

2. Model for Training

Input to this phase are the pickle files that have been obtained as output in the previous phase.

In [1]:

```
FN = 'train'
```

Import all the necessary package for building the model, this project works on Theano backend on Keras version 1. So if you wish to replicate the project in higher versions please make necessary changes as some of the packages have deprecated/replaced by newer ones.

In [3]:

```
import os
# os.environ['THEANO_FLAGS'] = 'device=cpu,floatX=float32'

import keras
keras.__version__
```

Using Theano backend.

Out[3]:

```
'1.0.4'
```

Do all the necessary initialization

In [4]:

```
FN0 = 'vocabulary-embedding'

FN1 = 'train'

maxlend=25 # 0 - if we dont want to use description at all
maxlenh=25
maxlen = maxlend + maxlenh
rnn_size = 512 # must be same as 160330-word-gen
rnn_layers = 3 # match FN1
batch_norm=False

activation_rnn_size = 40 if maxlend else 0

# training parameters
seed=42
p_W, p_U, p_dense, p_emb, weight_decay = 0, 0, 0, 0, 0
optimizer = 'adam'
LR = 1e-4
batch_size=64
nflips=10

nb_train_samples = 30000
nb_val_samples = 3000
```

2.1 Read word embeddings

Change the path to the files according to the file location in your system, in the code the files are assumed to be in the location ./data/

In [11]:

```
with open('data/%s.pkl'%FN0, 'rb') as fp:  
    embedding, idx2word, word2idx, glove_idx2idx = pickle.load(fp)  
vocab_size, embedding_size = embedding.shape
```

In [12]:

```
with open('data/%s.data.pkl'%FN0, 'rb') as fp:  
    X, Y = pickle.load(fp)
```

In [13]:

```
nb_unknown_words = 10
```

In [14]:

```
print 'number of examples',len(X),len(Y)  
print 'dimension of embedding space for words',embedding_size  
print 'vocabulary size', vocab_size, 'the last %d words can be used as place holders for  
unknown/oov words'%nb_unknown_words  
print 'total number of different words',len(idx2word), len(word2idx)  
print 'number of words outside vocabulary which we can substitue using glove similarity',  
len(glove_idx2idx)  
print 'number of words that will be regarded as unknonw(unk) /out-of-vocabulary(oov)',len(idx2word)  
-vocab_size-len(glove_idx2idx)
```

```
number of examples 684114 684114  
dimension of embedding space for words 100  
vocabulary size 40000 the last 10 words can be used as place holders for unknown/oov words  
total number of different words 523734 523734  
number of words outside vocabulary which we can substitue using glove similarity 122377  
number of words that will be regarded as unknonw(unk) /out-of-vocabulary(oov) 361357
```

In [15]:

```
for i in range(nb_unknown_words):  
    idx2word[vocab_size-1-i] = '<%d>'%i
```

when printing mark words outside vocabulary with ^ at their end

In [16]:

```
oov0 = vocab_size-nb_unknown_words
```

In [17]:

```
for i in range(oov0, len(idx2word)):  
    idx2word[i] = idx2word[i]+'^'
```

In [18]:

```
from sklearn.model_selection import train_test_split  
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=nb_val_samples, random_state=seed)  
len(X_train), len(Y_train), len(X_test), len(Y_test)
```

Out[18]:

```
(681114, 681114, 3000, 3000)
```

In [19]:

```
del X  
del Y
```

```
empty = 0
eos = 1
idx2word[empty] = '_'
idx2word[eos] = '~'
```

In [21]:

```
import numpy as np
from keras.preprocessing import sequence
from keras.utils import np_utils
import random, sys
```

In [22]:

```
def prt(label, x):
    print label+':',
    for w in x:
        print idx2word[w],
    print
```

In [23]:

```
i = 334
prt('H',Y_train[i])
prt('D',X_train[i])
```

H: Behold the Fastest CF Card Ever Made (For Now)

D: Are you worried your compact flash card wo n't be able to keep up with the demands of the upcoming 4K revolution ? Toshiba 's got your back with its new Exceria^ Pro series CF cards that boast read and write speeds of 160MB/s^ and 150MB/s^ respectively , making them the world 's fastest -- for the time being .

In [24]:

```
i = 334
prt('H',Y_test[i])
prt('D',X_test[i])
```

H: Obama : Bad for Black People

D: Are you one of the 48 % of Americans who is `` hearing too much about Barack Obama '' ? Then you certainly wo n't like this Sunday 's story by professional Democratic Party Underminer Matt Bai^ . It 's about how Barack Obama represents the End of Black Politics , because he 's a black person who white people do n't feel threatened by . In the story , Bai^ harangues^ Philadelphia mayor Michael Nutter^ about why he di

2.2 Model Definition

In [25]:

```
from keras.models import Sequential
from keras.layers.core import Dense, Activation, Dropout, RepeatVector, Merge
from keras.layers.wrappers import TimeDistributed
from keras.layers.recurrent import LSTM
from keras.layers.embeddings import Embedding
from keras.regularizers import l2
```

In [26]:

```
# seed weight initialization
random.seed(seed)
np.random.seed(seed)
```

In [27]:

start with a standard stacked LSTM

In [28]:

```
model = Sequential()
model.add(Embedding(vocab_size, embedding_size,
                    input_length=maxlen,
                    W_regularizer=regularizer, dropout=p_emb, weights=[embedding], mask_zero=True,
                    name='embedding_1'))
for i in range(rnn_layers):
    lstm = LSTM(rnn_size, return_sequences=True, # batch_norm=batch_norm,
                W_regularizer=regularizer, U_regularizer=regularizer,
                b_regularizer=regularizer, dropout_W=p_W, dropout_U=p_U,
                name='lstm_%d'%(i+1))
    model.add(lstm)
model.add(Dropout(p_dense, name='dropout_%d'%(i+1)))
```

A special layer that reduces the input just to its headline part (second half). For each word in this part it concatenate the output of the previous layer (RNN) with a weighted average of the outputs of the description part. In this only the last rnn_size - activation_rnn_size are used from each output. The first activation_rnn_size output is used to computer the weights for the averaging.

In [29]:

```
from keras.layers.core import Lambda
import keras.backend as K

def simple_context(X, mask, n=activation_rnn_size, maxlen=maxlen, maxlenh=maxlenh):
    desc, head = X[:, :maxlen, :], X[:, maxlen:, :]
    head_activations, head_words = head[:, :, :n], head[:, :, n:]
    desc_activations, desc_words = desc[:, :, :n], desc[:, :, n:]

    # RTFM
    http://deeplearning.net/software/theano/library/tensor/basic.html#theano.tensor.batched_tensordot
    # activation for every head word and every desc word
    activation_energies = K.batch_dot(head_activations, desc_activations, axes=(2, 2))
    # make sure we dont use description words that are masked out
    activation_energies = activation_energies + -1e20*K.expand_dims(1.-K.cast(mask[:, :maxlen], 'float32'), 1)

    # for every head word compute weights for every desc word
    activation_energies = K.reshape(activation_energies, (-1, maxlen))
    activation_weights = K.softmax(activation_energies)
    activation_weights = K.reshape(activation_weights, (-1, maxlenh, maxlen))

    # for every head word compute weighted average of desc words
    desc_avg_word = K.batch_dot(activation_weights, desc_words, axes=(2, 1))
    return K.concatenate((desc_avg_word, head_words))

class SimpleContext(Lambda):
    def __init__(self, **kwargs):
        super(SimpleContext, self).__init__(simple_context, **kwargs)
        self.supports_masking = True

    def compute_mask(self, input, input_mask=None):
        return input_mask[:, maxlen:]

    def get_output_shape_for(self, input_shape):
        nb_samples = input_shape[0]
        n = 2*(rnn_size - activation_rnn_size)
        return (nb_samples, maxlenh, n)
```

In [30]:

```
if activation_rnn_size:
    model.add(SimpleContext(name='simplecontext_1'))
model.add(TimeDistributed(Dense(vocab_size,
                               W_regularizer=regularizer, b_regularizer=regularizer,
```

```
In [31]:
```

```
from keras.optimizers import Adam, RMSprop # usually I prefer Adam but article used rmsprop
# opt = Adam(lr=LR) # keep calm and reduce learning rate
model.compile(loss='categorical_crossentropy', optimizer=optimizer)
```

```
In [33]:
```

```
K.set_value(model.optimizer.lr,np.float32(LR))
```

```
In [34]:
```

```
def str_shape(x):
    return 'x'.join(map(str,x.shape))

def inspect_model(model):
    for i,l in enumerate(model.layers):
        print i, 'cls=%s name=%s'%(type(l).__name__, l.name)
        weights = l.get_weights()
        for weight in weights:
            print str_shape(weight),
        print
```

```
In [35]:
```

```
inspect_model(model)
```

```
0 cls=Embedding name=embedding_1
40000x100
1 cls=LSTM name=lstm_1
100x512 512x512 512 512x512 512 100x512 512x512 512 100x512 512x512 512
2 cls=Dropout name=dropout_1

3 cls=LSTM name=lstm_2
512x512 512x512 512 512x512 512x512 512 512x512 512 512x512 512x512 512
4 cls=Dropout name=dropout_2

5 cls=LSTM name=lstm_3
512x512 512x512 512 512x512 512x512 512 512x512 512 512x512 512x512 512
6 cls=Dropout name=dropout_3

7 cls=SimpleContext name=simplecontext_1

8 cls=TimeDistributed name=timedistributed_1
944x40000 40000
9 cls=Activation name=activation_1
```

2.3 Load

Load any pre-trained model if any

```
In [36]:
```

```
if FN1 and os.path.exists('data/%s.hdf5'%FN1):
    model.load_weights('data/%s.hdf5'%FN1)
```

2.4 Test

```
In [36]:
```

```
def lpadd(x, maxlen=maxlen, eos=eos):
    """Pad sequence x to maxlen with eos at the end. maxlen >= len(x)."
```

```

assert maxlen >= 0
if maxlen == 0:
    return [eos]
n = len(x)
if n > maxlen:
    x = x[-maxlen:]
    n = maxlen
return [empty]*(maxlen-n) + x + [eos]

```

In [37]:

```

samples = [lpadd([3]*26)]
# pad from right (post) so the first maxlen will be description followed by headline
data = sequence.pad_sequences(samples, maxlen=maxlen, value=empty, padding='post', truncating='post')

```

In [38]:

```
np.all(data[:,maxlen] == eos)
```

Out[38]:

True

In [39]:

```
data.shape, map(len, samples)
```

Out[39]:

((1, 50), [26])

2.5 Data generator

Data generator generates batches of inputs and outputs/labels for training. The inputs are each made from two parts. The first maxlen words are the original description, followed by eos followed by the headline which we want to predict, except for the last word in the headline which is always eos and then empty padding until maxlen words.

For each, input, the output is the headline words (without the start eos but with the ending eos) padded with empty words up to maxlen words. The output is also expanded to be y-hot encoding of each word.

To be more realistic, the second part of the input should be the result of generation and not the original headline. Instead we will flip just nflips words to be from the generator, but even this is too hard and instead implement flipping in a naive way (which consumes less time.) Using the full input (description + eos + headline) generate predictions for outputs. For nflips random words from the output, replace the original word with the word with highest probability from the prediction.

In [48]:

```

def flip_headline(x, nflips=None, model=None, debug=False):
    """given a vectorized input (after `pad_sequences`) flip some of the words in the second half (headline)
    with words predicted by the model
    """
    if nflips is None or model is None or nflips <= 0:
        return x

    batch_size = len(x)
    assert np.all(x[:,maxlen] == eos)
    probs = model.predict(x, verbose=0, batch_size=batch_size)
    x_out = x.copy()
    for b in range(batch_size):
        # pick locations we want to flip
        # 0...maxlen-1 are descriptions and should be fixed
        # maxlen is eos and should be fixed
        flips = sorted(random.sample(xrange(maxlen+1,maxlen), nflips))
        if debug and b < debug:

```

```

    if x[b, input_idx] == empty or x[b, input_idx] == eos:
        continue
    # convert from input location to label location
    # the output at maxlen (when input is eos) is feed as input at maxlen+1
    label_idx = input_idx - (maxlen+1)
    prob = probs[b, label_idx]
    w = prob.argmax()
    if w == empty: # replace accidental empty with oov
        w = oov0
    if debug and b < debug:
        print '%s => %s'%(idx2word[x_out[b, input_idx]], idx2word[w]),
        x_out[b, input_idx] = w
if debug and b < debug:
    print
return x_out

```

In [49]:

```

def conv_seq_labels(xds, xhs, nflips=None, model=None, debug=False):
    """description and headlines are converted to padded input vectors. headlines are one-hot to label"""
    batch_size = len(xhs)
    assert len(xds) == batch_size
    x = [vocab_fold(lpadd(xd)+xh) for xd,xh in zip(xds,xhs)] # the input does not have 2nd eos
    x = sequence.pad_sequences(x, maxlen=maxlen, value=empty, padding='post', truncating='post')
    x = flip_headline(x, nflips=nflips, model=model, debug=debug)

    y = np.zeros((batch_size, maxlen, vocab_size))
    for i, xh in enumerate(xhs):
        xh = vocab_fold(xh) + [eos] + [empty]*maxlen # output does have a eos at end
        xh = xh[:maxlen]
        y[i,:,:] = np_utils.to_categorical(xh, vocab_size)

    return x, y

```

In [50]:

```

def gen(Xd, Xh, batch_size=batch_size, nb_batches=None, nflips=None, model=None, debug=False, seed=seed):
    """yield batches. for training use nb_batches=None
    for validation generate deterministic results repeating every nb_batches

    while training it is good idea to flip once in a while the values of the headlines from the
    value taken from Xh to value generated by the model.
    """
    c = nb_batches if nb_batches else 0
    while True:
        xds = []
        xhs = []
        if nb_batches and c >= nb_batches:
            c = 0
        new_seed = random.randint(0, sys.maxint)
        random.seed(c+123456789+seed)
        for b in range(batch_size):
            t = random.randint(0,len(Xd)-1)

            xd = Xd[t]
            s = random.randint(min(maxlen,len(xd)), max(maxlen,len(xd)))
            xds.append(xd[:s])

            xh = Xh[t]
            s = random.randint(min(maxlen,len(xh)), max(maxlen,len(xh)))
            xhs.append(xh[:s])

        # undo the seeding before we yield inorder not to affect the caller
        c+= 1
        random.seed(new_seed)

        yield conv_seq_labels(xds, xhs, nflips=nflips, model=model, debug=debug)

```

In [51]:

Out [51] :

((64, 50), (64, 25, 40000), 2)

In [52]:

```
def test_gen(gen, n=5):
    Xtr,Ytr = next(gen)
    for i in range(n):
        assert Xtr[i,maxlend] == eos
        x = Xtr[i,:maxlend]
        y = Xtr[i,maxlend:]
        yy = Ytr[i,:]
        yy = np.where(yy)[1]
        prt('L',yy)
        prt('H',y)
        if maxlend:
            prt('D',x)
```

In [53]:

```
test_gen(gen(X_train, Y_train, batch_size=batch_size))
```

L: Rockets GM Daryl O'Keefe Has A implausible And geometry Reason For Giving Money To Mitt Romney ~
H: ~ Rockets GM Daryl O'Keefe Has A implausible And geometry Reason For Giving Money To Mitt Romne Y
D: personal views , but no one really seems to like Mitt for his essential <0>^ . But Daryl O'Keef e , the Houston Rockets general manager
L: Today 's Celebrity Twitter Chatter ~
H: ~ Today 's Celebrity Twitter Chatter
D: Your daily look into the <0>^ lives of our favorite celebs .
L: Natalie Wood Case : Boat Captain Did n't Hear Her Cries For Help (VIDEO) ~
H: ~ Natalie Wood Case : Boat Captain Did n't Hear Her Cries For Help (VIDEO)
D: Many questions still remain in the case surrounding actress Natalie Wood 's mysterious death , but one person who was on the boat is opening
L: The Beautiful Water Cube In Beijing Is Now A Water Theme Park ~
H: ~ The Beautiful Water Cube In Beijing Is Now A Water Theme Park
D: They 're marveled and <0>^ at during the games , but after the festivities , they 're forgotten and left to Rust . China is
L: Real Housewives of New Jersey : The revelation of Kim G . ~
H: ~ Real Housewives of New Jersey : The revelation of Kim G .
D: that the fights and clashes on the Real Housewives of New Jersey come out of nowhere , but they have all been resurrection by a

test fliping

In [54]:

```
test gen(gen(X_train, Y_train, nflips=6, model=model, debug=False, batch_size=batch_size))
```

L: Rockets GM Daryl O'Keefe Has A implausible And geometry Reason For Giving Money To Mitt Romney ~
H: ~ The GM Daryl O'Keefe Has A <0>^ And geometry Reason For Giving Him To Mitt Romney _ _ _ _ _
D: personal views , but no one really seems to like Mitt for his essential <0>^ . But Daryl O'Keefe , the Houston Rockets general manager
L: Today 's Celebrity Twitter Chatter ~
H: ~ Today 's Celebrity Twitter Chatter _ _ _ _ _
D: _ _ _ _ _ Your daily look into the <0>^ lives of our favorite celebs .
L: Natalie Wood Case : Boat Captain Did n't Hear Her Cries For Help (VIDEO) ~
H: ~ Natalie Wood : : Natalie Captain ~ n't Hear Her Cries For Help (VIDEO) _ _ _ _ _
D: Many questions still remain in the case surrounding actress Natalie Wood 's mysterious death , but one person who was on the boat is opening
L: The Beautiful Water Cube In Beijing Is Now A Water Theme Park ~
H: ~ The World Water Cube In Beijing Is Now A ~ Theme Park _ _ _ _ _
D: They 're marveled and <0>^ at during the games , but after the festivities , they 're forgotten and left to Rust . China is

D: that the fights and crashes on the real housewives of new jersey come out of nowhere , but they have all been resurrection by a

In [55]:

```
valgen = gen(X_test, Y_test,nb_batches=3, batch_size=batch_size)
```

check that valgen repeats itself after nb_batches

In [56]:

```
for i in range(4):
    test_gen(valgen, n=1)
```

L: Dell Releases Draft 802.11n Wireless Card ~ -----
H: ~ Dell Releases Draft 802.11n Wireless Card -----
D: the all the problems that draft N has been facing , we 're glad they 're releasing this as an W
hatsApp card and not embedding
L: Can You Identify a Mystery Cocoon Which Has the Whole Internet stumped ? ~ -----
H: ~ Can You Identify a Mystery Cocoon Which Has the Whole Internet stumped ? -----
D: Usually , throw the internet an image of something you ca n't identify and it 's only a couple
of minutes until you 're bombarded
L: <0>^ ~ -----
H: ~ <0>^ ~ -----
D: Whoa whoa whoa . Wait a second . You 're telling me there 's no liquid cooling ?
L: Dell Releases Draft 802.11n Wireless Card ~ -----
H: ~ Dell Releases Draft 802.11n Wireless Card -----
D: the all the problems that draft N has been facing , we 're glad they 're releasing this as an W
hatsApp card and not embedding

2.6 Train

In [57]:

```
history = {}
```

In [58]:

```
traingen = gen(X_train, Y_train, batch_size=batch_size, nflips=nflips, model=model)
valgen = gen(X_test, Y_test, nb_batches=nb_val_samples//batch_size, batch_size=batch_size)
```

In [59]:

```
r = next(traingen)
r[0].shape, r[1].shape, len(r)
```

Out[59]:

```
((64, 50), (64, 25, 40000), 2)
```

In []:

```
for iteration in range(500):
    print 'Iteration', iteration
    h = model.fit_generator(traingen, samples_per_epoch=nb_train_samples,
                           nb_epoch=1, validation_data=valgen, nb_val_samples=nb_val_samples
                           )
    for k,v in h.history.iteritems():
        history[k] = history.get(k,[]) + v
    with open('data/%s.history.pkl'%FN,'wb') as fp:
        pickle.dump(history,fp,-1)
    model.save_weights('data/%s.hdf5'%FN, overwrite=True)
    gensamples(batch_size=batch_size)
```