

Neural Machine Translation

Welcome to your first programming assignment for this week!

You will build a Neural Machine Translation (NMT) model to translate human readable dates ("25th of June, 2009") into machine readable dates ("2009-06-25"). You will do this using an attention model, one of the most sophisticated sequence to sequence models.

This notebook was produced together with NVIDIA's Deep Learning Institute.

Let's load all the packages you will need for this assignment.

```
In [1]: from keras.layers import Bidirectional, Concatenate, Permute, Dot, Input, LSTM, Multiply
        from keras.layers import RepeatVector, Dense, Activation, Lambda
        from keras.optimizers import Adam
        from keras.utils import to_categorical
        from keras.models import load_model, Model
        import keras.backend as K
        import numpy as np

        from faker import Faker
        import random
        from tqdm import tqdm
        from babel.dates import format_date
        from nmt_utils import *
        import matplotlib.pyplot as plt
        %matplotlib inline
```

1 - Translating human readable dates into machine readable dates

The model you will build here could be used to translate from one language to another, such as translating from English to Hindi. However, language translation requires massive datasets and usually takes days of training on GPUs. To give you a place to experiment with these models even without using massive datasets, we will instead use a simpler "date translation" task.

The network will input a date written in a variety of possible formats (e.g. "the 29th of August 1958", "03/30/1968", "24 JUNE 1987") and translate them into standardized, machine readable dates (e.g. "1958-08-29", "1968-03-30", "1987-06-24"). We will have the network learn to output dates in the common machine-readable format YYYY-MM-DD.

1.1 - Dataset

We will train the model on a dataset of 10000 human readable dates and their equivalent, standardized, machine readable dates. Let's run the following cells to load the dataset and print some examples.

```
In [2]: m = 10000
dataset, human_vocab, machine_vocab, inv_machine_vocab = load_dataset(m)
```

```
In [3]: dataset[:10]
```

```
Out[3]: [('9 may 1998', '1998-05-09'),
 ('10.09.70', '1970-09-10'),
 ('4/28/90', '1990-04-28'),
 ('thursday january 26 1995', '1995-01-26'),
 ('monday march 7 1983', '1983-03-07'),
 ('sunday may 22 1988', '1988-05-22'),
 ('tuesday july 8 2008', '2008-07-08'),
 ('08 sep 1999', '1999-09-08'),
 ('1 jan 1981', '1981-01-01'),
 ('monday may 22 1995', '1995-05-22')]
```

You've loaded:

- `dataset`: a list of tuples of (human readable date, machine readable date)
- `human_vocab`: a python dictionary mapping all characters used in the human readable dates to an integer-valued index
- `machine_vocab`: a python dictionary mapping all characters used in machine readable dates to an integer-valued index. These indices are not necessarily consistent with `human_vocab`.
- `inv_machine_vocab`: the inverse dictionary of `machine_vocab`, mapping from indices back to characters.

Let's preprocess the data and map the raw text data into the index values. We will also use `Tx=30` (which we assume is the maximum length of the human readable date; if we get a longer input, we would have to truncate it) and `Ty=10` (since "YYYY-MM-DD" is 10 characters long).

```
In [4]: Tx = 30
        Ty = 10
        X, Y, Xoh, Yoh = preprocess_data(dataset, human_vocab, machine_vocab,
        Tx, Ty)

        print("X.shape:", X.shape)
        print("Y.shape:", Y.shape)
        print("Xoh.shape:", Xoh.shape)
        print("Yoh.shape:", Yoh.shape)

X.shape: (10000, 30)
Y.shape: (10000, 10)
Xoh.shape: (10000, 30, 37)
Yoh.shape: (10000, 10, 11)
```

You now have:

- X: a processed version of the human readable dates in the training set, where each character is replaced by an index mapped to the character via `human_vocab`. Each date is further padded to T_x values with a special character (< pad >). `X.shape = (m, Tx)`
- Y: a processed version of the machine readable dates in the training set, where each character is replaced by the index it is mapped to in `machine_vocab`. You should have `Y.shape = (m, Ty)`.
- Xoh: one-hot version of X, the "1" entry's index is mapped to the character thanks to `human_vocab`. `Xoh.shape = (m, Tx, len(human_vocab))`
- Yoh: one-hot version of Y, the "1" entry's index is mapped to the character thanks to `machine_vocab`. `Yoh.shape = (m, Ty, len(machine_vocab))`. Here, `len(machine_vocab) = 11` since there are 11 characters ('-' as well as 0-9).

Lets also look at some examples of preprocessed training examples. Feel free to play with `index` in the cell below to navigate the dataset and see how source/target dates are preprocessed.

```
Source after preprocessing (indices): [12  0 24 13 34  0  4 12 12 11
36 36 36 36 36 36 36 36 36 36 36 36 36 36 36
 36 36 36 36 36]
Target after preprocessing (indices): [ 2 10 10  9  0  1  6  0  1 10]

Source after preprocessing (one-hot): [[ 0.  0.  0. ...,  0.  0.  0.]
[ 1.  0.  0. ...,  0.  0.  0.]
[ 0.  0.  0. ...,  0.  0.  0.]
...,
[ 0.  0.  0. ...,  0.  0.  1.]
[ 0.  0.  0. ...,  0.  0.  1.]
[ 0.  0.  0. ...,  0.  0.  1.]]
Target after preprocessing (one-hot): [[ 0.  0.  1.  0.  0.  0.  0.
0.  0.  0.  0.]
[ 0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  1.]
[ 0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  1.]
[ 1.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.]
[ 0.  1.  0.  0.  0.  0.  0.  0.  0.  0.  0.]
[ 0.  0.  0.  0.  0.  0.  1.  0.  0.  0.  0.]
[ 1.  0.  0.  0.  0.  0.  0.  0.  0.  0.  0.]
[ 0.  1.  0.  0.  0.  0.  0.  0.  0.  0.  0.]
[ 0.  0.  0.  0.  0.  0.  0.  0.  0.  0.  1.]]
```

2 - Neural machine translation with attention

If you had to translate a book's paragraph from French to English, you would not read the whole paragraph, then close the book and translate. Even during the translation process, you would read/re-read and focus on the parts of the French paragraph corresponding to the parts of the English you are writing down.

The attention mechanism tells a Neural Machine Translation model where it should pay attention to at any step.

2.1 - Attention mechanism

In this part, you will implement the attention mechanism presented in the lecture videos. Here is a figure to remind you how the model works. The diagram on the left shows the attention model. The diagram on the right shows what one "Attention" step does to calculate the attention variables $\alpha^{(t,t')}$, which are used to compute the context variable $context^{(t)}$ for each timestep in the output ($t = 1, \dots, T_y$).

</table>

Here are some properties of the model that you may notice:

- There are two separate LSTMs in this model (see diagram on the left). Because the one at the bottom of the picture is a Bi-directional LSTM and comes *before* the attention mechanism, we will call it *pre-attention* Bi-LSTM. The LSTM at the top of the diagram comes *after* the attention mechanism, so we will call it the *post-attention* LSTM. The pre-attention Bi-LSTM goes through T_x time steps; the post-attention LSTM goes through T_y time steps.
- The post-attention LSTM passes $s^{(t)}$, $c^{(t)}$ from one time step to the next. In the lecture videos, we were using only a basic RNN for the post-activation sequence model, so the state captured by the RNN output activations $s^{(t)}$. But since we are using an LSTM here, the LSTM has both the output activation $s^{(t)}$ and the hidden cell state $c^{(t)}$. However, unlike previous text generation examples (such as Dinosaur in week 1), in this model the post-activation LSTM at time t does not take the specific generated $y^{(t-1)}$ as input; it only takes $s^{(t)}$ and $c^{(t)}$ as input. We have designed the model this way, because (unlike language generation where adjacent characters are highly correlated) there isn't as strong a dependency between the previous character and the next character in a YYYY-MM-DD date.
- We use $a^{(t)} = [\vec{a}^{(t)}; \overleftarrow{a}^{(t)}]$ to represent the concatenation of the activations of both the forward-direction and backward-directions of the pre-attention Bi-LSTM.
- The diagram on the right uses a RepeatVector node to copy $s^{(t-1)}$'s value T_x times, and then Concatenation to concatenate $s^{(t-1)}$ and $a^{(t)}$ to compute $e^{(t,t')}$, which is then passed through a softmax to compute $\alpha^{(t,t')}$. We'll explain how to use RepeatVector and Concatenation in Keras below.

Lets implement this model. You will start by implementing two functions: `one_step_attention()` and `model()`.

1) one_step_attention(): At step t , given all the hidden states of the Bi-LSTM ($[a^{(1)}, a^{(2)}, \dots, a^{(T_x)}]$) and the previous hidden state of the second LSTM ($s^{(t)}$), `one_step_attention()` will compute the attention weights ($[\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(T_x)}]$) and output the context vector (see Figure 1 (right) for details): $context^{(t)} = \sum_{t'=0}^{T_x} \alpha^{(t,t')} a^{(t')}$

Note that we are denoting the attention in this notebook $context^{(t)}$. In the lecture videos, the context was denoted $c^{(t)}$, but here we are calling it $context^{(t)}$ to avoid confusion with the (post-attention) LSTM's internal memory cell variable, which is sometimes also denoted $c^{(t)}$.

2) model(): Implements the entire model. It first runs the input through a Bi-LSTM to get back $[a^{(1)}, a^{(2)}, \dots, a^{(T_x)}]$. Then, it calls `one_step_attention()` T_y times (for loop). At each iteration of this loop, it gives the computed context vector $c^{(t)}$ to the second LSTM, and runs the output of the LSTM through a dense layer with softmax activation to generate

Exercise: Implement `one_step_attention()`. The function `model()` will call the layers in `one_step_attention()` T_y using a for-loop, and it is important that all T_y copies have the same weights. I.e., it should not re-initialize the weights every time. In other words, all T_y steps should have shared weights. Here's how you can implement layers with shareable weights in Keras:

1. Define the layer objects (as global variables for examples).
2. Call these objects when propagating the input.

We have defined the layers you need as global variables. Please run the following cells to create them. Please check the Keras documentation to make sure you understand what these layers are: `RepeatVector()` (<https://keras.io/layers/core/#repeatvector>), `Concatenate()` (<https://keras.io/layers/merge/#concatenate>), `Dense()` (<https://keras.io/layers/core/#dense>), `Activation()` (<https://keras.io/layers/core/#activation>), `Dot()` (<https://keras.io/layers/merge/#dot>).

In [6]:

Now you can use these layers to implement `one_step_attention()`. In order to propagate a Keras tensor object `X` through one of these layers, use `layer(X)` (or `layer([X,Y])` if it requires multiple inputs.), e.g. `dense(X)` will propagate `X` through the `Dense(1)` layer defined above.

In [51]:

You will be able to check the expected output of `one_step_attention()` after you've coded the `model()` function.

Exercise: Implement `model()` as explained in figure 2 and the text above. Again, we have defined global layers that will share weights to be used in `model()`.

In [52]:

Now you can use these layers T_y times in a for loop to generate the outputs, and their parameters will not be reinitialized. You will have to carry out the following steps:

1. Propagate the input into a `Bidirectional LSTM` (<https://keras.io/layers/wrappers/#bidirectional>) (<https://keras.io/layers/recurrent/#lstm>)
2. Iterate for $t = 0, \dots, T_y - 1$:
 - A. Call `one_step_attention()` on $\{\alpha^0, \alpha^1, \dots, \alpha^{T_y-1}\}$ and s^t to get the context vector $context^t$.
 - B. Give $context^t$ to the post-attention LSTM cell. Remember pass in the previous hidden states $\langle s^{t-1} \rangle$ and cell states $\langle c^{t-1} \rangle$ of this LSTM using `initial_state=[previous hidden state, previous cell state]`. Get back the new hidden state s^t and the new cell state c^t .
 - C. Apply a softmax layer to s^t , get the output.
 - D. Save the output by adding it to the list of outputs.
3. Create your Keras model instance, it should have three inputs ("inputs", $s^{<0>}$ and $c^{<0>}$) and output the list of "outputs".

In [59]:

Run the following cell to create your model.

In [60]:

```
30
  (?, 30, 37)
Tensor("bidirectional_20/concat_2:0", shape=(?, ?, 64), dtype=float32)
  (?, 30, 64)
  (?, ?, 64)
  (?, 30, 64)
  (?, ?, 64)
  (?, 30, 64)
  (?, ?, 64)
  (?, 30, 64)
  (?, ?, 64)
  (?, 30, 64)
  (?, ?, 64)
  (?, 30, 64)
  (?, ?, 64)
  (?, 30, 64)
  (?, ?, 64)
  (?, 30, 64)
  (?, ?, 64)
  (?, 30, 64)
  (?, ?, 64)
  (?, 30, 64)
  (?, ?, 64)
```

Let's get a summary of the model to check if it matches the expected output.

In [61]:

Layer (type) Connected to	Output Shape	Param #
=====		
input_22 (InputLayer)	(None, 30, 37)	0
=====		
s0 (InputLayer)	(None, 64)	0
=====		
bidirectional_20 (Bidirectional) input_22[0][0]	(None, 30, 64)	17920
=====		
repeat_vector_1 (RepeatVector) s0[0][0]	(None, 30, 64)	0
lstm_21[11][0]		
lstm_21[12][0]		
lstm_21[13][0]		
lstm_21[14][0]		
lstm_21[15][0]		
lstm_21[16][0]		
lstm_21[17][0]		
lstm_21[18][0]		
lstm_21[19][0]		
=====		
concatenate_1 (Concatenate) bidirectional_20[0][0]	(None, 30, 128)	0
repeat_vector_1[28][0]		
bidirectional_20[0][0]		
repeat_vector_1[29][0]		
bidirectional_20[0][0]		
repeat_vector_1[30][0]		
bidirectional_20[0][0]		
repeat_vector_1[31][0]		

bidirectional_20[0][0]

repeat_vector_1[32][0]

bidirectional_20[0][0]

repeat_vector_1[33][0]

bidirectional_20[0][0]

repeat_vector_1[34][0]

bidirectional_20[0][0]

repeat_vector_1[35][0]

bidirectional_20[0][0]

repeat_vector_1[36][0]

bidirectional_20[0][0]

repeat_vector_1[37][0]

dense_1 (Dense)	(None, 30, 10)	1290
-----------------	----------------	------

concatenate_1[15][0]

concatenate_1[16][0]

concatenate_1[17][0]

concatenate_1[18][0]

concatenate_1[19][0]

concatenate_1[20][0]

concatenate_1[21][0]

concatenate_1[22][0]

concatenate_1[23][0]

concatenate_1[24][0]

dense_2 (Dense)	(None, 30, 1)	11
-----------------	---------------	----

dense_1[15][0]

dense_1[16][0]

dense_1[17][0]

dense_1[18][0]

dense_1[19][0]

dense_1[20][0]

dense_1[21][0]

dense_1[22][0]

dense_1[23][0]

dense_1[24][0]

attention_weights (Activation)	(None, 30, 1)	0
--------------------------------	---------------	---

dense_2[15][0]

dense_2[16][0]

dense_2[17][0]

dense_2[18][0]

dense_2[19][0]

dense_2[20][0]

dense_2[21][0]

dense_2[22][0]

dense_2[23][0]

dense_2[24][0]

dot_1 (Dot)	(None, 1, 64)	0
-------------	---------------	---

attention_weights[15][0]

bidirectional_20[0][0]

attention_weights[16][0]

bidirectional_20[0][0]

attention_weights[17][0]

bidirectional_20[0][0]

attention_weights[18][0]

bidirectional_20[0][0]

attention_weights[19][0]

bidirectional_20[0][0]

attention_weights[20][0]

bidirectional_20[0][0]

attention_weights[21][0]

bidirectional_20[0][0]

attention_weights[22][0]

bidirectional_20[0][0]

attention_weights[23][0]

bidirectional_20[0][0]

attention_weights[24][0]

bidirectional_20[0][0]

c0 (InputLayer)

(None, 64)

0

lstm_21 (LSTM)

[(None, 64), (None, 64)]

dot_1[14][0]

s0[0][0]

c0[0][0]

dot_1[15][0]

lstm_21[11][0]

lstm_21[11][2]

dot_1[16][0]

lstm_21[12][0]

lstm_21[12][2]

dot_1[17][0]

lstm_21[13][0]

lstm_21[13][2]

dot_1[18][0]

lstm_21[14][0]

lstm_21[14][2]

dot_1[19][0]

lstm_21[15][0]

lstm_21[15][2]

dot_1[20][0]

lstm_21[16][0]

lstm_21[16][2]

dot_1[21][0]

lstm_21[17][0]

lstm_21[17][2]

dot_1[22][0]

lstm_21[18][0]

lstm_21[18][2]

dot_1[23][0]

lstm_21[19][0]

lstm_21[19][2]

dense_7 (Dense)

(None, 11)

715

lstm_21[11][0]

lstm_21[12][0]

lstm_21[13][0]

lstm_21[14][0]

lstm_21[15][0]

lstm_21[16][0]

lstm_21[17][0]

lstm_21[18][0]

lstm_21[19][0]

lstm_21[20][0]

=====

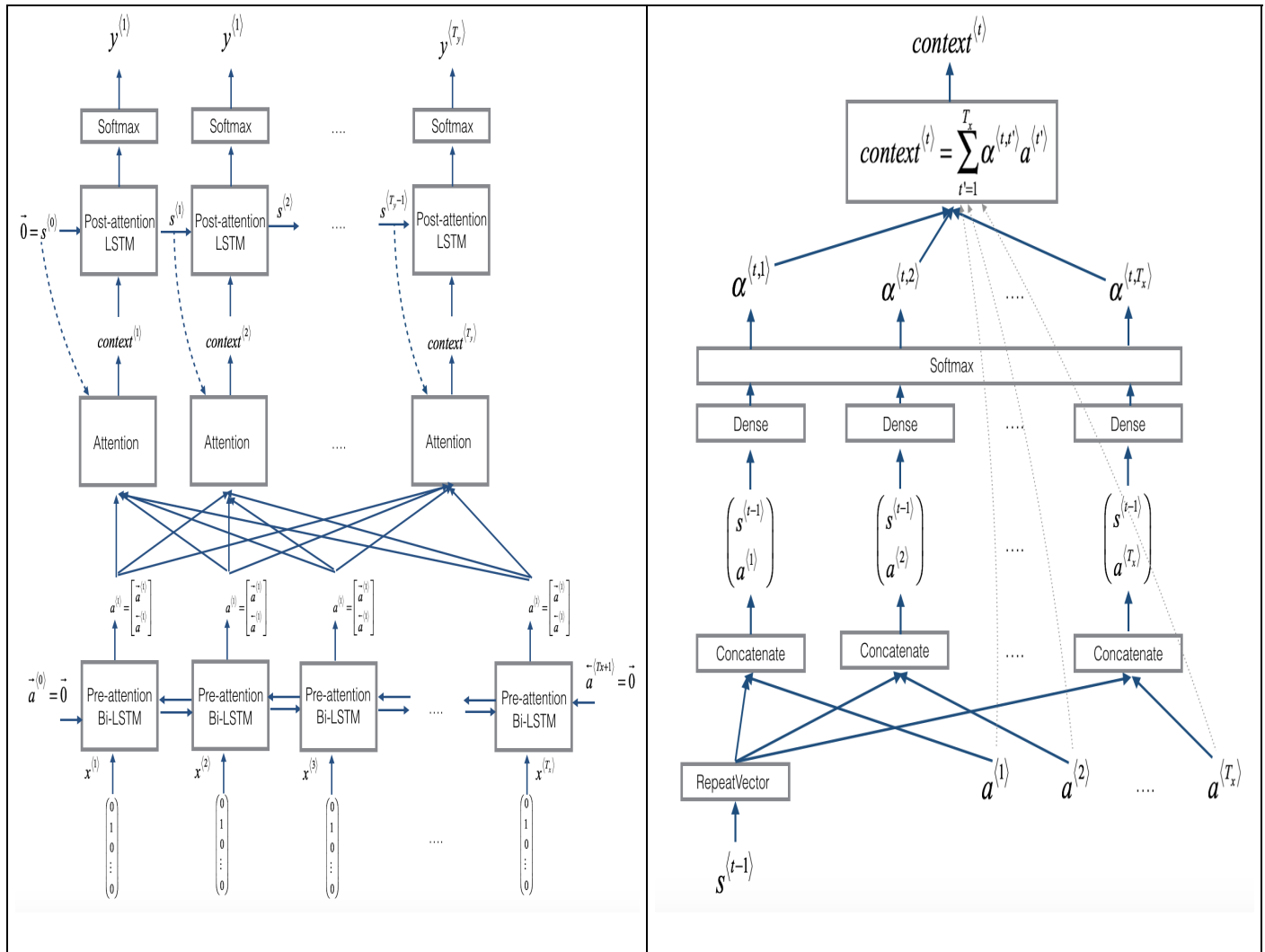
Total params: 52,960

Trainable params: 52,960

Non-trainable params: 0

Expected Output:

Here is the summary you should see

****Figure 1****: Neural machine translation with attention

Total params:	52,960
Trainable params:	52,960
Non-trainable params:	0
**bidirectional_1's output shape **	(None, 30, 64)
**repeat_vector_1's output shape **	(None, 30, 64)
**concatenate_1's output shape **	(None, 30, 128)
**attention_weights's output shape **	(None, 30, 1)
**dot_1's output shape **	(None, 1, 64)
**dense_3's output shape **	(None, 11)

As usual, after creating your model in Keras, you need to compile it and define what loss, optimizer and metrics you want to use. Compile your model using `categorical_crossentropy` loss, a custom [Adam](https://keras.io/optimizers/#adam) (<https://keras.io/optimizers/#usage-of-optimizers>) (learning rate = 0.005, $\beta_1 = 0.9$, $\beta_2 = 0.999$, decay = 0.01) and ['accuracy'] metrics:

In []:

```
In [66]: ### START CODE HERE ### (~2 lines)
opt = Adam(beta_1=0.9,beta_2=0.999,decay=0.01,lr=0.005,)
model.compile(loss='categorical_crossentropy',optimizer=opt,metrics=[
'accuracy'])
### END CODE HERE ###
```

The last step is to define all your inputs and outputs to fit the model:

- You already have X of shape ($m = 10000$, $T_x = 30$) containing the training examples.
- You need to create `s0` and `c0` to initialize your `post_activation_LSTM_cell` with 0s.
- Given the `model()` you coded, you need the "outputs" to be a list of 11 elements of shape (m , T_y). So that: `outputs[i][0]`, ..., `outputs[i][Ty]` represent the true labels (characters) corresponding to the i^{th} training example ($X[i]$). More generally, `outputs[i][j]` is the true label of the j^{th} character in the i^{th} training example.

```
In [67]: s0 = np.zeros((m, n_s))
c0 = np.zeros((m, n_s))
outputs = list(Yoh.swapaxes(0,1))
```

Let's now fit the model and run it for one epoch.

```
In [68]: model.fit([Xoh, s0, c0], outputs, epochs=1, batch_size=100)
```

```
Epoch 1/1
10000/10000 [=====] - 39s - loss: 16.6293 -
dense_7_loss_1: 1.2717 - dense_7_loss_2: 1.0764 - dense_7_loss_3: 1.7
963 - dense_7_loss_4: 2.6848 - dense_7_loss_5: 0.7456 - dense_7_loss_
6: 1.2381 - dense_7_loss_7: 2.6284 - dense_7_loss_8: 0.9221 - dense_7
_loss_9: 1.6960 - dense_7_loss_10: 2.5699 - dense_7_acc_1: 0.4799 - d
ense_7_acc_2: 0.6420 - dense_7_acc_3: 0.3003 - dense_7_acc_4: 0.0911
- dense_7_acc_5: 0.9086 - dense_7_acc_6: 0.3890 - dense_7_acc_7: 0.06
29 - dense_7_acc_8: 0.9523 - dense_7_acc_9: 0.2510 - dense_7_acc_10:
0.1005
```

```
Out[68]: <keras.callbacks.History at 0x7fdbcea39a58>
```


While training you can see the loss as well as the accuracy on each of the 10 positions of the output. The table below gives you an example of what the accuracies could be if the batch had 2 examples:

True labels 1	1	9	9	5	-	1	2	-	0	4
Predictions 1	1	9	9	5	-	1	0	-	0	5
True labels 1	1	9	6	8	-	0	1	-	0	4
Predictions 2	1	9	7	8	-	0	3	-	0	4
Index	1	2	3	4	5	6	7	8	9	10
Accuracy	1.0	1.0	0.5	1.0	1.0	1.0	0.0	1.0	1.0	0.5

Thus, `dense_2_acc_8: 0.89` means that you are predicting the 7th character of the output correctly 89% of the time in the current batch of data.

We have run this model for longer, and saved the weights. Run the next cell to load our weights. (By training a model for several minutes, you should be able to obtain a model of similar accuracy, but loading our model will save you time.)

```
In [69]: model.load_weights('models/model.h5')
```

You can now see the results on new examples.

```
In [70]: EXAMPLES = ['3 May 1979', '5 April 09', '21th of August 2016', 'Tue 1
0 Jul 2007', 'Saturday May 9 2018', 'March 3 2001', 'March 3rd 2001',
'1 March 2001']
for example in EXAMPLES:

    source = string_to_int(example, Tx, human_vocab)
    source = np.array(list(map(lambda x: to_categorical(x, num_classes=len(human_vocab)), source))).swapaxes(0,1)
    prediction = model.predict([source, s0, c0])
    prediction = np.argmax(prediction, axis = -1)
    output = [inv_machine_vocab[int(i)] for i in prediction]

    print("source:", example)
    print("output:", ''.join(output))

source: 3 May 1979
output: 1979-05-03
source: 5 April 09
output: 2009-05-05
source: 21th of August 2016
output: 2016-08-21
source: Tue 10 Jul 2007
output: 2007-07-10
source: Saturday May 9 2018
output: 2018-05-09
source: March 3 2001
output: 2001-03-03
source: March 3rd 2001
output: 2001-03-03
source: 1 March 2001
output: 2001-03-01
```

You can also change these examples to test with your own examples. The next part will give you a better sense on what the attention mechanism is doing--i.e., what part of the input the network is paying attention to when generating a particular output character.

3 - Visualizing Attention (Optional / Ungraded)

Since the problem has a fixed output length of 10, it is also possible to carry out this task using 10 different softmax units to generate the 10 characters of the output. But one advantage of the attention model is that each part of the output (say the month) knows it needs to depend only on a small part of the input (the characters in the input giving the month). We can visualize what part of the output is looking at what part of the input.

Consider the task of translating "Saturday 9 May 2018" to "2018-05-09". If we visualize the computed $\alpha^{(t,t')}$ we get this:

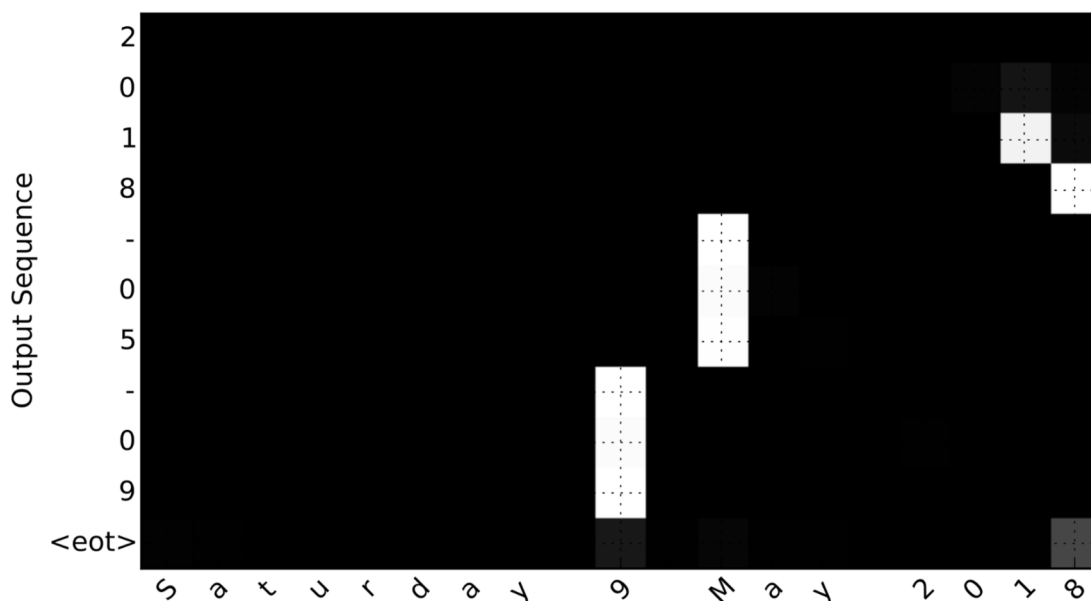


Figure 8: Full Attention Map

Notice how the output ignores the "Saturday" portion of the input. None of the output timesteps are paying much attention to that portion of the input. We see also that 9 has been translated as 09 and May has been correctly translated into 05, with the output paying attention to the parts of the input it needs to to make the translation. The year mostly requires it to pay attention to the input's "18" in order to generate "2018."

3.1 - Getting the activations from the network

Lets now visualize the attention values in your network. We'll propagate an example through the network, then visualize the values of $\alpha^{(t,t')}$.

To figure out where the attention values are located, let's start by printing a summary of the model .

```
In [ ]: model.summary()
```