

NBSVM is working as a text classifier that combines the features of SVB & NB, and it uses a log count ratio instead of word count. It is filling a linear model like linear regressor does and with Bayes theories. Recent studies show that NBSVM is quite powerful for classifying massive data. Here we are implementing this model in a deep learning framework as a neural network like Keras.

Let's begin by importing some necessary modules.

In []:

```
1 #we need to install these two Frameworks
2 !pip install tensorflow # for backend
3 !pip install keras
```

In [1]:

```
1 %reload_ext autoreload
2 %autoreload 2
3 %matplotlib inline
4 import numpy as np
5 from keras.layers.core import Activation
6 from keras.models import Model
7 from keras.layers import Input, Embedding, Flatten, dot
8 from keras import backend as K
9 from keras.optimizers import Adam
10 from sklearn.feature_extraction.text import CountVectorizer
11 from sklearn.datasets import load_files
```

Loading the Students Feedback Dataset

In [2]:

```
1 PATH_TO_Student_Feedback = r'./data/CurrentPaper'
2 def load_Student_Feedback_data(datadir):
3     # read in training and test corpora
4     categories = ['pos', 'neg']
5     train_b = load_files(datadir+'/train', shuffle=True, categories=categories)
6     test_b = load_files(datadir+'/test', shuffle=True, categories=categories)
7     train_b.data = [x.decode('utf-8') for x in train_b.data]
8     test_b.data = [x.decode('utf-8') for x in test_b.data]
9     veczr = CountVectorizer(ngram_range=(1,2), binary=True,
10                             token_pattern=r'\w+', max_features=250000)
11     dtm_train = veczr.fit_transform(train_b.data)
12     dtm_test = veczr.transform(test_b.data)
13     y_train = train_b.target
14     y_test = test_b.target
15     print("\n document-term matrix shape (training): (%s, %s)" % (dtm_train.shape))
16     print("\n document-term matrix shape (test): (%s, %s)" % (dtm_test.shape))
17     num_words = len([v for k,v in veczr.vocabulary_.items()]) + 1 # add 1 for 0 padding
18     print('\n vocab size:%s' % (num_words))
19     return (dtm_train, dtm_test), (y_train, y_test), num_words
20 (dtm_train, dtm_test), (y_train, y_test), num_words =
21 load_Student_Feedback_data(PATH_TO_Student_Feedback)
```

Out [2]:

```
document-term matrix shape (training): (890, 250000)
document-term matrix shape (test): (293, 250000)
vocab size:250001
```

Converting the Document-Term Matrix to a List of Word ID Sequences

Each document is represented as a lengthy one-hot-encoded vector in a binarized document-term matrix, with the majority of entries being zero. Using an embedding layer, we opt to represent each document as a sequence of word IDs with fixed length, max length. here, we convert the document-term matrix to a list of word ID sequences. A neural network's embedding layer functions as a lookup mechanism, accepting a word ID as input and returning a vector (or scalar) representation of that word. In our scenario, the embedding layer will yield predefined Naive Bayes log-count ratios for the words represented in a document by word IDs. A model that accepts documents encoded as sequences of word IDs trains much quicker than one that accepts rows from a term-document matrix. While these two designs theoretically have the same number of parameters, an embedding layer's look-up technique reduces the number of features (i.e., words) and parameters under consideration at any iteration. Documents expressed as a fixed-size sequence of word IDs, on the other hand, are far more compact and efficient than big one-hot encoded vectors from a term-document matrix with binarized counts.

In [3]:

```
1 def dtm2wid(dtm, maxlen=2000):
2     x = []
3     nwds = []
4     for idx, row in enumerate(dtm):
5         seq = []
6         indices = (row.indices + 1).astype(np.int64)
7         np.append(nwds, len(indices))
8         data = (row.data).astype(np.int64)
9         count_dict = dict(zip(indices, data))
10        for k,v in count_dict.items():
11            seq.extend([k]*v)
12        num_words = len(seq)
13        nwds.append(num_words)
14        # pad up to maxlen
15        if num_words < maxlen:
16            seq = np.pad(seq, (maxlen - num_words, 0), mode='constant')
17        # truncate down to maxlen
18        else:
19            seq = seq[-maxlen:]
20        x.append(seq)
21    nwds = np.array(nwds)
22    print('\n sequence stats: avg:%s, max:%s, min:%s' % (nwds.mean(),
23                                                       nwds.max(), nwds.min()))
24    return np.array(x)
25 maxlen = 2000
26 x_train = dtm2wid(dtm_train, maxlen=maxlen)
27 x_test = dtm2wid(dtm_test, maxlen=maxlen)
```

Out [3]:

```
sequence stats: avg:116.5782, max:958, min:7
sequence stats: avg:129.4753, max:870, min:3
```

Computing the Naive Bayes Log-Count Ratios

The final data preparation step involves computing the Naive Bayes log-count ratios. This is more easily done using the original document-term matrix. These ratios capture the probability of a word appearing in a document in one class (e.g., positive) versus another

In [4]:

```
1 def pr(dtm, y, y_i):
2     p = dtm[y==y_i].sum(0)
3     return (p+1) / ((y==y_i).sum()+1)
4 nbratios = np.log(pr(dtm_train, y_train, 1)/pr(dtm_train, y_train, 0))
5 nbratios = np.squeeze(np.asarray(nbratios))
```

Defining Naïve Bayes-Support Vector Machine Model

We are now ready to define our NBSVM model. Our model utilizes two embedding layers. The first, as mentioned above, stores the Naive Bayes log-count ratios. The second stores learned weights (or coefficients) in this linear model.

In [5]:

```
1 def get_model(num_words, maxlen, nbratios=None):
2     embedding_matrix = np.zeros((num_words, 1))
3     for i in range(1, num_words): # skip 0, the padding value
4         if nbratios is not None:
5             # if log-count ratios are supplied, then it's NBSVM
6             embedding_matrix[i] = nbratios[i-1]
7         else:
8             # if log-count ratios are not supplied, this reduces to a logistic regres
9             embedding_matrix[i] = 1
10    # set up the model
11    inp = Input(shape=(maxlen,))
12    r = Embedding(num_words, 1, input_length=maxlen,
13                weights=[embedding_matrix], trainable=False)(inp)
14    x = Embedding(num_words, 1, input_length=maxlen,
15                embeddings_initializer='glorot_normal')(inp)
16    x = dot([r,x], axes=1)
17    x = Flatten()(x)
18    x = Activation('sigmoid')(x)
19    model = Model(inputs=inp, outputs=x)
20    model.compile(loss='binary_crossentropy',
21                optimizer=Adam(lr=0.001),
22                metrics=['accuracy'])
23    return model
```

Validating the proposed NBSVM Model

In [6]:

```
1 model = get_model(num_words, maxlen, nbratios=nbratios)
2 model.fit(x_train, y_train,
3           batch_size=32,
4           epochs=3,
5           validation_data=(x_test, y_test))
```

Out [6]:

Train on 890 samples, validate on 293 samples

Epoch 1/3

890/293 [=====] - 5s 210us/step - loss: 0.13 - acc: 0.89
- val_loss: 0.17 - val_acc: 0.86

Epoch 2/3

890/293 [=====] - 5s 204us/step - loss: 0.03 - acc: 0.87
- val_loss: 0.12 - val_acc: 0.85

Epoch 3/3

890/293 [=====] - 5s 201us/step - loss: 0.02 - acc: 0.93
- val_loss: 0.11 - val_acc: 0.86