- Mean Word Count:
 - **18.7 Words** (+/- 6.1)
- Mean Character Count:
 - **123 Characters** (+/- 33.7)

| | tweet | word_count | character_count |
|---|---|------------|-----------------|
| 0 | 'RT @SethAbramson: In case you missed it: what... | 25 | 155 |
| 1 | 'RT @RonWyden: Incredible. More luxury travel ... | 22 | 149 |
| 2 | 'President Trump Directed Michael Cohen To Lie... | 16 | 125 |
| 3 | 'RT @TeaPainUSA: Tea would wager that Trump in... | 26 | 147 |
| 4 | 'RT @BruceBartlett: There is one person in Ame... | 24 | 147 |

| | word_count | character_count |
|-------|---------------|-----------------|
| count | 430228.000000 | 430228.000000 |
| mean | 18.674238 | 122.972138 |
| std | 6.134861 | 33.739280 |
| min | 1.000000 | 1.000000 |
| 25% | 16.000000 | 122.000000 |
| 50% | 20.000000 | 140.000000 |
| 75% | 23.000000 | 140.000000 |
| max | 37.000000 | 258.000000 |

Data Cleaning

I compared two methods of data cleaning, the principal difference being the presence of stop-words.

| filtered_tweet | clean |
|---|---|
| case missed fuss tonight retweet think followe... | case you missed what all the fuss about ton... |
| incredible luxury travel trump administration ... | incredible more luxury travel from the trump a... |
| president trump directed michael cohen lie con... | president trump directed michael cohen lie c... |
| tea wager trump instructed folks lie congress ... | tea would wager that trump instructed all his ... |
| person america trump chose rupert murdoch medi... | there one person america who could somethin... |

```
#Second Cleaning Function that keeps stopwords
def extract_text(text):

    # Convert to string
    text = text.astype(str)

    # Remove URLs
    text = text.str.replace('https?://[A-Za-z0-9./]+','')

    # Keep Hashtag text
    text = text.str.replace("[^a-zA-Z]", " ")

    # Make lowercase
    text = text.apply(lambda x: " ".join(x.lower() for x in x.split()))

    # Remove whitespaces
    text = text.apply(lambda x: " ".join(x.strip() for x in x.split()))

    # Remove special characters
    text = text.apply(lambda x: "".join(
        [" " if ord(i) < 32 or ord(i) > 126 else i for i in x]))

    # Remove punctuation
    text = text.str.replace('[^\w\s]', '')

    # Remove numbers
    text = text.str.replace('\d+', '')

    #Remove 1-2 letter clutter remnants
    text = text.apply(lambda x: re.sub(r'\b\w{1,2}\b', '', x))

    return text

# remove RT: @user
df['clean'] = np.vectorize(remove_pattern)(df['tweet'], "RT @[\\w]*")
```

Sentiment Labeling with TextBlob

The sentiment property of Textblob returns a namedtuple of the form Sentiment(polarity, subjectivity). The polarity score is a float within the range [-1.0, 1.0]. We will use this to label tweets as positive, neutral, or negative.

```
def analyze_sentiment(tweet):  
    '''  
    Utility function to classify the polarity of a tweet  
    using textblob.  
    '''  
    analysis = TextBlob(tweet)  
    if analysis.sentiment.polarity > 0:  
        return 1  
    elif analysis.sentiment.polarity == 0:  
        return 0  
    else:  
        return -1  
  
# Create a column with the result of the analysis:  
df['SA'] = np.array(  
    [analyze_sentiment(tweet) for tweet in df['filtered_tweet']])  
  
df['SA2'] = np.array(  
    [analyze_sentiment(tweet) for tweet in df['clean']])  
  
df.head()
```

Sentiment Labeling Discrepancy

It appears TextBlob performs better on text containing stopwords.

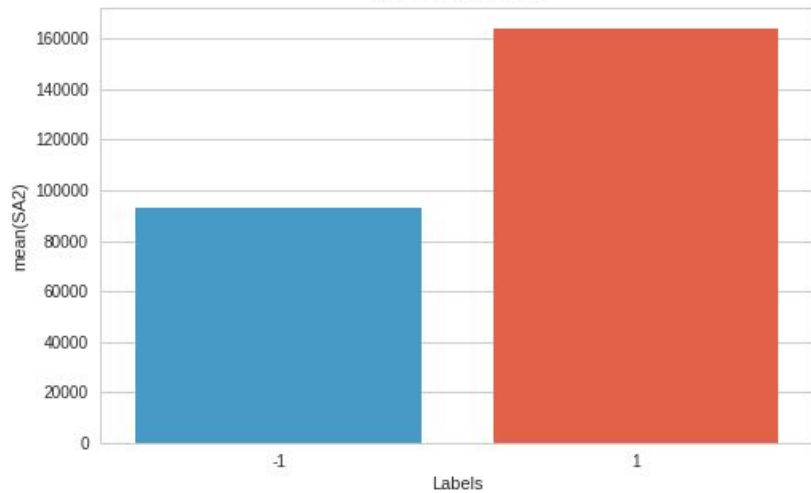
As such, I used 'clean' and SA2 as my data and target labels (respectively).

“RT @TeaPainUSA: Tea would wager that Trump instructed all his folks to lie to Congress because he knew Nunes and his other GOP imps and dem\u2026”

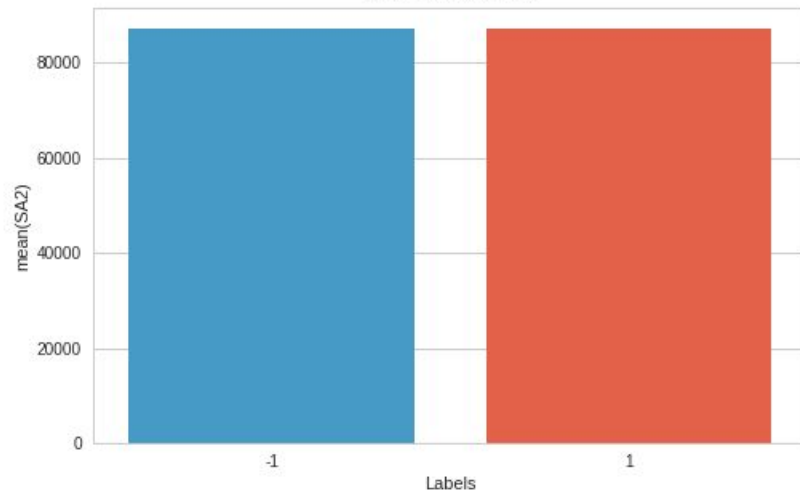
| tweet | word_count | character_count | filtered_tweet | clean | SA | SA2 |
|--|------------|-----------------|---|---|----|-----|
| 'RT @TeaPainUSA: Tea would wager that Trump in... | 26 | 147 | tea wager trump instructed folks lie congress ... | tea would wager that trump instructed all his ... | 0 | -1 |

EDA: Class Balance

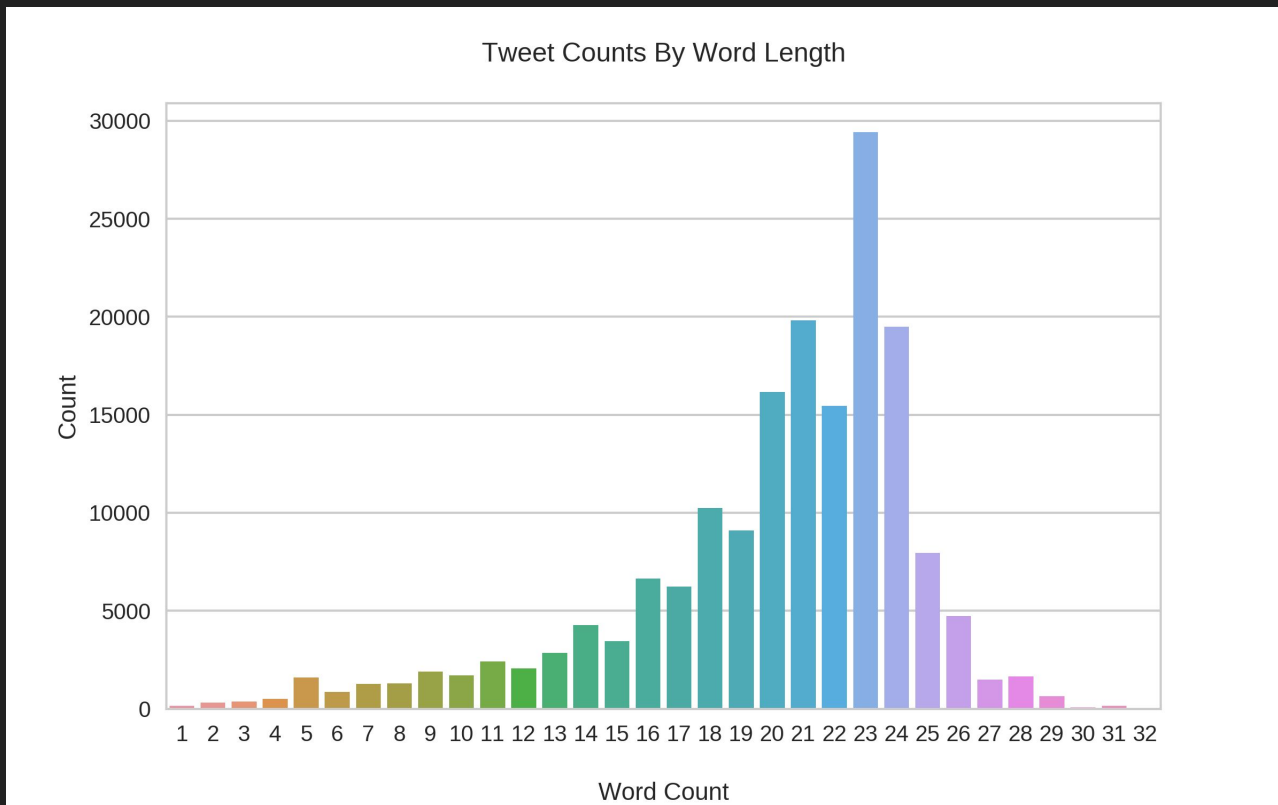
DF2 Class Balance



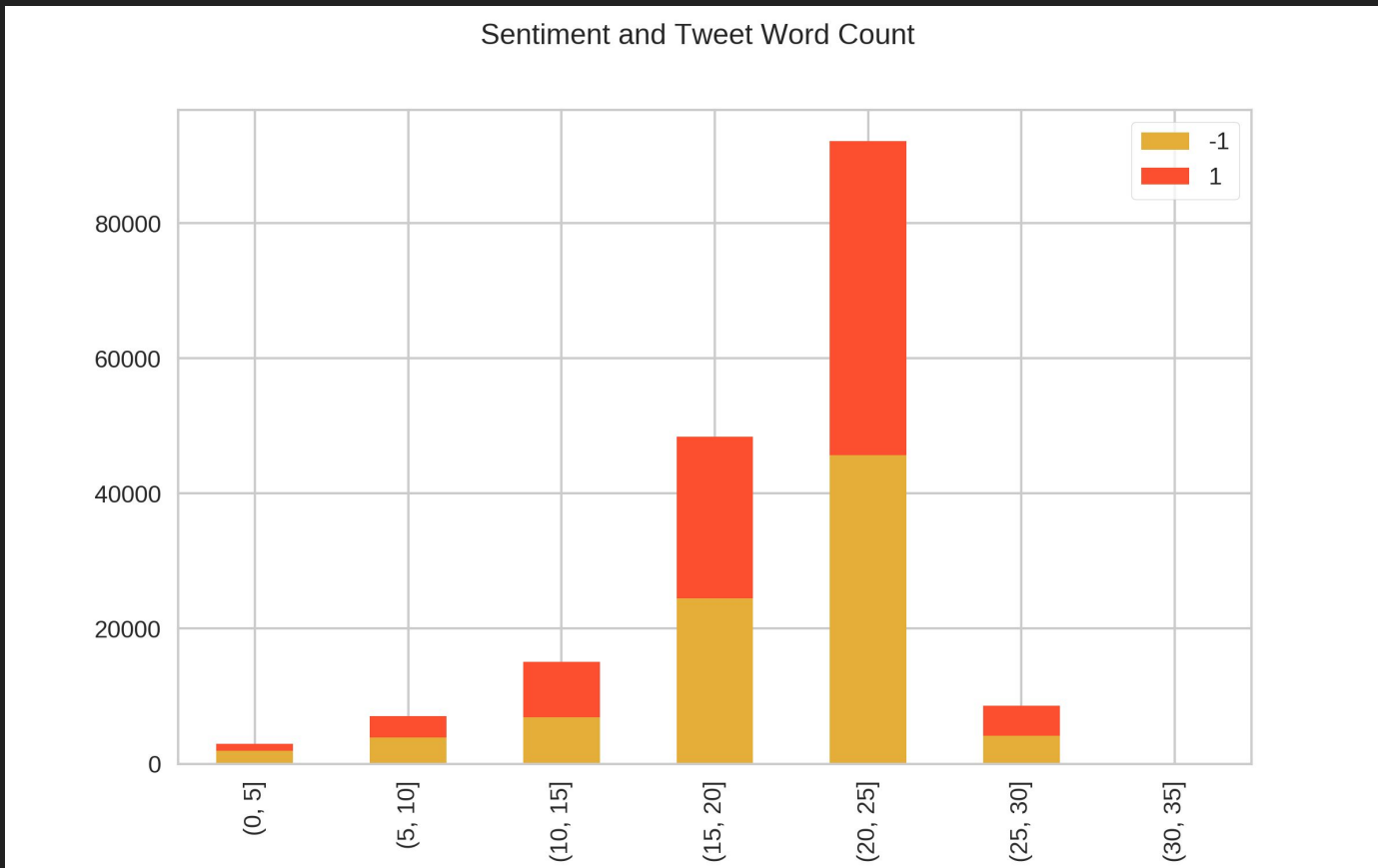
DF3 Class Balance



EDA:



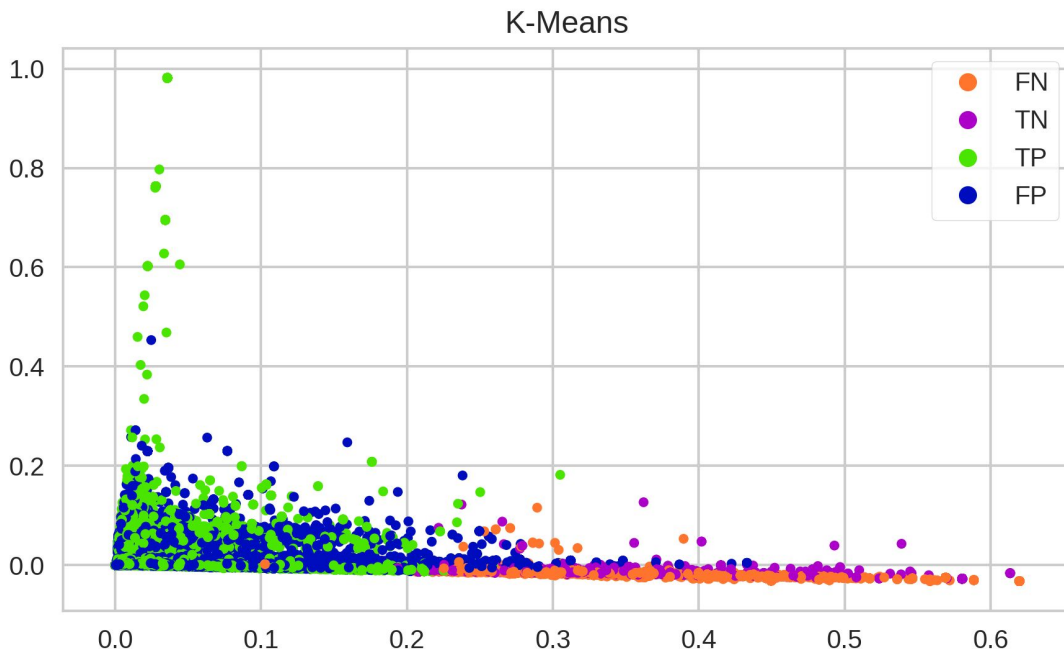
EDA:



EDA: K-Means Cluster Analysis

Using TF-IDF I vectorized the cleaned text.

Next I reduced the dimensionality of the data down to 600 principal components using Truncated SVD.



EDA: Word2Vec

Word2vec is a two-layer neural net that processes text.

The input is a text corpus (clean tweets) and its output is a set of vectors: **feature vectors for words in that corpus.**

```
from gensim.models import Word2Vec

w2v_model = Word2Vec(
    sentences=sentences, size=300, window=5, min_count=5, workers=4, sg=0)

w2v_model.wv.most_similar('bad')

[('busy', 0.42365583777427673),
 ('good', 0.4175521433353424),
 ('aware', 0.408443421125412),
 ('funny', 0.39820361137390137),
 ('irresponsible', 0.38138046860694885),
 ('sorry', 0.3618493378162384),
 ('badly', 0.35371482372283936),
 ('expensive', 0.34147143363952637),
 ('perfect', 0.33978307247161865),
 ('dead', 0.33889153599739075)]
```

EDA: FastText

FastText is an extension to Word2Vec proposed by Facebook in 2016.

Instead of feeding individual words into the Neural Network, **FastText breaks words into several n-grams** (sub-words).

```
from gensim.models import FastText
fast_model = FastText(
    sentences, size=300, window=5, min_count=5, workers=4, sg=0)

fast_model.wv.most_similar("bad")

[('badly', 0.6813015341758728),
 ('badass', 0.6336409449577332),
 ('vlad', 0.6131807565689087),
 ('load', 0.5794876217842102),
 ('dad', 0.5165011882781982),
 ('knead', 0.503669023513794),
 ('glad', 0.49534744024276733),
 ('bath', 0.48387226462364197),
 ('bags', 0.47988539934158325),
 ('brad', 0.46785059571266174)]
```

EDA: Doc2Vec

Doc2Vec is a small extension to the CBOW Word2Vec model.

Instead of using just words to predict the next word, **we also add another feature vector, which is document-unique.**

```
doc2vec_model.wv.most_similar("bad").  
  
[('all', 0.617525041103363),  
 ('before', 0.5723784565925598),  
 ('funny', 0.5426432490348816),  
 ('travisallen', 0.4876490533351898),  
 ('sure', 0.48046356439590454),  
 ('brave', 0.4789738953113556),  
 ('dead', 0.47510266304016113),  
 ('what', 0.47376853227615356),  
 ('lunatic', 0.472905695438385),  
 ('when', 0.462916761636734)]
```

Pytorch: LSTM RNN

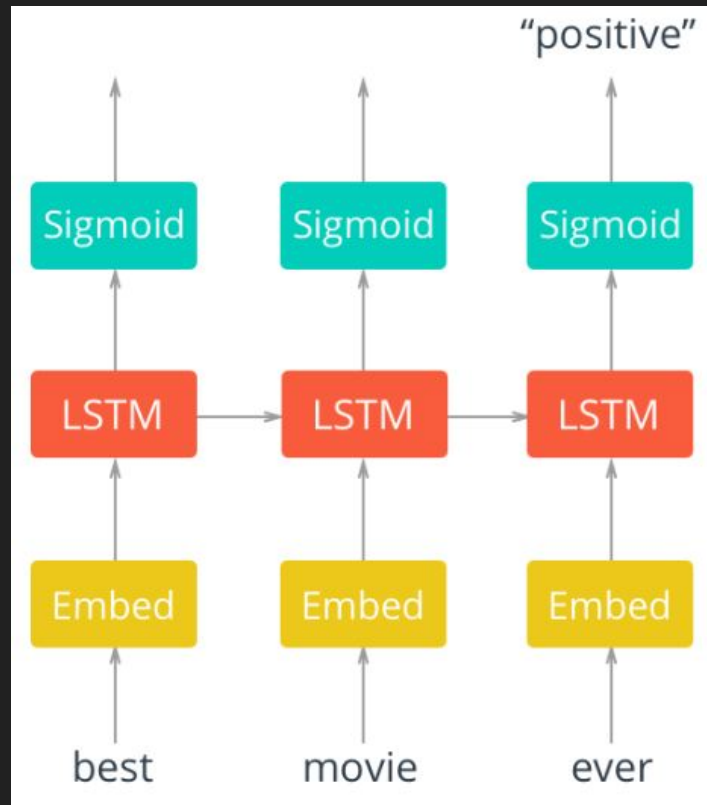
Long Short-Term Memory (LSTM) recurrent neural networks manage to keep contextual information of inputs by integrating a loop that allows information to flow from one step to the next.

Since tweets are composed from a sequence of words, and the specific order of those words provide context to the sentiment of the tweet, it stands to reason [this Neural Network should predict tweet sentiment well.](#)

Pytorch: LSTM RNN

The layers are as follows:

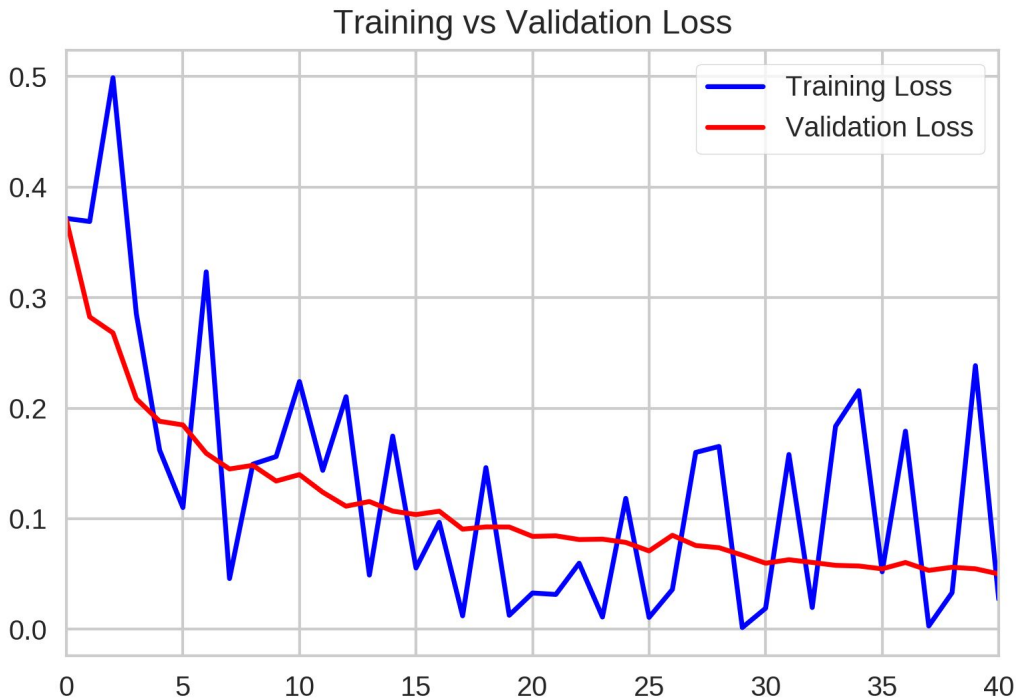
1. An embedding layer that converts our word tokens (integers) into embeddings of a specific size.
2. An LSTM layer defined by a hidden_state size and number of layers
3. A fully-connected output layer that maps the LSTM layer outputs to a desired output_size
4. A sigmoid activation layer which turns all outputs into a value 0-1; return only the last sigmoid output as the output of this network.



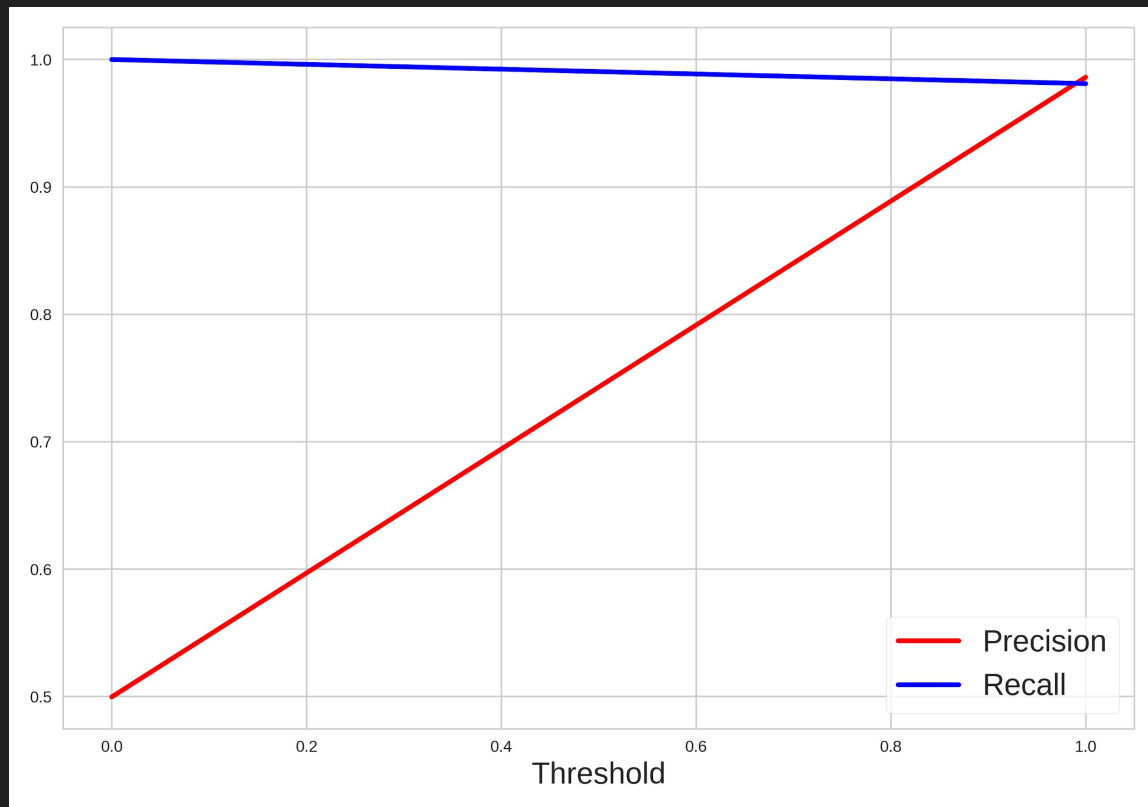
Pytorch: Word Embeddings

1. Convert column of cleaned tweets to a giant list of words (words), and also a list of tweets (tweet_split).
2. Build dictionary to pair words to integers.
3. Use the dictionary to tokenize each tweet in tweet_split
 - susansarandon nytimes what wrong with you why are you doing this again you must closet trump supporter wtf
4. Store the tokenized tweet in list (tweet_ints).
 - [3157, 739, 34, 135, 18, 8, 86, 11, 8, 437, 6, 145, 8, 221, 7626, 2, 991, 519]
5. Pad the tokenized tweets with zeros so all tokenized tweets are equal dimensions.
 - [0, 0, 0, 0, 0, 0, 0, 0, 3157, 739, 34, 135, 18, 8, 86, 11, 8, 437, 6, 145, 8, 221, 7626, 2, 991, 519]

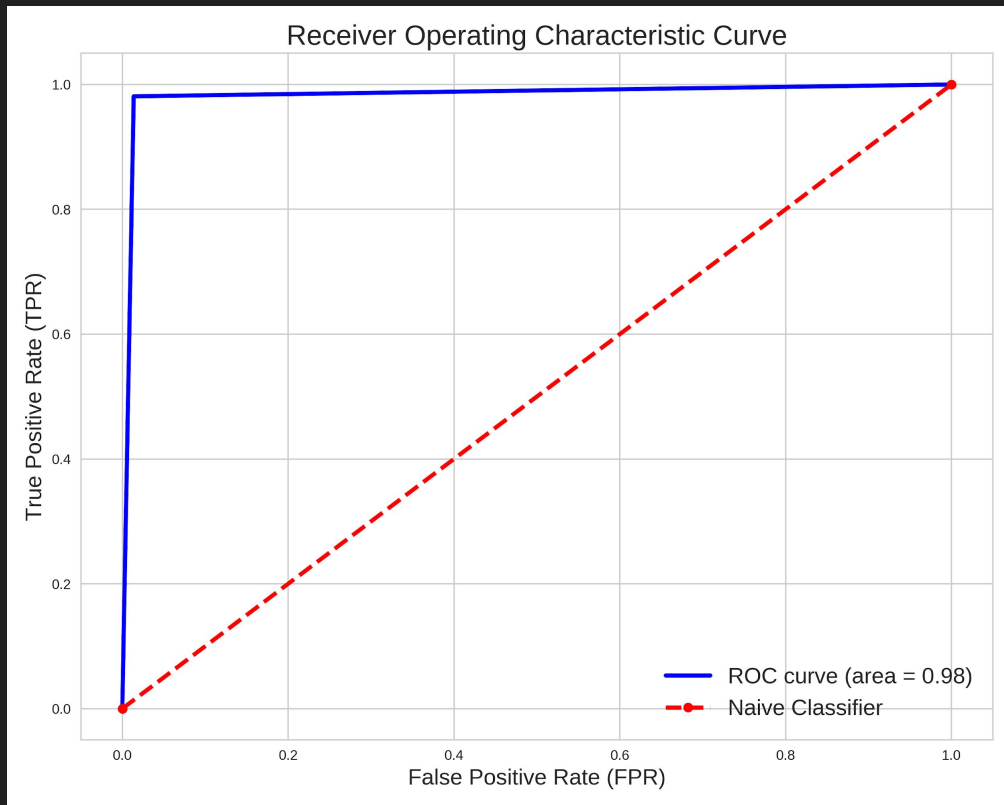
Pytorch: Model Training



Pytorch: Model Validation



Pytorch: Model Validation



Live Demo

Questions?



Extra Slides