

# Libraries

In [1]:

```
1 #Data:
2 import pandas as pd
3 import numpy as np
4 from collections import Counter
5
6 # Text preprocessing:
7 from keras.preprocessing.sequence import pad_sequences
8 from keras.preprocessing.text import Tokenizer, one_hot
9 from sklearn.feature_extraction.text import CountVectorizer
10
11 # Layers:
12 from keras.layers import Input, Embedding, LSTM, Dense, Flatten, concatenate, Dropout,
13
14 # Model:
15 from keras.models import Model
16
17 # Metrics:
18 from sklearn.metrics import roc_auc_score
19
20 from time import time
21 import keras
22 import matplotlib.pyplot as plt
23 from keras.callbacks import TensorBoard
24
25 from sklearn.utils import class_weight
26
27 %%matplotlibinline
```

Using TensorFlow backend.

## Load data

In [2]:

```
1 data = pd.read_csv('data/final_features.csv')
2
3 data.sort_values(by='project_submitted_datetime',inplace=True)
4
5 Counter(data['project_is_approved'])
6 print('ratio excepted: ', round(Counter(data['project_is_approved']).get(1)/data.shape
7 print('ratio rejected: ', round(Counter(data['project_is_approved']).get(0)/data.shape
8
```

ratio excepted: 85.0 %

ratio rejected: 15.0 %

In [3]:

```
1 for feature in data.iteritems():
2     print(feature[0], ': ', 'has', str(data[data[feature[0]].isnull().values][feature[0]]
```

```
project_submitted_datetime : has 0 missing values
clean_teacher_prefix : has 0 missing values
clean_school_state : has 0 missing values
clean_grade_categories : has 0 missing values
clean_subject_categories : has 0 missing values
clean_subject_subcategories : has 0 missing values
clean_project_title : has 43 missing values
clean_essay : has 0 missing values
clean_resource_summary : has 0 missing values
resource_summary_contains_numerical_digits : has 0 missing values
std_price : has 0 missing values
std_quantity : has 0 missing values
std_teacher_number_of_previously_posted_projects : has 0 missing values
nrm_price : has 0 missing values
nrm_quantity : has 0 missing values
nrm_teacher_number_of_previously_posted_projects : has 0 missing values
project_is_approved : has 0 missing values
```

In [4]:

```
1 data.fillna(value={'clean_project_title':''}, inplace=True)
```

In [5]:

```
1 for feature in data.iteritems():
2     print(feature[0], ': ', 'has', str(data[data[feature[0]].isnull().values][feature[0]]
```

```
project_submitted_datetime : has 0 missing values
clean_teacher_prefix : has 0 missing values
clean_school_state : has 0 missing values
clean_grade_categories : has 0 missing values
clean_subject_categories : has 0 missing values
clean_subject_subcategories : has 0 missing values
clean_project_title : has 0 missing values
clean_essay : has 0 missing values
clean_resource_summary : has 0 missing values
resource_summary_contains_numerical_digits : has 0 missing values
std_price : has 0 missing values
std_quantity : has 0 missing values
std_teacher_number_of_previously_posted_projects : has 0 missing values
nrm_price : has 0 missing values
nrm_quantity : has 0 missing values
nrm_teacher_number_of_previously_posted_projects : has 0 missing values
project_is_approved : has 0 missing values
```

In [6]:

```

1 data['total_text_data'] = data['clean_essay']
2
3 data.drop(labels=['clean_project_title','clean_essay','clean_resource_summary'],axis=1
4

```

## Split data into train,CV and test

In [7]:

```

1 # data = data.iloc[0:1000,: ]
2
3 data_train = data.iloc[0:int(data.shape[0]*0.8),:]
4 data_train = data_train.iloc[0:int(data_train.shape[0]*0.8),:]
5 data_cv = data_train.iloc[int(data_train.shape[0]*0.8):,:]
6 data_test = data.iloc[int(data.shape[0]*0.8):,:]
7
8 Y_train = data_train['project_is_approved']
9 data_train.drop(labels=['project_is_approved'],axis=1,inplace=True)
10 X_train = data_train
11
12 Y_cv = data_cv['project_is_approved']
13 data_cv.drop(labels=['project_is_approved'],axis=1,inplace=True)
14 X_cv = data_cv
15
16 Y_test = data_test['project_is_approved']
17 data_test.drop(labels=['project_is_approved'],axis=1,inplace=True)
18 X_test = data_test
19

```

c:\users\byron\applications\pythonmaster\lib\site-packages\pandas\core\frame.py:3697: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy> (<http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>)  
errors=errors)

In [8]:

```

1 print('X_train shape: ',X_train.shape, 'Y_train shape: ',Y_train.shape)
2 print('X_cv shape: ',X_cv.shape, 'Y_cv shape: ',Y_cv.shape)
3 print('X_test shape: ',X_test.shape, 'Y_test shape: ',Y_test.shape)
4

```

X\_train shape: (69918, 14) Y\_train shape: (69918,)  
X\_cv shape: (13984, 14) Y\_cv shape: (13984,)  
X\_test shape: (21850, 14) Y\_test shape: (21850,)

# Turn text data into sequence data - text preprocessing

In [9]:

```
1  # Use training data:
2  text = Tokenizer()
3  text.fit_on_texts(X_train['total_text_data'])
4  text_sequences_train = text.texts_to_sequences(X_train['total_text_data'])
5
6  def max_length(observation_text):
7      observation_text_lengths = list()
8      for obs in observation_text:
9          observation_text_lengths.append(len(obs.split()))
10     return np.mean(observation_text_lengths)
11
12 max_len = int(max_length(X_train['total_text_data'].values))
13
14 text_sequences_train = pad_sequences(sequences=text_sequences_train, maxlen=max_len, padding='left')
15
16 dictionary = text.word_index
17
18 frequencies = text.word_counts
19 frequencies = dict(frequencies)
20 vocab_size = len(dictionary.keys()) + 1
21
22 # Transform cv and test data
23 text_sequences_cv = text.texts_to_sequences(X_cv['total_text_data'])
24 text_sequences_cv = pad_sequences(sequences=text_sequences_cv, maxlen=max_len, padding='left')
25 text_sequences_test = text.texts_to_sequences(X_test['total_text_data'])
26 text_sequences_test = pad_sequences(sequences=text_sequences_test, maxlen=max_len, padding='left')
27
```

## Get pretrained word embeddings

In [10]:

```

1 embeddings_index = dict()
2 with open('glove/glove.6B.100d.txt', 'r', encoding="utf-8") as f:
3     for line in f:
4         values = line.split()
5         word = values[0]
6         coefs = np.array(values[1:], dtype='float32')
7         embeddings_index[word] = coefs
8 print('Loaded {} word vectors.'.format(len(embeddings_index)))
9
10 embedding_matrix = np.zeros((vocab_size, 100))
11 for word, i in dictionary.items():
12     embedding_vector = embeddings_index.get(word)
13     if embedding_vector is not None:
14         embedding_matrix[i] = embedding_vector
15

```

Loaded 400000 word vectors.

## One hot encode categorical features

In [11]:

```

1 # Encode teacher prefix:
2 teacher_vec = CountVectorizer(binary=True)
3 clean_teacher_prefix_ohe_train = teacher_vec.fit_transform(X_train['clean_teacher_prefix'])
4 clean_teacher_prefix_ohe_cv = teacher_vec.transform(X_cv['clean_teacher_prefix'])
5 clean_teacher_prefix_ohe_test = teacher_vec.transform(X_test['clean_teacher_prefix'])
6
7 # Encode school state:
8 school_vec = CountVectorizer(binary=True)
9 clean_school_state_ohe_train = school_vec.fit_transform(X_train['clean_school_state'])
10 clean_school_state_ohe_cv = school_vec.transform(X_cv['clean_school_state'])
11 clean_school_state_ohe_test = school_vec.transform(X_test['clean_school_state'])
12
13 # Encode grade categories:
14 grade_vec = CountVectorizer(binary=True)
15 clean_grade_categories_ohe_train = grade_vec.fit_transform(X_train['clean_grade_categories'])
16 clean_grade_categories_ohe_cv = grade_vec.transform(X_cv['clean_grade_categories'])
17 clean_grade_categories_ohe_test = grade_vec.transform(X_test['clean_grade_categories'])
18
19 # Encode subject categories:
20 cat_vec = CountVectorizer(binary=True)
21 clean_subject_categories_ohe_train = cat_vec.fit_transform(X_train['clean_subject_categories'])
22 clean_subject_categories_ohe_cv = cat_vec.transform(X_cv['clean_subject_categories'])
23 clean_subject_categories_ohe_test = cat_vec.transform(X_test['clean_subject_categories'])
24
25 # Encode subject subcategories:
26 subcat_vec = CountVectorizer(binary=True)
27 clean_subject_subcategories_ohe_train = subcat_vec.fit_transform(X_train['clean_subject_subcategories'])
28 clean_subject_subcategories_ohe_cv = subcat_vec.transform(X_cv['clean_subject_subcategories'])
29 clean_subject_subcategories_ohe_test = subcat_vec.transform(X_test['clean_subject_subcategories'])

```

In [12]:

```

1  def concat_numeric_features(dataset):
2      if dataset == 'train':
3          a = pd.DataFrame(clean_teacher_prefix_ohe_train.toarray())
4          b = pd.DataFrame(clean_school_state_ohe_train.toarray())
5          c = pd.DataFrame(clean_grade_categories_ohe_train.toarray())
6          d = pd.DataFrame(clean_subject_categories_ohe_train.toarray())
7          e = pd.DataFrame(clean_subject_subcategories_ohe_train.toarray())
8          f = pd.DataFrame(X_train['resource_summary_contains_numerical_digits'].values)
9          g = pd.DataFrame(X_train['nrm_teacher_number_of_previously_posted_projects'].values)
10         h = pd.DataFrame(X_train['nrm_price'].values)
11         i = pd.DataFrame(X_train['nrm_quantity'].values)
12     elif dataset == 'cv':
13         a = pd.DataFrame(clean_teacher_prefix_ohe_cv.toarray())
14         b = pd.DataFrame(clean_school_state_ohe_cv.toarray())
15         c = pd.DataFrame(clean_grade_categories_ohe_cv.toarray())
16         d = pd.DataFrame(clean_subject_categories_ohe_cv.toarray())
17         e = pd.DataFrame(clean_subject_subcategories_ohe_cv.toarray())
18         f = pd.DataFrame(X_cv['resource_summary_contains_numerical_digits'].values)
19         g = pd.DataFrame(X_cv['nrm_teacher_number_of_previously_posted_projects'].values)
20         h = pd.DataFrame(X_cv['nrm_price'].values)
21         i = pd.DataFrame(X_cv['nrm_quantity'].values)
22
23     else:
24         a = pd.DataFrame(clean_teacher_prefix_ohe_test.toarray())
25         b = pd.DataFrame(clean_school_state_ohe_test.toarray())
26         c = pd.DataFrame(clean_grade_categories_ohe_test.toarray())
27         d = pd.DataFrame(clean_subject_categories_ohe_test.toarray())
28         e = pd.DataFrame(clean_subject_subcategories_ohe_test.toarray())
29         f = pd.DataFrame(X_test['resource_summary_contains_numerical_digits'].values)
30         g = pd.DataFrame(X_test['nrm_teacher_number_of_previously_posted_projects'].values)
31         h = pd.DataFrame(X_test['nrm_price'].values)
32         i = pd.DataFrame(X_test['nrm_quantity'].values)
33
34     concat = pd.concat(objs=[a,b,c,d,e,f,g,h,i],axis=1)
35     del a,b,c,d,e,f,g,h,i
36     return concat
37
38 numeric_train = concat_numeric_features(dataset='train')
39 numeric_cv = concat_numeric_features(dataset='cv')
40 numeric_test = concat_numeric_features(dataset='test')
41
42 numeric_dim = numeric_train.shape[1]
43
44 numeric_train = np.array(numeric_train).reshape((numeric_train.shape[0],numeric_train.shape[1]))
45 numeric_cv = np.array(numeric_cv).reshape((numeric_cv.shape[0],numeric_cv.shape[1],1))
46 numeric_test = np.array(numeric_test).reshape((numeric_test.shape[0],numeric_test.shape[1],1))

```

## Building the Neural Network

In [13]:

```

1 text_input = Input(shape=(max_len,), dtype='int32', name='text_input')
2 embedding = Embedding(input_dim=vocab_size, output_dim=100, weights=[embedding_matrix])
3 lstm = LSTM(units=100, activation='relu', kernel_initializer='glorot_normal', bias_initializer='zeros')
4 flatten_1 = Flatten()(lstm)
5 #####
6 numeric_input = Input(shape=(numeric_dim,1), dtype='float32', name='numeric_input')
7 conv1D_1 = Conv1D(filters=30, kernel_size=4, strides=1, padding='valid', activation='relu')
8 conv1D_2 = Conv1D(filters=30, kernel_size=4, strides=1, padding='valid', activation='relu')
9 flatten_2 = Flatten()(conv1D_2)
10 #####
11 concat = concatenate([flatten_1, flatten_2])
12 #####
13 dense_1 = Dense(units=100, activation='relu', kernel_initializer='he_normal', bias_initializer='zeros')
14 drop_1 = Dropout(rate=0.5)(dense_1)
15 dense_2 = Dense(units=100, activation='relu', kernel_initializer='he_normal', bias_initializer='zeros')
16 drop_2 = Dropout(rate=0.5)(dense_2)
17 dense_3 = Dense(units=100, activation='relu', kernel_initializer='he_normal', bias_initializer='zeros')
18 output = Dense(units=1, activation='sigmoid', name='output')(dense_3)
19

```

WARNING:tensorflow:From c:\users\byron\applications\pythonmaster\lib\site-packages\tensorflow\python\framework\op\_def\_library.py:263: colocate\_with (from tensorflow.python.framework.ops) is deprecated and will be removed in a future version.

Instructions for updating:

Colocations handled automatically by placer.

WARNING:tensorflow:From c:\users\byron\applications\pythonmaster\lib\site-packages\keras\backend\tensorflow\_backend.py:3445: calling dropout (from tensorflow.python.ops.nn\_ops) with keep\_prob is deprecated and will be removed in a future version.

Instructions for updating:

Please use `rate` instead of `keep\_prob`. Rate should be set to `rate = 1 - keep\_prob`.

In [14]:

```

1 model = Model(inputs=[text_input, numeric_input], outputs=[output])


```

In [15]:

```
1 model.summary()
```

Layer (type)	Output Shape	Param #	Connected to
text_input (InputLayer)	(None, 137)	0	
numeric_input (InputLayer)	(None, 501, 1)	0	
embedding_1 (Embedding) [0][0]	(None, 137, 100)	4668500	text_input
conv1d_1 (Conv1D) ut[0][0]	(None, 498, 30)	150	numeric_in
lstm_1 (LSTM) [0][0]	(None, 137, 100)	80400	embedding_1
conv1d_2 (Conv1D) [0]	(None, 495, 30)	3630	conv1d_1[0]
flatten_1 (Flatten) [0]	(None, 13700)	0	lstm_1[0]
flatten_2 (Flatten) [0]	(None, 14850)	0	conv1d_2[0]
concatenate_1 (Concatenate) [0][0]	(None, 28550)	0	flatten_1 flatten_2
dense_1 (Dense) _1[0][0]	(None, 100)	2855100	concatenate
dropout_1 (Dropout) [0]	(None, 100)	0	dense_1[0]
dense_2 (Dense) [0][0]	(None, 100)	10100	dropout_1



8/18/2019	LSTM_model3		
dropout_2 (Dropout) [0]	(None, 100)	0	dense_2[0]
<hr/>			
dense_3 (Dense) [0][0]	(None, 100)	10100	dropout_2
<hr/>			
output (Dense) [0]	(None, 1)	101	dense_3[0]
<hr/>			
=====			
=====			
Total params: 7,628,081			
Trainable params: 2,959,581			
Non-trainable params: 4,668,500			
<hr/>			
<hr/>			
			

In [16]:

```
1 tensorboard = TensorBoard(log_dir='.logs/{}'.format(time()))
```

In [17]:

```
1 model.compile(optimizer=keras.optimizers.Adam(lr=0.0001), loss='binary_crossentropy',
```

In [18]:

```
1 class roc_callback(keras.callbacks.Callback):
2     def __init__(self, training_data, validation_data):
3         self.x = training_data[0]
4         self.y = training_data[1]
5         self.x_val = validation_data[0]
6         self.y_val = validation_data[1]
7
8     def on_train_begin(self, logs={}):
9         return
10
11    def on_train_end(self, logs={}):
12        return
13
14    def on_epoch_begin(self, epoch, logs={}):
15        return
16
17    def on_epoch_end(self, epoch, logs={}):
18        y_pred = self.model.predict(self.x)
19        roc = roc_auc_score(self.y, y_pred)
20        y_pred_val = self.model.predict(self.x_val)
21        roc_val = roc_auc_score(self.y_val, y_pred_val)
22        print('\rroc-auc: %s - roc-auc_val: %s' % (str(round(roc,4)),str(round(roc_val
23        return
24
25    def on_batch_begin(self, batch, logs={}):
26        return
27
28    def on_batch_end(self, batch, logs={}):
29        return
30
```

In [19]:

```
1 class_weights = class_weight.compute_class_weight('balanced',
2                                                    np.unique(Y_train),
3                                                    Y_train)
```

In [20]:

```

1 batch_size = 100
2 epochs = 15
3 history = model.fit(x=[text_sequences_train,numeric_train],y=[Y_train],validation_data=
4                     callbacks=[tensorboard,
5                               roc_callback(training_data=([text_sequences_train,numeric_train],
6                                                         validation_data=([text_sequences_cv,numeric_cv],
7                               class_weight=class_weights
8                     )

```

WARNING:tensorflow:From c:\users\byron\applications\pythonmaster\lib\site-packages\tensorflow\python\ops\math\_ops.py:3066: to\_int32 (from tensorflow.python.ops.math\_ops) is deprecated and will be removed in a future version.

Instructions for updating:

Use tf.cast instead.

Train on 69918 samples, validate on 13984 samples

Epoch 1/15

69918/69918 [=====] - 159s 2ms/step - loss: 0.4381

- acc: 0.8415 - val\_loss: 0.3998 - val\_acc: 0.8593

roc-auc: 0.7031 - roc-auc\_val: 0.7003

Epoch 2/15

69918/69918 [=====] - 166s 2ms/step - loss: 0.4125

- acc: 0.8441 - val\_loss: 0.3893 - val\_acc: 0.8593

roc-auc: 0.7194 - roc-auc\_val: 0.7175

Epoch 3/15

69918/69918 [=====] - 167s 2ms/step - loss: 0.4027

- acc: 0.8440 - val\_loss: 0.3716 - val\_acc: 0.8594

roc-auc: 0.7322 - roc-auc\_val: 0.7298

Epoch 4/15

69918/69918 [=====] - 165s 2ms/step - loss: 0.3958

- acc: 0.8442 - val\_loss: 0.3805 - val\_acc: 0.8601

roc-auc: 0.7439 - roc-auc\_val: 0.743

Epoch 5/15

69918/69918 [=====] - 162s 2ms/step - loss: 0.3901

- acc: 0.8449 - val\_loss: 0.3789 - val\_acc: 0.8666

roc-auc: 0.7514 - roc-auc\_val: 0.7512

Epoch 6/15

69918/69918 [=====] - 163s 2ms/step - loss: 0.3846

- acc: 0.8463 - val\_loss: 0.3594 - val\_acc: 0.8666

roc-auc: 0.7628 - roc-auc\_val: 0.7628

Epoch 7/15

69918/69918 [=====] - 164s 2ms/step - loss: 0.3804

- acc: 0.8474 - val\_loss: 0.3564 - val\_acc: 0.8696

roc-auc: 0.7701 - roc-auc\_val: 0.7707

Epoch 8/15

69918/69918 [=====] - 163s 2ms/step - loss: 0.3769

- acc: 0.8489 - val\_loss: 0.3511 - val\_acc: 0.8703

roc-auc: 0.7796 - roc-auc\_val: 0.7814

Epoch 9/15

69918/69918 [=====] - 165s 2ms/step - loss: 0.3706

- acc: 0.8505 - val\_loss: 0.3477 - val\_acc: 0.8712

roc-auc: 0.788 - roc-auc\_val: 0.7896

Epoch 10/15

69918/69918 [=====] - 170s 2ms/step - loss: 0.3658

- acc: 0.8521 - val\_loss: 0.3315 - val\_acc: 0.8734

roc-auc: 0.7984 - roc-auc\_val: 0.8017

Epoch 11/15

69918/69918 [=====] - 167s 2ms/step - loss: 0.3611

- acc: 0.8550 - val\_loss: 0.3324 - val\_acc: 0.8770

roc-auc: 0.8089 - roc-auc\_val: 0.8117

Epoch 12/15

69918/69918 [=====] - 169s 2ms/step - loss: 0.3557

- acc: 0.8572 - val\_loss: 0.3362 - val\_acc: 0.8771

roc-auc: 0.8177 - roc-auc\_val: 0.8191

Epoch 13/15

69918/69918 [=====] - 171s 2ms/step - loss: 0.3506

- acc: 0.8590 - val\_loss: 0.3211 - val\_acc: 0.8837

roc-auc: 0.8297 - roc-auc\_val: 0.8324

Epoch 14/15

69918/69918 [=====] - 168s 2ms/step - loss: 0.3424

- acc: 0.8630 - val\_loss: 0.3082 - val\_acc: 0.8854

roc-auc: 0.839 - roc-auc\_val: 0.8421

Epoch 15/15

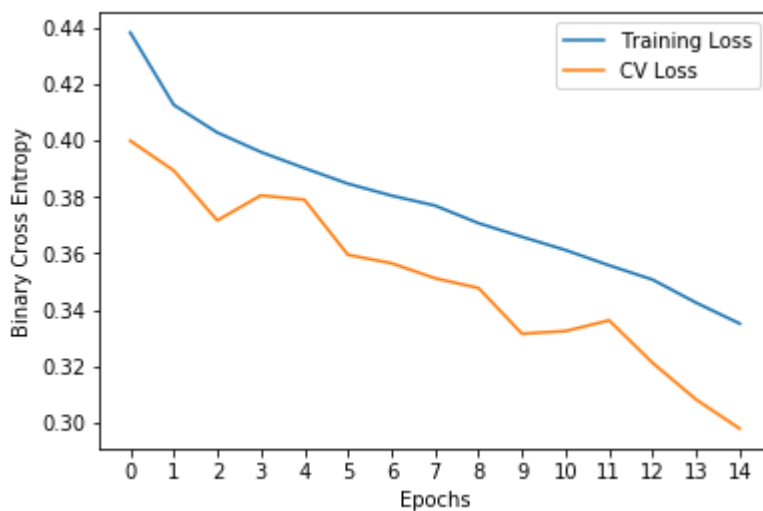
69918/69918 [=====] - 166s 2ms/step - loss: 0.3351

- acc: 0.8665 - val\_loss: 0.2979 - val\_acc: 0.8907

roc-auc: 0.8503 - roc-auc\_val: 0.852

In [21]:

```
1 plt.plot(np.arange(epochs),history.history['loss'],label='Training Loss')
2 plt.plot(np.arange(epochs),history.history['val_loss'],label='CV Loss')
3 plt.xlabel('Epochs')
4 plt.ylabel('Binary Cross Entropy')
5 plt.legend()
6 plt.xticks(np.arange(epochs))
7 plt.show()
```



In [22]:

```
1 model.evaluate(x=[text_sequences_test,numeric_test],y=[Y_test],batch_size=batch_size)
```

21850/21850 [=====] - 15s 696us/step

Out[22]:

[0.3769638879348266, 0.8518077769323127]

## AUC Score on test data

In [23]:

```
1 Y_pred = model.predict(x=[text_sequences_test,numeric_test],  
2                           batch_size=batch_size)  
3 roc_auc_score(Y_test, Y_pred)
```

Out[23]:

0.7229880277013565