

# LSTM classifier for abusive/sarcastic language

## Import libraries

```
In [3]: import pandas as pd
import numpy as np

from pymongo import MongoClient
from Preprocessing import config
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import CountVectorizer
from keras.preprocessing.text import Tokenizer
from keras.preprocessing.sequence import pad_sequences
from keras.models import Sequential
from keras.layers import Dense, Embedding, LSTM, SpatialDropout1D
from keras.utils.np_utils import to_categorical
```

```
c:\users\evaarevalo\appdata\local\programs\python\python36\lib\site-packages
\h5py\__init__.py:36: FutureWarning: Conversion of the second argument of iss
ubdtype from `float` to `np.floating` is deprecated. In future, it will be tr
eated as `np.float64 == np.dtype(float).type`.
```

```
    from ._conv import register_converters as _register_converters
Using TensorFlow backend.
```

## MongoDB connection

```
In [4]: client = MongoClient(config.MONGODB['hostname'], config.MONGODB['port'])
db = client[config.MONGODB['db']]
collection = db[config.MONGODB['collection_news_and_sarcasm']]
```

## Get datasets

```
In [5]: def getDatasetsFromMongoDB():
        ''' mongodb to pandas dataframe, export to csv and return'''
        results=collection.find()
        #strip and reshuffle
        df = pd.DataFrame(list(results))
        df=df[['label','text']]
        df=df.reindex(np.random.permutation(df.index))
        df.to_csv('sarcasm_and_news_dataset.csv',encoding='utf-8-sig')
        return df
```

```
In [13]: def getDatasetsFromCsv():  
        '''import csv, reshuffle and return it'''  
        df=pd.read_csv('sarcasm_and_news_dataset.csv')  
        df=df.reindex(np.random.permutation(df.index))  
        return df[['label','text']]
```

```
In [14]: data=getDatasetsFromCsv()  
        train,test = train_test_split(data,test_size=0.1)
```

## Embedding layer parameters

```
In [27]: max_features=2000  
        embedding_dim=128  
        lstm_out_dim=196  
        batch_size=32  
        epochs=10
```

## Tokenize

```
In [28]: tokenizer = Tokenizer(num_words=max_features, split=' ')  
        tokenizer.fit_on_texts(data['text'].values)  
        X = tokenizer.texts_to_sequences(data['text'].values)  
        X = pad_sequences(X)
```

## Build model

```
In [29]: model = Sequential()
model.add(Embedding(max_features, embedding_dim,input_length=X.shape[1]))
model.add(SpatialDropout1D(0.4))
model.add(LSTM(lstm_out_dim, dropout=0.2,recurrent_dropout=0.2))
model.add(Dense(2,activation='softmax'))
model.compile(loss='categorical_crossentropy', optimizer='adam',
              metrics=['accuracy'])
print(model.summary())
```

Layer (type)	Output Shape	Param #
embedding_3 (Embedding)	(None, 41, 128)	256000
spatial_dropout1d_3 (Spatial	(None, 41, 128)	0
lstm_3 (LSTM)	(None, 196)	254800
dense_3 (Dense)	(None, 2)	394
Total params: 511,194		
Trainable params: 511,194		
Non-trainable params: 0		
None		

## Train

```

In [30]: Y= pd.get_dummies(data['label']).values
X_train, X_test, Y_train, Y_test = train_test_split(X,Y, test_size=0.33,
    random_state=42)
#Get shapes
print(X_train.shape,Y_train.shape)
print(X_test.shape,Y_test.shape)
model.fit(X_train, Y_train, epochs = epochs, batch_size=batch_size, verbose=2)

(17592, 41) (17592, 2)
(8665, 41) (8665, 2)
Epoch 1/10
  - 64s - loss: 0.1196 - acc: 0.9589
Epoch 2/10
  - 64s - loss: 0.0516 - acc: 0.9805
Epoch 3/10
  - 68s - loss: 0.0339 - acc: 0.9877
Epoch 4/10
  - 60s - loss: 0.0277 - acc: 0.9907
Epoch 5/10
  - 55s - loss: 0.0207 - acc: 0.9931
Epoch 6/10
  - 58s - loss: 0.0167 - acc: 0.9943
Epoch 7/10
  - 55s - loss: 0.0150 - acc: 0.9948
Epoch 8/10
  - 53s - loss: 0.0125 - acc: 0.9964
Epoch 9/10
  - 54s - loss: 0.0094 - acc: 0.9970
Epoch 10/10
  - 59s - loss: 0.0107 - acc: 0.9965

Out[30]: <keras.callbacks.History at 0x1c1686df588>

```

## Validate

```

In [31]: validation_size = 1500

X_validate = X_test[-validation_size:]
Y_validate = Y_test[-validation_size:]
X_test = X_test[:-validation_size]
Y_test = Y_test[:-validation_size]
score,acc = model.evaluate(X_test, Y_test, verbose = 2, batch_size = batch_size)

print("score: %.2f" % (score))
print("acc: %.2f" % (acc))

score: 0.08
acc: 0.98

```

```
In [32]: positive_count=0
negative_count=0
positive_correct=0
negative_correct = 0

for x in range(len(X_validate)):

    result = model.predict(X_validate[x].reshape(1,X_test.shape[1]),batch_size=1,verbose = 2)[0]

    if np.argmax(result) == np.argmax(Y_validate[x]):
        if np.argmax(Y_validate[x]) == 0:
            negative_correct += 1
        else:
            positive_correct += 1

    if np.argmax(Y_validate[x]) == 0:
        negative_count += 1
    else:
        positive_count += 1

print("Positive Accuracy", positive_correct/positive_count*100, "%")
print("Negative Accuracy", negative_correct/negative_count*100, "%")
```

Positive Accuracy 96.61016949152543 %

Negative Accuracy 98.82491186839013 %