

```
In [1]: 1 import tensorflow as tf
        2 import numpy as np
        3 import os, re
        4 from tensorflow.python.layers.core import Dense
```

`load_sentences` reads the text data into a list, where each list element is a sentence, cropped at `MAX_CHAR_PER_LINE` characters

```
In [2]: 1 MAX_CHAR_PER_LINE = 20
        2
        3 def load_sentences(path):
        4     with open(path, 'r', encoding="ISO-8859-1") as f:
        5         data_raw = f.read().encode('ascii', 'ignore').decode('UTF-8').lower()
        6         data_alpha = re.sub('[^a-z\n]+', ' ', data_raw)
        7         data = []
        8         for line in data_alpha.split('\n'):
        9             data.append(line[:MAX_CHAR_PER_LINE])
        10    return data
```

`extract_character_vocab` converts a list of sentences to vocabularies which are python dictionaries of `int_to_symbol` {index:character} and `symbol_to_int` {character:index}.

The list comprehension `[character for line in data for character in line]` works like this (probably): `[subListItem for ListItem in List for subListItem in ListItem]`

The creation of the vocabularies work using a dictionary comprehension: `dict = {key:item for key, item in enumerate(list)}`

In [3]:

```
1 def extract_character_vocab(data):
2     special_symbols = ['<PAD>', '<UNK>', '<GO>', '<EOS>']
3     # extract unique characters from all sentences in data
4     set_symbols = set([character for line in data for character in line])
5     # add special symbols
6     all_symbols = special_symbols + list(set_symbols)
7     # create vocabularies that match symbol to index, and index to symbol
8     int_to_symbol = {word_i: word for word_i, word in enumerate(all_symbols)}
9     symbol_to_int = {word: word_i for word_i, word in int_to_symbol.items()}
10    return int_to_symbol, symbol_to_int
11
12 input_sentences = load_sentences('data/words_input.txt')
13 output_sentences = load_sentences('data/words_output.txt')
14
15 input_int_to_symbol, input_symbol_to_int = extract_character_vocab(input_sentences)
16 output_int_to_symbol, output_symbol_to_int = extract_character_vocab(output_sentences)
```

► In [4]: 1 input_int_to_symbol

```
Out[4]: {0: '<PAD>',
1: '<UNK>',
2: '<GO>',
3: '<EOS>',
4: 'a',
5: 'm',
6: 'q',
7: 'r',
8: 'u',
9: 'l',
10: 'y',
11: 'e',
12: 't',
13: 'v',
14: 'z',
15: 'd',
16: 'h',
17: 'f',
18: 'k',
19: 's',
20: 'o',
21: 'n',
22: 'p',
23: 'b',
24: 'i',
25: 'g',
26: 'j',
27: 'w',
28: 'x',
29: 'c',
30: ' '}
```

```
In [*]: 1 output_int_to_symbol
```

```
Out[5]: {0: '<PAD>',  
1: '<UNK>',  
2: '<GO>',  
3: '<EOS>',  
4: 'a',  
5: 'm',  
6: 'r',  
7: 'q',  
8: 'u',  
9: 'l',  
10: 'y',  
11: 'e',  
12: 't',  
13: 'v',  
14: 'z',  
15: 'd',  
16: 'h',  
17: 'f',  
18: 'k',  
19: 's',  
20: 'o',  
21: 'n',  
22: 'p',  
23: 'b',  
24: 'g',  
25: 'i',  
26: 'j',  
27: 'w',  
28: 'x',  
29: 'c',  
30: ' '}
```

In [*]:

```
1 NUM_EPOCHS = 300
2 RNN_STATE_DIM = 512
3 RNN_NUM_LAYERS = 2
4 ENCODER_EMBEDDING_DIM = DECODER_EMBEDDING_DIM = 64
5
6 BATCH_SIZE = int(32)
7 LEARNING_RATE = 0.0003
8
9 INPUT_NUM_VOCAB = len(input_symbol_to_int)
10 OUTPUT_NUM_VOCAB = len(output_symbol_to_int)
```

```
In [*]: 1 # Encoder placeholders
2 encoder_input_seq = tf.placeholder(
3     tf.int32,
4     [None, None],
5     name='encoder_input_seq'
6 )
7
8 encoder_seq_len = tf.placeholder(
9     tf.int32,
10    (None,),
11    name='encoder_seq_len'
12 )
13
14 # Decoder placeholders
15 decoder_output_seq = tf.placeholder(
16     tf.int32,
17     [None, None],
18     name='decoder_output_seq'
19 )
20
21 decoder_seq_len = tf.placeholder(
22     tf.int32,
23     (None,),
24     name='decoder_seq_len'
25 )
26
27 max_decoder_seq_len = tf.reduce_max(
28     decoder_seq_len,
29     name='max_decoder_seq_len'
30 )
```

```
In [*]: 1 def make_cell(state_dim):
2     lstm_initializer = tf.random_uniform_initializer(-0.1, 0.1)
3     return tf.contrib.rnn.LSTMCell(state_dim, initializer=lstm_initializer)
4
5 def make_multi_cell(state_dim, num_layers):
6     cells = [make_cell(state_dim) for _ in range(num_layers)]
7     return tf.contrib.rnn.MultiRNNCell(cells)
```

```

In [*]: 1 # Encoder embedding
2
3 encoder_input_embedded = tf.contrib.layers.embed_sequence(
4     encoder_input_seq,
5     INPUT_NUM_VOCAB,
6     ENCODER_EMBEDDING_DIM
7 )
8
9
10 # Encoder output
11
12 encoder_multi_cell = make_multi_cell(RNN_STATE_DIM, RNN_NUM_LAYERS)
13
14 encoder_output, encoder_state = tf.nn.dynamic_rnn(
15     encoder_multi_cell,
16     encoder_input_embedded,
17     sequence_length=encoder_seq_len,
18     dtype=tf.float32
19 )
20
21 del(encoder_output)

```

```

In [*]: 1 decoder_raw_seq = decoder_output_seq[:, :-1]
2 go_prefixes = tf.fill([BATCH_SIZE, 1], output_symbol_to_int['<GO>'])
3 decoder_input_seq = tf.concat([go_prefixes, decoder_raw_seq], 1)

```

```

In [*]: 1 decoder_embedding = tf.Variable(tf.random_uniform([OUTPUT_NUM_VOCAB,
2                                                     DECODER_EMBEDDING_DIM]))
3 decoder_input_embedded = tf.nn.embedding_lookup(decoder_embedding,
4                                                  decoder_input_seq)
5
6 decoder_multi_cell = make_multi_cell(RNN_STATE_DIM, RNN_NUM_LAYERS)
7
8 output_layer_kernel_initializer = tf.truncated_normal_initializer(mean=0.0, stddev=0.1)
9 output_layer = Dense(
10     OUTPUT_NUM_VOCAB,
11     kernel_initializer = output_layer_kernel_initializer
12 )

```

```
In [*]: 1 with tf.variable_scope("decode"):
2
3         training_helper = tf.contrib.seq2seq.TrainingHelper(
4             inputs=decoder_input_embedded,
5             sequence_length=decoder_seq_len,
6             time_major=False
7         )
8
9         training_decoder = tf.contrib.seq2seq.BasicDecoder(
10            decoder_multi_cell,
11            training_helper,
12            encoder_state,
13            output_layer
14        )
15
16        training_decoder_output_seq, _, _ = tf.contrib.seq2seq.dynamic_decode(
17            training_decoder,
18            impute_finished=True,
19            maximum_iterations=max_decoder_seq_len
20        )
```



```
In [*]: 1 with tf.variable_scope("decode", reuse=True):
2         start_tokens = tf.tile(
3             tf.constant([output_symbol_to_int['<GO>']],
4                         dtype=tf.int32),
5             [BATCH_SIZE],
6             name='start_tokens')
7
8         # Helper for the inference process.
9         inference_helper = tf.contrib.seq2seq.GreedyEmbeddingHelper(
10             embedding=decoder_embedding,
11             start_tokens=start_tokens,
12             end_token=output_symbol_to_int['<EOS>']
13         )
14
15         # Basic decoder
16         inference_decoder = tf.contrib.seq2seq.BasicDecoder(
17             decoder_multi_cell,
18             inference_helper,
19             encoder_state,
20             output_layer
21         )
22
23         # Perform dynamic decoding using the decoder
24         inference_decoder_output_seq, _, _ = tf.contrib.seq2seq.dynamic_decode(
25             inference_decoder,
26             impute_finished=True,
27             maximum_iterations=max_decoder_seq_len
28         )
```

```
In [*]: 1 # rename the tensor for our convenience
2 training_logits = tf.identity(training_decoder_output_seq.rnn_output, name='logits')
3 inference_logits = tf.identity(inference_decoder_output_seq.sample_id, name='predictions')
4
5 # Create the weights for sequence_loss
6 masks = tf.sequence_mask(
7     decoder_seq_len,
8     max_decoder_seq_len,
9     dtype=tf.float32,
10    name='masks'
11 )
12
13 cost = tf.contrib.seq2seq.sequence_loss(
14     training_logits,
15     decoder_output_seq,
16     masks
17 )
```

```
In [*]: 1 optimizer = tf.train.AdamOptimizer(LEARNING_RATE)
2
3 gradients = optimizer.compute_gradients(cost)
4 capped_gradients = [(tf.clip_by_value(grad, -5., 5.), var)
5                     for grad, var in gradients if grad is not None]
6 train_op = optimizer.apply_gradients(capped_gradients)
```

```
In [*]: 1 def pad(xs, size, pad):
2     return xs + [pad] * (size - len(xs))
```

```

In [*]: 1 input_seq = [
2         [input_symbol_to_int.get(symbol, input_symbol_to_int['<UNK>'])
3           for symbol in line]
4         for line in input_sentences
5     ]
6
7     output_seq = [
8         [output_symbol_to_int.get(symbol, output_symbol_to_int['<UNK>'])
9           for symbol in line] + [output_symbol_to_int['<EOS>']]
10        for line in output_sentences
11    ]
12
13    sess = tf.InteractiveSession()
14    sess.run(tf.global_variables_initializer())
15    saver = tf.train.Saver()
16
17    for epoch in range(NUM_EPOCHS + 1):
18        for batch_idx in range(len(input_sentences) // BATCH_SIZE):
19
20            input_batch, input_lengths, output_batch, output_lengths = [], [], [], []
21            for sentence in input_sentences[batch_idx:batch_idx + BATCH_SIZE]:
22                symbol_sent = [input_symbol_to_int[symbol] for symbol in sentence]
23                padded_symbol_sent = pad(symbol_sent, MAX_CHAR_PER_LINE, input_symbol_to_int['<PAD>'])
24                input_batch.append(padded_symbol_sent)
25                input_lengths.append(len(sentence))
26            for sentence in output_sentences[batch_idx:batch_idx + BATCH_SIZE]:
27                symbol_sent = [output_symbol_to_int[symbol] for symbol in sentence]
28                padded_symbol_sent = pad(symbol_sent, MAX_CHAR_PER_LINE, output_symbol_to_int['<PAD>'])
29                output_batch.append(padded_symbol_sent)
30                output_lengths.append(len(sentence))
31
32            _, cost_val = sess.run(
33                [train_op, cost],
34                feed_dict={
35                    encoder_input_seq: input_batch,
36                    encoder_seq_len: input_lengths,
37                    decoder_output_seq: output_batch,
38                    decoder_seq_len: output_lengths
39                }
40            )
41

```

```

42         if batch_idx % 629 == 0:
43             print('Epoch {}. Batch {}/{}. Cost {}'.format(epoch, batch_idx, len(input_sentences) // BATCH_SIZE, cost))
44
45         saver.save(sess, 'model.ckpt')
46     sess.close()

```

```

Epoch 0. Batch 0/6919. Cost 3.582240104675293
Epoch 0. Batch 629/6919. Cost 0.8955135941505432
Epoch 0. Batch 1258/6919. Cost 0.7081497311592102
Epoch 0. Batch 1887/6919. Cost 0.7692012786865234
Epoch 0. Batch 2516/6919. Cost 0.8782562017440796
Epoch 0. Batch 3145/6919. Cost 0.699303925037384
Epoch 0. Batch 3774/6919. Cost 0.726224958896637
Epoch 0. Batch 4403/6919. Cost 0.6380016207695007
Epoch 0. Batch 5032/6919. Cost 0.6510648727416992
Epoch 0. Batch 5661/6919. Cost 0.6080302596092224
Epoch 0. Batch 6290/6919. Cost 0.5920670628547668

```

In [*]:

```

1  sess = tf.InteractiveSession()
2  saver.restore(sess, 'model.ckpt')
3
4  example_input_sent = "do you want to play games"
5  example_input_symb = [input_symbol_to_int[symbol] for symbol in example_input_sent]
6  example_input_batch = [pad(example_input_symb, MAX_CHAR_PER_LINE, input_symbol_to_int['<PAD>'])] * BATCH_SIZE
7  example_input_lengths = [len(example_input_sent)] * BATCH_SIZE
8
9  output_ints = sess.run(inference_logits, feed_dict={
10     encoder_input_seq: example_input_batch,
11     encoder_seq_len: example_input_lengths,
12     decoder_seq_len: example_input_lengths
13 })[0]
14
15 output_str = ''.join([output_int_to_symbol[i] for i in output_ints])
16 print(output_str)

```



Present



Slides



Themes



Help