Google's strongest NLP model: BERT

https://github.com/google-research/bert (https://github.com/google-research/bert)

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Nov 08, 2018	BERT (single model) Google Al Language	80.005	83.061
1 Nov 16, 2018	Candi-Net+BERT (single model) 42Maru NLP Team	80.106	82.862
2 Nov 09, 2018	L6Net + BERT (single model) Layer 6 Al	79.181	82.259
3 Nov 06, 2018	SLQA+BERT (single model) Alibaba DAMO NLP http://www.aclweb.org/anthology/P18-1158	77.003	80.20
4 Nov 08, 2018	BERT_base_aug (ensemble) GammaLab	76.721	79.61
5 Nov 05, 2018	MIR-MRC(F-Net) (single model) Kangwon National University, Natural Language Processing Lab. & ForceWin, KP Lab.	74.791	77.98

BERT can be used in tasks such as question answering systems, sentiment analysis, spam filtering, named entity recognition, document clustering, etc., as the infrastructure or language model for these tasks.

Importing Libraries

```
In [1]:
```

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
```

Load Data

```
In [2]:
```

```
data = pd.read_csv('Tweets.csv')
print('Dataframe:')
data.head(1)
```

Dataframe:

Out[2]:

tweet_id airline_sentiment airline_sentiment_confidence negativereason negative

0 570306133677760513 neutral 1.0 NaN

```
In [3]:
```

```
# Select features
df = data[['text', 'airline_sentiment']]
print('Feature selected DataFrame:')
df.head()
```

Feature selected DataFrame:

Out[3]:

text airline sentiment

0	@VirginAmerica What @dhepburn said.	neutral
1	@VirginAmerica plus you've added commercials t	positive
2	@VirginAmerica I didn't today Must mean I n	neutral
3	@VirginAmerica it's really aggressive to blast	negative
4	@VirginAmerica and it's a really big bad thing	negative

In [4]:

```
df.shape
```

Out[4]:

(14640, 2)

Test and Train dataframes

```
In [5]:
```

```
train_df,eval_df = train_test_split(df,test_size = 0.2)
```

Building a model

In [8]:

```
#pip install simpletransformers
from simpletransformers.classification import ClassificationModel
import torch
# Create a TransformerModel
model = ClassificationModel('bert', 'bert-base-cased', num_labels=3, args={'repr
ocess_input_data': True, 'overwrite_output_dir': True}, use_cuda=False)
```

wandb: WARNING W&B installed but not logged in. Run `wandb login` or set the WANDB API KEY env variable.

Some weights of the model checkpoint at bert-base-cased were not use d when initializing BertForSequenceClassification: ['cls.prediction s.bias', 'cls.predictions.transform.dense.weight', 'cls.predictions.transform.dense.bias', 'cls.predictions.decoder.weight', 'cls.seq_re lationship.weight', 'cls.seq_relationship.bias', 'cls.predictions.tr ansform.LayerNorm.weight', 'cls.predictions.transform.LayerNorm.bia s']

- This IS expected if you are initializing BertForSequenceClassifica tion from the checkpoint of a model trained on another task or with another architecture (e.g. initializing a BertForSequenceClassifica tion model from a BertForPretraining model).
- This IS NOT expected if you are initializing BertForSequenceClassi fication from the checkpoint of a model that you expect to be exactly identical (initializing a BertForSequenceClassification model from a BertForSequenceClassification model).

Some weights of BertForSequenceClassification were not initialized f rom the model checkpoint at bert-base-cased and are newly initialize d: ['classifier.weight', 'classifier.bias']

You should probably TRAIN this model on a down-stream task to be abl e to use it for predictions and inference.

2020/11/9 Lab Bert

```
In [9]:
```

```
def making label(airline sentiment):
    if(airline sentiment=='positive'):
        return 0
    elif(airline sentiment=='neutral'):
        return 2
    else:
        return 1
train df['label'] = train df['airline sentiment'].apply(making label)
eval df['label'] = eval df['airline sentiment'].apply(making label)
print(train df.head())
                                                           ... label
                                                     text
11001
       @USAirways we already spoke to someone several...
                                                           . . .
                                                                   1
1561
                 Qunited I need help with a missing bag.
                                                                   1
2557
       @united my flight out of BGM Cancelled Flightl...
                                                                   1
4148
       @united here I was thinking how I could say so...
                                                                   1
       @SouthwestAir 3 hours and 80 degree difference...
4616
                                                                   0
[5 rows x 3 columns]
/usr/local/lib/python3.6/dist-packages/ipykernel launcher.py:9: Sett
ingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row indexer,col indexer] = value instead
See the caveats in the documentation: https://pandas.pydata.org/pand
as-docs/stable/user guide/indexing.html#returning-a-view-versus-a-co
  if name == ' main ':
/usr/local/lib/python3.6/dist-packages/ipykernel launcher.py:10: Set
tingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row indexer,col indexer] = value instead
See the caveats in the documentation: https://pandas.pydata.org/pand
as-docs/stable/user guide/indexing.html#returning-a-view-versus-a-co
ру
  # Remove the CWD from sys.path while we load stuff.
```

In [10]:

```
train df2 = pd.DataFrame({
    'text': train df['text'].replace(r'\n', ' ', regex=True),
    'label': train df['label']
})
eval df2 = pd.DataFrame({
    'text': eval df['text'].replace(r'\n', ' ', regex=True),
    'label': eval df['label']
})
```

In [11]:

```
model.train model(train df2)
```

/usr/local/lib/python3.6/dist-packages/simpletransformers/classifica tion/classification_model.py:353: UserWarning: Dataframe headers not specified. Falling back to using column 0 as text and column 1 as la bels.

"Dataframe headers not specified. Falling back to using column 0 a s text and column 1 as labels."

/usr/local/lib/python3.6/dist-packages/torch/optim/lr_scheduler.py:2 16: UserWarning: Please also save or load the state of the optimizer when saving or loading the scheduler.

warnings.warn(SAVE STATE WARNING, UserWarning)

Out[11]:

(1464, 0.522169746533596)

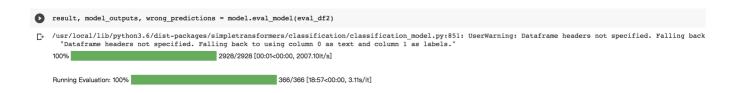


In [12]:

result, model outputs, wrong predictions = model.eval model(eval df2)

/usr/local/lib/python3.6/dist-packages/simpletransformers/classifica tion/classification_model.py:851: UserWarning: Dataframe headers not specified. Falling back to using column 0 as text and column 1 as la bels.

"Dataframe headers not specified. Falling back to using column 0 as text and column 1 as labels."



Model Evaluation

```
In [13]:
print(result)
print(model outputs)
#print(wrong predictions)
{'mcc': 0.6898161979010954, 'eval loss': 0.4371291710053637}
[[-3.34957337 2.45363379 -1.35096169]
 [-2.83464003
                3.16320467 -2.11451244]
 [-3.41498327 \quad 2.94220018 \quad -1.60261059]
 [-2.8468976]
                3.07814193 -2.157468321
 [-3.25913382 \quad 2.70254469 \quad -1.55403459]
 [-2.71210432 \quad 3.05955839 \quad -2.21712756]]
In [14]:
lst = []
for arr in model outputs:
    lst.append(np.argmax(arr))
In [15]:
true = eval_df2['label'].tolist()
predicted = 1st
In [16]:
import sklearn
mat = sklearn.metrics.confusion matrix(true , predicted)
mat
Out[16]:
array([[ 345,
                 84,
                       45],
       [ 33, 1687,
                       83],
          50, 186, 415]])
In [17]:
print(sklearn.metrics.classification report(true,predicted,target names=['positi
ve','neutral','negative']))
               precision
                             recall f1-score
                                                 support
                               0.73
                                          0.76
                                                     474
                    0.81
    positive
     neutral
                    0.86
                               0.94
                                          0.90
                                                    1803
                    0.76
                               0.64
                                          0.70
                                                     651
    negative
    accuracy
                                          0.84
                                                    2928
                    0.81
                               0.77
   macro avg
                                          0.79
                                                    2928
weighted avg
                    0.83
                               0.84
                                          0.83
                                                    2928
```

Our model has an accuracy rate of 83.5%!

```
In [18]:
sklearn.metrics.accuracy_score(true,predicted)
Out[18]:
0.835724043715847
```

Test statement

```
In [19]:

def get_result(text):
    result = model.predict([text])
    pos = np.where(result[1][0] == np.amax(result[1][0]))
    pos = int(pos[0])
    sentiment_dict = {0:'positive',1:'negative',2:'neutral'}
    print(sentiment_dict[pos])
    return

In [20]:

## positive statement
get_result("You are so nice.")

positive

In [21]:

## negative statement
get_result('I hate you.')
```

negative

Save model

```
In [22]:
```

```
from sklearn.externals import joblib
# Save model
joblib.dump(model, 'ML-Model.pkl')
```

/usr/local/lib/python3.6/dist-packages/sklearn/externals/joblib/__in it__.py:15: FutureWarning: sklearn.externals.joblib is deprecated in 0.21 and will be removed in 0.23. Please import this functionality d irectly from joblib, which can be installed with: pip install jobli b. If this warning is raised when loading pickled models, you may ne ed to re-serialize those models with scikit-learn 0.21+.

warnings.warn(msg, category=FutureWarning)

```
Out[22]:
['ML-Model.pkl']
```