

Loading the Data

We will be working with the Babi Data Set from Facebook Research.

Full Details: <https://research.fb.com/downloads/babi/>
(<https://research.fb.com/downloads/babi/>).

- Jason Weston, Antoine Bordes, Sumit Chopra, Tomas Mikolov, Alexander M. Rush, "Towards AI-Complete Question Answering: A Set of Prerequisite Toy Tasks", <http://arxiv.org/abs/1502.05698> (<http://arxiv.org/abs/1502.05698>).

```
In [1]: 1 import pickle
        2 import numpy as np
```

```
In [2]: 1 with open("train_qa.txt", "rb") as fp:  # Unpickling
        2     train_data = pickle.load(fp)
```

```
In [3]: 1 with open("test_qa.txt", "rb") as fp:  # Unpickling
        2     test_data = pickle.load(fp)
```

Exploring the Format of the Data

```
In [4]: 1 type(test_data)
```

```
Out[4]: list
```

```
In [5]: 1 type(train_data)
```

```
Out[5]: list
```

```
In [6]: 1 len(test_data)
```

```
Out[6]: 1000
```

```
In [7]: 1 len(train_data)
```

```
Out[7]: 10000
```

```
In [8]: 1 train_data[0]
```

```
Out[8]: (['Mary',  
         'moved',  
         'to',  
         'the',  
         'bathroom',  
         '.',  
         'Sandra',  
         'journeyed',  
         'to',  
         'the',  
         'bedroom',  
         '.'],  
 ['Is', 'Sandra', 'in', 'the', 'hallway', '?'],  
 'no')
```

```
In [9]: 1 ' '.join(train_data[0][0])
```

```
Out[9]: 'Mary moved to the bathroom . Sandra journeyed to the bedroom .'
```

```
In [10]: 1 ' '.join(train_data[0][1])
```

```
Out[10]: 'Is Sandra in the hallway ?'
```

```
In [11]: 1 train_data[0][2]
```

```
Out[11]: 'no'
```

Setting up Vocabulary of All Words

```
In [12]: 1 # Create a set that holds the vocab words  
2 vocab = set()
```

```
In [13]: 1 all_data = test_data + train_data
```

```
In [14]: 1 for story, question , answer in all_data:  
2     # In case you don't know what a union of sets is:  
3     # https://www.programiz.com/python-programming/methods/set/union  
4     vocab = vocab.union(set(story))  
5     vocab = vocab.union(set(question))
```

```
In [15]: 1 vocab.add('no')  
2 vocab.add('yes')
```

```
In [16]: 1 vocab
```

```
Out[16]: {'.',  
'?',  
'Daniel',  
'Is',  
'John',  
'Mary',  
'Sandra',  
'apple',  
'back',  
'bathroom',  
'bedroom',  
'discarded',  
'down',  
'dropped',  
'football',  
'garden',  
'got',  
'grabbed',  
'hallway',  
'in',  
'journeyed',  
'kitchen',  
'left',  
'milk',  
'moved',  
'no',  
'office',  
'picked',  
'put',  
'the',  
'there',  
'to',  
'took',  
'travelled',  
'up',  
'went',  
'yes'}
```

```
In [17]: 1 vocab_len = len(vocab) + 1 #we add an extra space to hold a 0 for
```

```
In [18]: 1 max_story_len = max([len(data[0]) for data in all_data])
```

```
In [19]: 1 max_story_len
```

```
Out[19]: 156
```

```
In [20]: 1 max_question_len = max([len(data[1]) for data in all_data])
```

```
In [21]: 1 max_question_len
```

```
Out[21]: 6
```

Vectorizing the Data

In [22]:

```
1 vocab
```

Out [22]:

```
{'.',  
'?',  
'Daniel',  
'Is',  
'John',  
'Mary',  
'Sandra',  
'apple',  
'back',  
'bathroom',  
'bedroom',  
'discarded',  
'down',  
'dropped',  
'football',  
'garden',  
'got',  
'grabbed',  
'hallway',  
'in',  
'journeyed',  
'kitchen',  
'left',  
'milk',  
'moved',  
'no',  
'office',  
'picked',  
'put',  
'the',  
'there',  
'to',  
'took',  
'travelled',  
'up',  
'went',  
'yes'}
```

In [23]:

```
1 # Reserve 0 for pad_sequences  
2 vocab_size = len(vocab) + 1
```

In [24]:

```
1 from keras.preprocessing.sequence import pad_sequences  
2 from keras.preprocessing.text import Tokenizer
```

Using TensorFlow backend.

```
In [25]: 1 # integer encode sequences of words
          2 tokenizer = Tokenizer(filters=[])
          3 tokenizer.fit_on_texts(vocab)
```

```
In [26]: 1 tokenizer.word_index
```

```
Out[26]: {'.': 13,
          '?': 12,
          'apple': 3,
          'back': 18,
          'bathroom': 23,
          'bedroom': 29,
          'daniel': 21,
          'discarded': 5,
          'down': 22,
          'dropped': 4,
          'football': 37,
          'garden': 36,
          'got': 34,
          'grabbed': 15,
          'hallway': 27,
          'in': 8,
          'is': 30,
          'john': 10,
          'journeyed': 35,
          'kitchen': 11,
          'left': 31,
          'mary': 25,
          'milk': 7,
          'moved': 28,
          'no': 6,
          'office': 24,
          'picked': 1,
          'put': 33,
          'sandra': 2,
          'the': 32,
          'there': 16,
          'to': 14,
          'took': 9,
          'travelled': 20,
          'up': 17,
          'went': 26,
          'yes': 19}
```

```
In [27]: 1 train_story_text = []
          2 train_question_text = []
          3 train_answers = []
          4
          5 for story, question, answer in train_data:
          6     train_story_text.append(story)
          7     train_question_text.append(question)
```

```
In [28]: 1 train_story_seq = tokenizer.texts_to_sequences(train_story_text)
```

```
In [29]: 1 len(train_story_text)
```

```
Out[29]: 10000
```

```
In [30]: 1 len(train_story_seq)
```

```
Out[30]: 10000
```

```
In [31]: 1 # word_index = tokenizer.word_index
```

Functionalize Vectorization

```

In [32]: 1 def vectorize_stories(data, word_index=tokenizer.word_index, max_s
2         '''
3         INPUT:
4
5         data: consisting of Stories,Queries,and Answers
6         word_index: word index dictionary from tokenizer
7         max_story_len: the length of the longest story (used for pad_s
8         max_question_len: length of the longest question (used for pad
9
10
11        OUTPUT:
12
13        Vectorizes the stories,questions, and answers into padded sequ
14        answer in the data. Then we convert the raw words to an word i
15        output list. Then once we have converted the words to numbers,
16
17        Returns this in the form of a tuple (X,Xq,Y) (padded based on
18        '''
19
20
21        # X = STORIES
22        X = []
23        # Xq = QUERY/QUESTION
24        Xq = []
25        # Y = CORRECT ANSWER
26        Y = []
27
28
29        for story, query, answer in data:
30
31            # Grab the word index for every word in story
32            x = [word_index[word.lower()] for word in story]
33            # Grab the word index for every word in query
34            xq = [word_index[word.lower()] for word in query]
35
36            # Grab the Answers (either Yes/No so we don't need to use
37            # Index 0 is reserved so we're going to use + 1
38            y = np.zeros(len(word_index) + 1)
39
40            # Now that y is all zeros and we know its just Yes/No , we
41            #
42            y[word_index[answer]] = 1
43
44            # Append each set of story,query, and answer to their resp
45            X.append(x)
46            Xq.append(xq)
47            Y.append(y)
48
49            # Finally, pad the sequences based on their max length so the
50
51            # RETURN TUPLE FOR UNPACKING
52            return (pad_sequences(X, maxlen=max_story_len),pad_sequences(Xq, maxlen=max_question_len),Y)

```

```
In [33]: 1 inputs_train, queries_train, answers_train = vectorize_stories(train_data)
```

```
In [34]: 1 inputs_test, queries_test, answers_test = vectorize_stories(test_data)
```

```
In [35]: 1 inputs_test
```

```
Out[35]: array([[ 0,  0,  0, ..., 32, 29, 13],
                [ 0,  0,  0, ..., 32, 36, 13],
                [ 0,  0,  0, ..., 32, 36, 13],
                ...,
                [ 0,  0,  0, ..., 32,  3, 13],
                [ 0,  0,  0, ..., 32, 36, 13],
                [ 0,  0,  0, ...,  3, 16, 13]])
```

```
In [36]: 1 queries_test
```

```
Out[36]: array([[30, 10,  8, 32, 11, 12],
                [30, 10,  8, 32, 11, 12],
                [30, 10,  8, 32, 36, 12],
                ...,
                [30, 25,  8, 32, 29, 12],
                [30,  2,  8, 32, 36, 12],
                [30, 25,  8, 32, 36, 12]])
```

```
In [37]: 1 answers_test
```

```
Out[37]: array([[0., 0., 0., ..., 0., 0., 0.],
                [0., 0., 0., ..., 0., 0., 0.],
                [0., 0., 0., ..., 0., 0., 0.],
                ...,
                [0., 0., 0., ..., 0., 0., 0.],
                [0., 0., 0., ..., 0., 0., 0.],
                [0., 0., 0., ..., 0., 0., 0.]])
```

```
In [38]: 1 sum(answers_test)
```

```
Out[38]: array([ 0.,  0.,  0.,  0.,  0.,  0., 503.,  0.,  0.,  0.,
                0.,  0.,  0.,  0.,  0.,  0.,  0.,  0., 497.,  0.,
                0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,
                0.,  0.,  0.,  0.,  0.] )
```

```
In [39]: 1 tokenizer.word_index['yes']
```

```
Out[39]: 19
```

```
In [40]: 1 tokenizer.word_index['no']
```

```
Out[40]: 6
```


Creating the Model

```
In [41]: 1 from keras.models import Sequential, Model
2 from keras.layers.embeddings import Embedding
3 from keras.layers import Input, Activation, Dense, Permute, Dropout
4 from keras.layers import add, dot, concatenate
5 from keras.layers import LSTM
```

Placeholders for Inputs

Recall we technically have two inputs, stories and questions. So we need to use placeholders. `Input()` is used to instantiate a Keras tensor.

```
In [42]: 1 input_sequence = Input((max_story_len,))
2 question = Input((max_question_len,))
```

Building the Networks

To understand why we chose this setup, make sure to read the paper we are using:

- Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, Rob Fergus, "End-To-End Memory Networks", <http://arxiv.org/abs/1503.08895> (<http://arxiv.org/abs/1503.08895>)

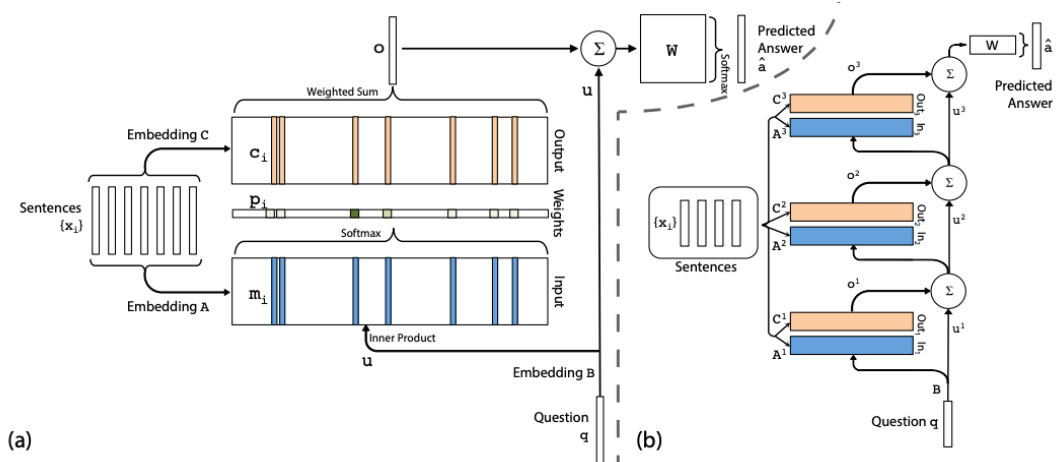


Figure 1: (a): A single layer version of our model. (b): A three layer version of our model. In practice, we can constrain several of the embedding matrices to be the same (see Section 2.2).

Encoders

Input Encoder m

```
In [43]: 1 # Input gets embedded to a sequence of vectors
2 input_encoder_m = Sequential()
3 input_encoder_m.add(Embedding(input_dim=vocab_size,output_dim=64))
4 input_encoder_m.add(Dropout(0.3))
5
6 # This encoder will output:
7 # (samples, story_maxlen, embedding_dim)
```

Input Encoder c

```
In [44]: 1 # embed the input into a sequence of vectors of size query_maxlen
2 input_encoder_c = Sequential()
3 input_encoder_c.add(Embedding(input_dim=vocab_size,
4                               output_dim=max_question_len))
5 input_encoder_c.add(Dropout(0.3))
6 # output: (samples, story_maxlen, query_maxlen)
```

Question Encoder

```
In [45]: 1 # embed the question into a sequence of vectors
2 question_encoder = Sequential()
3 question_encoder.add(Embedding(input_dim=vocab_size,
4                                output_dim=64,
5                                input_length=max_question_len))
6 question_encoder.add(Dropout(0.3))
7 # output: (samples, query_maxlen, embedding_dim)
```

Encode the Sequences

```
In [46]: 1 # encode input sequence and questions (which are indices)
2 # to sequences of dense vectors
3 input_encoded_m = input_encoder_m(input_sequence)
4 input_encoded_c = input_encoder_c(input_sequence)
5 question_encoded = question_encoder(question)
```

Use dot product to compute the match between first input vector seq and the query

```
In [47]: 1 # shape: `(samples, story_maxlen, query_maxlen)`
2 match = dot([input_encoded_m, question_encoded], axes=(2, 2))
3 match = Activation('softmax')(match)
```

Add this match matrix with the second input vector sequence

```
In [48]: 1 # add the match matrix with the second input vector sequence
2 response = add([match, input_encoded_c]) # (samples, story_maxlen, vocab_size)
3 response = Permute((2, 1))(response) # (samples, query_maxlen, vocab_size)
```

Concatenate

```
In [49]: 1 # concatenate the match matrix with the question vector sequence
2 answer = concatenate([response, question_encoded])
```

```
In [50]: 1 answer
```

```
Out[50]: <tf.Tensor 'concatenate_1/concat:0' shape=(?, 6, 220) dtype=float32>
```

```
In [51]: 1 # Reduce with RNN (LSTM)
2 answer = LSTM(32)(answer) # (samples, 32)
```

```
In [52]: 1 # Regularization with Dropout
2 answer = Dropout(0.5)(answer)
3 answer = Dense(vocab_size)(answer) # (samples, vocab_size)
```

```
In [53]: 1 # we output a probability distribution over the vocabulary
2 answer = Activation('softmax')(answer)
3
4 # build the final model
5 model = Model([input_sequence, question], answer)
6 model.compile(optimizer='rmsprop', loss='categorical_crossentropy',
7               metrics=['accuracy'])
```

In [54]:

1	<code>model.summary()</code>
---	------------------------------

Layer (type) connected to	Output Shape	Param #	Con
input_1 (InputLayer)	(None, 156)	0	
input_2 (InputLayer)	(None, 6)	0	
sequential_1 (Sequential) ut_1[0][0]	multiple	2432	inp
sequential_3 (Sequential) ut_2[0][0]	(None, 6, 64)	2432	inp
dot_1 (Dot) quential_1[1][0] quential_3[1][0]	(None, 156, 6)	0	seq seq
activation_1 (Activation) _1[0][0]	(None, 156, 6)	0	dot
sequential_2 (Sequential) ut_1[0][0]	multiple	228	inp
add_1 (Add) ivation_1[0][0] quential_2[1][0]	(None, 156, 6)	0	act seq
permute_1 (Permute) _1[0][0]	(None, 6, 156)	0	add
concatenate_1 (Concatenate) mute_1[0][0] quential_3[1][0]	(None, 6, 220)	0	per seq
lstm_1 (LSTM) catenate_1[0][0]	(None, 32)	32384	con
dropout_4 (Dropout) m_1[0][0]	(None, 32)	0	lst
dense_1 (Dense) pout_4[0][0]	(None, 38)	1254	dro

```

activation_2 (Activation)      (None, 38)      0      den
se_1[0][0]
=====
=====
Total params: 38,730
Trainable params: 38,730
Non-trainable params: 0

```

```

In [55]: # train
history = model.fit([inputs_train, queries_train], answers_train, batch_size=100, epochs=120, validation_data=([inputs_val, queries_val], answers_val), verbose=1)

0.1131 - acc: 0.9560 - val_loss: 0.2496 - val_acc: 0.9180
Epoch 115/120
10000/10000 [=====] - 4s 364us/step - loss: 0.1050 - acc: 0.9599 - val_loss: 0.3021 - val_acc: 0.9170
Epoch 116/120
10000/10000 [=====] - 3s 342us/step - loss: 0.1038 - acc: 0.9619 - val_loss: 0.2673 - val_acc: 0.9160
Epoch 117/120
10000/10000 [=====] - 4s 350us/step - loss: 0.1100 - acc: 0.9599 - val_loss: 0.2855 - val_acc: 0.9160
Epoch 118/120
10000/10000 [=====] - 3s 349us/step - loss: 0.1012 - acc: 0.9610 - val_loss: 0.2887 - val_acc: 0.9200
Epoch 119/120
10000/10000 [=====] - 4s 352us/step - loss: 0.1040 - acc: 0.9618 - val_loss: 0.2774 - val_acc: 0.9120
Epoch 120/120
10000/10000 [=====] - 3s 349us/step - loss: 0.1010 - acc: 0.9615 - val_loss: 0.3139 - val_acc: 0.9120

```

Saving the Model

```

In [72]: 1 filename = 'chatbot_120_epochs.h5'
          2 model.save(filename)

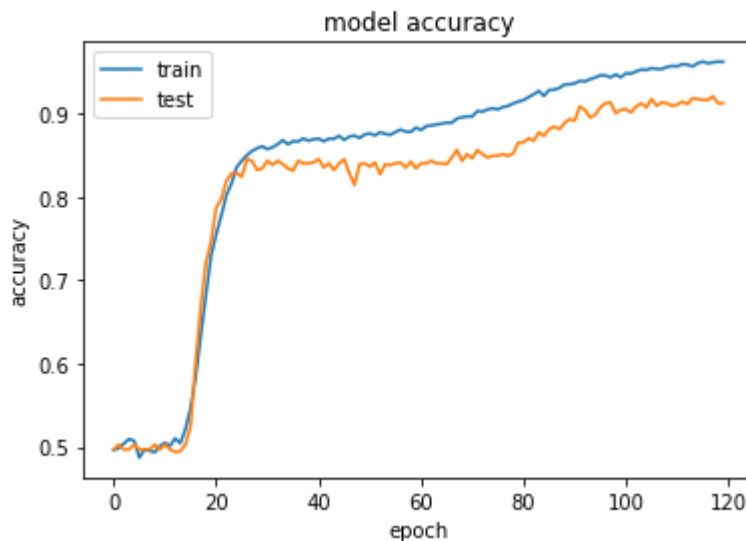
```

Evaluating the Model

Plotting Out Training History

```
In [57]: 1 import matplotlib.pyplot as plt
2 %matplotlib inline
3 print(history.history.keys())
4 # summarize history for accuracy
5 plt.plot(history.history['acc'])
6 plt.plot(history.history['val_acc'])
7 plt.title('model accuracy')
8 plt.ylabel('accuracy')
9 plt.xlabel('epoch')
10 plt.legend(['train', 'test'], loc='upper left')
11 plt.show()
```

```
dict_keys(['val_loss', 'val_acc', 'loss', 'acc'])
```



Evaluating on Given Test Set

```
In [73]: 1 model.load_weights(filename)
2 pred_results = model.predict([inputs_test, queries_test])
```

```
In [74]: 1 test_data[0][0]
```

```
Out[74]: ['Mary',
'got',
'the',
'milk',
'there',
'.',
'John',
'moved',
'to',
'the',
'bedroom',
'.']
```

```
In [75]: 1 story = ' '.join(word for word in test_data[0][0])
          2 print(story)
```

Mary got the milk there . John moved to the bedroom .

```
In [76]: 1 query = ' '.join(word for word in test_data[0][1])
          2 print(query)
```

Is John in the kitchen ?

```
In [77]: 1 print("True Test Answer from Data is:", test_data[0][2])
```

True Test Answer from Data is: no

```
In [78]: 1 #Generate prediction from model
          2 val_max = np.argmax(pred_results[0])
          3
          4 for key, val in tokenizer.word_index.items():
          5     if val == val_max:
          6         k = key
          7
          8 print("Predicted answer is: ", k)
          9 print("Probability of certainty was: ", pred_results[0][val_max])
```

Predicted answer is: no

Probability of certainty was: 0.9999999

Writing Your Own Stories and Questions

Remember you can only use words from the existing vocab

In [79]:

```
1 vocab
```

Out[79]:

```
{'.',  
'?',  
'Daniel',  
'Is',  
'John',  
'Mary',  
'Sandra',  
'apple',  
'back',  
'bathroom',  
'bedroom',  
'discarded',  
'down',  
'dropped',  
'football',  
'garden',  
'got',  
'grabbed',  
'hallway',  
'in',  
'journeyed',  
'kitchen',  
'left',  
'milk',  
'moved',  
'no',  
'office',  
'picked',  
'put',  
'the',  
'there',  
'to',  
'took',  
'travelled',  
'up',  
'went',  
'yes'}
```

In [80]:

```
1 # Note the whitespace of the periods  
2 my_story = "John left the kitchen . Sandra dropped the football in  
3 my_story.split()
```

Out[80]:

```
['John',  
'left',  
'the',  
'kitchen',  
'.',  
'Sandra',  
'dropped',  
'the',  
'football',  
'in',  
'the',  
'garden',  
'.']
```

```
In [81]: 1 my_question = "Is the football in the garden ?"
```

```
In [82]: 1 my_question.split()
```

```
Out[82]: ['Is', 'the', 'football', 'in', 'the', 'garden', '?']
```

```
In [83]: 1 mydata = [(my_story.split(),my_question.split(),'yes')]
```

```
In [84]: 1 my_story,my_ques,my_ans = vectorize_stories(mydata)
```

```
In [85]: 1 pred_results = model.predict([ my_story, my_ques])
```

```
In [86]: 1 #Generate prediction from model
2 val_max = np.argmax(pred_results[0])
3
4 for key, val in tokenizer.word_index.items():
5     if val == val_max:
6         k = key
7
8 print("Predicted answer is: ", k)
9 print("Probability of certainty was: ", pred_results[0][val_max])
```

Predicted answer is: yes

Probability of certainty was: 0.97079676

Great Job!