

Recognize named entities on Twitter with LSTMs

In this assignment, you will use a recurrent neural network to solve Named Entity Recognition (NER) problem. NER is a common task in natural language processing systems. It serves for extraction such entities from the text as persons, organizations, locations, etc. In this task you will experiment to recognize named entities from Twitter.

For example, we want to extract persons' and organizations' names from the text. Than for the input text:

Ian Goodfellow works for Google Brain

a NER model needs to provide the following sequence of tags:

B-PER I-PER O O B-ORG I-ORG

Where *B-* and *I-* prefixes stand for the beginning and inside of the entity, while *O* stands for out of tag or no tag. Markup with the prefix scheme is called *BIO markup*. This markup is introduced for distinguishing of consequent entities with similar types.

A solution of the task will be based on neural networks, particularly, on Bi-Directional Long Short-Term Memory Networks (Bi-LSTMs).

Libraries

For this task you will need the following libraries:

- [Tensorflow \(https://www.tensorflow.org\)](https://www.tensorflow.org) — an open-source software library for Machine Intelligence.
- [Numpy \(http://www.numpy.org\)](http://www.numpy.org) — a package for scientific computing.

If you have never worked with Tensorflow, you would probably need to read some tutorials during your work on this assignment, e.g. [this one \(https://www.tensorflow.org/tutorials/recurrent\)](https://www.tensorflow.org/tutorials/recurrent) could be a good starting point.

Data

The following cell will download all data required for this assignment into the folder `week2/data`.

```
In [1]: 1 import sys
        2 sys.path.append("..")
        3 from common.download_utils import download_week2_resources
        4
        5 download_week2_resources()
```

File data\train.txt is already downloaded.
File data\validation.txt is already downloaded.
File data\test.txt is already downloaded.

Load the Twitter Named Entity Recognition corpus

We will work with a corpus, which contains tweets with NE tags. Every line of a file contains a pair of a token (word/punctuation symbol) and a tag, separated by a whitespace. Different tweets are separated by an empty line.

The function `read_data` reads a corpus from the `file_path` and returns two lists: one with tokens and one with the corresponding tags. You need to complete this function by adding a code, which will replace a user's nickname to `<USR>` token and any URL to `<URL>` token. You could think that a URL and a nickname are just strings which start with `http://` or `https://` in case of URLs and a `@` symbol for nicknames.

```
In [2]: 1 def read_data(file_path):
2         tokens = []
3         tags = []
4
5         tweet_tokens = []
6         tweet_tags = []
7         for line in open(file_path, encoding='utf-8'):
8             line = line.strip()
9             if not line:
10                 if tweet_tokens:
11                     tokens.append(tweet_tokens)
12                     tags.append(tweet_tags)
13                     tweet_tokens = []
14                     tweet_tags = []
15             else:
16                 token, tag = line.split()
17                 # Replace all urls with <URL> token
18                 # Replace all users with <USR> token
19
20                 #####
21                 ##### YOUR CODE HERE #####
22                 #####
23                 if token.lower().startswith('http://') or token.lower().startswith('https://'):
24                     token = '<URL>'
25                 if token.startswith('@'):
26                     token = '<USR>'
27
28                 tweet_tokens.append(token)
29                 tweet_tags.append(tag)
30
31         return tokens, tags
```

And now we can load three separate parts of the dataset:

- *train* data for training the model;
- *validation* data for evaluation and hyperparameters tuning;
- *test* data for final evaluation of the model.

```
In [3]: 1 train_tokens, train_tags = read_data('data/train.txt')
2         validation_tokens, validation_tags = read_data('data/validation.txt')
3         test_tokens, test_tags = read_data('data/test.txt')
```

You should always understand what kind of data you deal with. For this purpose, you can print the data running the following cell:

```
In [4]: 1 for i in range(3):
2         for token, tag in zip(train_tokens[i], train_tags[i]):
3             print('%s\t%s' % (token, tag))
4         print()
```

```
RT      0
<USR>   0
:        0
Online  0
ticket  0
sales   0
for      0
Ghostland      B-musicartist
Observatory    I-musicartist
extended      0
until    0
6        0
PM       0
EST      0
due      0
to       0
high     0
demand   0
.        0
Get      0
```

Prepare dictionaries

To train a neural network, we will use two mappings:

- {token}→{token id}: address the row in embeddings matrix for the current token;
- {tag}→{tag id}: one-hot ground truth probability distribution vectors for computing the loss at the output of the network.

Now you need to implement the function *build_dict* which will return {token or tag}→{index} and vice versa.

```
In [5]: 1 from collections import defaultdict
```

```

In [6]: 1 def build_dict(tokens_or_tags, special_tokens):
2         """
3         tokens_or_tags: a list of lists of tokens or tags
4         special_tokens: some special tokens
5         """
6         # Create a dictionary with default value 0
7         tok2idx = defaultdict(lambda: 0)
8         idx2tok = defaultdict(lambda: 0)
9
10        # Create mappings from tokens (or tags) to indices and vice versa.
11        # At first, add special tokens (or tags) to the dictionaries.
12        # The first special token must have index 0.
13
14        # Mapping tok2idx should contain each token or tag only once.
15        # To do so, you should:
16        # 1. extract unique tokens/tags from the tokens_or_tags variable, which is not
17        #    occur in special_tokens (because they could have non-empty intersection)
18        # 2. index them (for example, you can add them into the list idx2tok
19        # 3. for each token/tag save the index into tok2idx).
20
21        #####
22        ##### YOUR CODE HERE #####
23        #####
24        tokens = [item for sublist in tokens_or_tags for item in sublist]
25        tokVocab = set(tokens) - set(special_tokens)
26        sortedTokVocab = sorted(tokVocab)
27
28        for i,u in enumerate(special_tokens + sortedTokVocab):
29            tok2idx[u] = i
30            idx2tok[i] = u
31
32        return tok2idx, idx2tok

```

After implementing the function *build_dict* you can make dictionaries for tokens and tags. Special tokens in our case will be:

- <UNK> token for out of vocabulary tokens;
- <PAD> token for padding sentence to the same length when we create batches of sentences.

```

In [7]: 1 special_tokens = ['<UNK>', '<PAD>']
2         special_tags = ['0']
3
4         # Create dictionaries
5         token2idx, idx2token = build_dict(train_tokens + validation_tokens, special_tokens)
6         tag2idx, idx2tag = build_dict(train_tags, special_tags)

```

```

In [8]: 1 tag2idx

```

```

Out[8]: defaultdict(<function __main__.build_dict.<locals>.<lambda>()>,
                    {'0': 0,
                     'B-company': 1,
                     'B-facility': 2,
                     'B-geo-loc': 3,
                     'B-movie': 4,
                     'B-musicartist': 5,
                     'B-other': 6,
                     'B-person': 7,
                     'B-product': 8,
                     'B-sportsteam': 9,
                     'B-tvshow': 10,
                     'I-company': 11,
                     'I-facility': 12,
                     'I-geo-loc': 13,
                     'I-movie': 14,
                     'I-musicartist': 15,
                     'I-other': 16,
                     'I-person': 17,
                     'I-product': 18,
                     'I-sportsteam': 19,
                     'I-tvshow': 20})

```

```
In [9]: 1 idx2tag
```

```
Out[9]: defaultdict(<function __main__.build_dict.<locals>.<lambda>()>,
                    {0: 'O',
                     1: 'B-company',
                     2: 'B-facility',
                     3: 'B-geo-loc',
                     4: 'B-movie',
                     5: 'B-musicartist',
                     6: 'B-other',
                     7: 'B-person',
                     8: 'B-product',
                     9: 'B-sportsteam',
                     10: 'B-tvshow',
                     11: 'I-company',
                     12: 'I-facility',
                     13: 'I-geo-loc',
                     14: 'I-movie',
                     15: 'I-musicartist',
                     16: 'I-other',
                     17: 'I-person',
                     18: 'I-product',
                     19: 'I-sportsteam',
                     20: 'I-tvshow'})})
```

The next additional functions will help you to create the mapping between tokens and ids for a sentence.

```
In [10]: 1 def words2idxs(tokens_list):
2         return [token2idx[word] for word in tokens_list]
3
4 def tags2idxs(tags_list):
5         return [tag2idx[tag] for tag in tags_list]
6
7 def idxs2words(idxs):
8         return [idx2token[idx] for idx in idxs]
9
10 def idxs2tags(idxs):
11         return [idx2tag[idx] for idx in idxs]
```

Generate batches

Neural Networks are usually trained with batches. It means that weight updates of the network are based on several sequences at every single time. The tricky part is that all sequences within a batch need to have the same length. So we will pad them with a special <PAD> token. It is also a good practice to provide RNN with sequence lengths, so it can skip computations for padding parts. We provide the batching function *batches_generator* readily available for you to save time.

```
In [11]: 1 def batches_generator(batch_size, tokens, tags,
2                           shuffle=True, allow_smaller_last_batch=True):
3         """Generates padded batches of tokens and tags."""
4
5         n_samples = len(tokens)
6         if shuffle:
7             order = np.random.permutation(n_samples)
8         else:
9             order = np.arange(n_samples)
10
11         n_batches = n_samples // batch_size
12         if allow_smaller_last_batch and n_samples % batch_size:
13             n_batches += 1
14
15         for k in range(n_batches):
16             batch_start = k * batch_size
17             batch_end = min((k + 1) * batch_size, n_samples)
18             current_batch_size = batch_end - batch_start
19             x_list = []
20             y_list = []
21             max_len_token = 0
22             for idx in order[batch_start: batch_end]:
23                 x_list.append(words2idxs(tokens[idx]))
24                 y_list.append(tags2idxs(tags[idx]))
25                 max_len_token = max(max_len_token, len(tags[idx]))
26
27             # Fill in the data into numpy nd-arrays filled with padding indices.
28             x = np.ones([current_batch_size, max_len_token], dtype=np.int32) * token2idx['<PAD>']
29             y = np.ones([current_batch_size, max_len_token], dtype=np.int32) * tag2idx['O']
30             lengths = np.zeros(current_batch_size, dtype=np.int32)
31             for n in range(current_batch_size):
32                 utt_len = len(x_list[n])
33                 x[n, :utt_len] = x_list[n]
34                 lengths[n] = utt_len
35                 y[n, :utt_len] = y_list[n]
36             yield x, y, lengths
```

Build a recurrent neural network

This is the most important part of the assignment. Here we will specify the network architecture based on TensorFlow building blocks. It's fun and easy as a lego constructor! We will create an LSTM network which will produce probability distribution over tags for each token in a sentence. To take into account both right and left contexts of the token, we will use Bi-Directional LSTM (Bi-LSTM). Dense layer will be used on top to perform tag classification.

```
In [12]: 1 import tensorflow as tf
          2 import numpy as np
```

```
C:\Users\XiaoWei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\framework\dtypes.py:516: FutureWarning:
Passing (type, 1) or '1type' as a synonym of type is deprecated; in a future version of numpy, it will be understood
as (type, (1,)) / '(1,)type'.
    _np_qint8 = np.dtype [("qint8", np.int8, 1)]
C:\Users\XiaoWei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\framework\dtypes.py:517: FutureWarning:
Passing (type, 1) or '1type' as a synonym of type is deprecated; in a future version of numpy, it will be understood
as (type, (1,)) / '(1,)type'.
    _np_quint8 = np.dtype [("quint8", np.uint8, 1)]
C:\Users\XiaoWei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\framework\dtypes.py:518: FutureWarning:
Passing (type, 1) or '1type' as a synonym of type is deprecated; in a future version of numpy, it will be understood
as (type, (1,)) / '(1,)type'.
    _np_qint16 = np.dtype [("qint16", np.int16, 1)]
C:\Users\XiaoWei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\framework\dtypes.py:519: FutureWarning:
Passing (type, 1) or '1type' as a synonym of type is deprecated; in a future version of numpy, it will be understood
as (type, (1,)) / '(1,)type'.
    _np_quint16 = np.dtype [("quint16", np.uint16, 1)]
C:\Users\XiaoWei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\framework\dtypes.py:520: FutureWarning:
Passing (type, 1) or '1type' as a synonym of type is deprecated; in a future version of numpy, it will be understood
as (type, (1,)) / '(1,)type'.
```

```
In [13]: 1 class BiLSTMModel():
          2     pass
```

First, we need to create [placeholders](https://www.tensorflow.org/api_docs/python/tf/compat/v1/placeholder) (https://www.tensorflow.org/api_docs/python/tf/compat/v1/placeholder) to specify what data we are going to feed into the network during the execution time. For this task we will need the following placeholders:

- *input_batch* — sequences of words (the shape equals to [batch_size, sequence_len]);
- *ground_truth_tags* — sequences of tags (the shape equals to [batch_size, sequence_len]);
- *lengths* — lengths of not padded sequences (the shape equals to [batch_size]);
- *dropout_ph* — dropout keep probability; this placeholder has a predefined value 1;
- *learning_rate_ph* — learning rate; we need this placeholder because we want to change the value during training.

It could be noticed that we use *None* in the shapes in the declaration, which means that data of any size can be feeded.

You need to complete the function *declare_placeholders*.

```
In [14]: 1 def declare_placeholders(self):
          2     """Specifies placeholders for the model."""
          3
          4     # Placeholders for input and ground truth output.
          5     self.input_batch = tf.placeholder(dtype=tf.int32, shape=[None, None], name='input_batch')
          6     ##### YOUR CODE HERE #####
          7     self.ground_truth_tags = tf.placeholder(dtype=tf.int32, shape=[None, None], name='ground_truth_tags')
          8
          9     # Placeholder for lengths of the sequences.
         10     self.lengths = tf.placeholder(dtype=tf.int32, shape=[None], name='lengths')
         11
         12     # Placeholder for a dropout keep probability. If we don't feed
         13     # a value for this placeholder, it will be equal to 1.0.
         14     self.dropout_ph = tf.placeholder_with_default(tf.cast(1.0, tf.float32), shape=[])
         15
         16     # Placeholder for a learning rate (tf.float32).
         17     ##### YOUR CODE HERE #####
         18     self.learning_rate_ph = tf.placeholder(dtype=tf.float32, shape=[], name='learning_rate_ph')
```

```
In [15]: 1 BiLSTMModel.__declare_placeholders = classmethod(declare_placeholders)
```

Now, let us specify the layers of the neural network. First, we need to perform some preparatory steps:

- Create embeddings matrix with [tf.Variable](https://www.tensorflow.org/api_docs/python/tf/Variable) (https://www.tensorflow.org/api_docs/python/tf/Variable). Specify its name (*embeddings_matrix*), type (*tf.float32*), and initialize with random values.
- Create forward and backward LSTM cells. TensorFlow provides a number of RNN cells ready for you. We suggest that you use *LSTMCell*, but you can also experiment with other types, e.g. GRU cells. [This](http://colah.github.io/posts/2015-08-Understanding-LSTMs/) (<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>) blogpost could be interesting if you want to learn more about the differences.
- Wrap your cells with [DropoutWrapper](https://www.tensorflow.org/api_docs/python/tf/contrib/rnn/DropoutWrapper) (https://www.tensorflow.org/api_docs/python/tf/contrib/rnn/DropoutWrapper). Dropout is an important regularization technique for neural networks. Specify all keep probabilities using the dropout placeholder that we created before.

After that, you can build the computation graph that transforms an input_batch:

- [Look up \(https://www.tensorflow.org/api_docs/python/tf/nn/embedding_lookup\)](https://www.tensorflow.org/api_docs/python/tf/nn/embedding_lookup) embeddings for an *input_batch* in the prepared *embedding_matrix*.
- Pass the embeddings through [Bidirectional Dynamic RNN \(https://www.tensorflow.org/api_docs/python/tf/nn/bidirectional_dynamic_rnn\)](https://www.tensorflow.org/api_docs/python/tf/nn/bidirectional_dynamic_rnn) with the specified forward and backward cells. Use the lengths placeholder here to avoid computations for padding tokens inside the RNN.
- Create a dense layer on top. Its output will be used directly in loss function.

Fill in the code below. In case you need to debug something, the easiest way is to check that tensor shapes of each step match the expected ones.

```
In [16]: 1 def build_layers(self, vocabulary_size, embedding_dim, n_hidden_rnn, n_tags):
2         """Specifies bi-LSTM architecture and computes logits for inputs."""
3
4         # Create embedding variable (tf.Variable) with dtype tf.float32
5         initial_embedding_matrix = np.random.randn(vocabulary_size, embedding_dim) / np.sqrt(embedding_dim)
6         ##### YOUR CODE HERE #####
7         embeddings_matrix_variable = tf.Variable(initial_value=initial_embedding_matrix,
8                                                  dtype=tf.float32, name='embeddings_matrix')
9
10        # Create RNN cells (for example, tf.nn.rnn_cell.BasicLSTMCell) with n_hidden_rnn number of units
11        # and dropout (tf.nn.rnn_cell.DropoutWrapper), initializing all *_keep_prob with dropout placeholder.
12        ##### YOUR CODE HERE #####
13        forward_cell = tf.nn.rnn_cell.BasicLSTMCell(num_units=n_hidden_rnn)
14        forward_cell = tf.nn.rnn_cell.DropoutWrapper(forward_cell, input_keep_prob=self.dropout_ph,
15                                                    output_keep_prob=self.dropout_ph, state_keep_prob=self.dropout_ph)
16
17        backward_cell = tf.nn.rnn_cell.BasicLSTMCell(num_units=n_hidden_rnn)
18        backward_cell = tf.nn.rnn_cell.DropoutWrapper(backward_cell, input_keep_prob=self.dropout_ph,
19                                                    output_keep_prob=self.dropout_ph, state_keep_prob=self.dropout_ph)
20
21        # Look up embeddings for self.input_batch (tf.nn.embedding_lookup).
22        # Shape: [batch_size, sequence_len, embedding_dim].
23        ##### YOUR CODE HERE #####
24        embeddings = tf.nn.embedding_lookup(embeddings_matrix_variable, self.input_batch)
25
26        # Pass them through Bidirectional Dynamic RNN (tf.nn.bidirectional_dynamic_rnn).
27        # Shape: [batch_size, sequence_len, 2 * n_hidden_rnn].
28        # Also don't forget to initialize sequence_length as self.lengths and dtype as tf.float32.
29        ##### YOUR CODE HERE #####
30        (rnn_output_fw, rnn_output_bw), _ = tf.nn.bidirectional_dynamic_rnn(forward_cell, backward_cell,
31                                                                           inputs = embeddings,
32                                                                           sequence_length = self.lengths,
33                                                                           dtype=tf.float32
34                                                                           )
35        rnn_output = tf.concat([rnn_output_fw, rnn_output_bw], axis=2)
36
37        # Dense layer on top.
38        # Shape: [batch_size, sequence_len, n_tags].
39        self.logits = tf.layers.dense(rnn_output, n_tags, activation=None)
```

```
In [17]: 1 BiLSTMModel.__build_layers = classmethod(build_layers)
```

To compute the actual predictions of the neural network, you need to apply [softmax \(https://www.tensorflow.org/api_docs/python/tf/nn/softmax\)](https://www.tensorflow.org/api_docs/python/tf/nn/softmax) to the last layer and find the most probable tags with [argmax \(https://www.tensorflow.org/api_docs/python/tf/argmax\)](https://www.tensorflow.org/api_docs/python/tf/argmax).

```
In [18]: 1 def compute_predictions(self):
2         """Transforms logits to probabilities and finds the most probable tags."""
3
4         # Create softmax (tf.nn.softmax) function
5         ##### YOUR CODE HERE #####
6         softmax_output = tf.nn.softmax(self.logits)
7
8         # Use argmax (tf.argmax) to get the most probable tags
9         # Don't forget to set axis=-1
10        # otherwise argmax will be calculated in a wrong way
11        ##### YOUR CODE HERE #####
12        self.predictions = tf.argmax(softmax_output, axis=-1)
```

```
In [19]: 1 BiLSTMModel.__compute_predictions = classmethod(compute_predictions)
```

During training we do not need predictions of the network, but we need a loss function. We will use [cross-entropy loss \(http://ml-cheatsheet.readthedocs.io/en/latest/loss_functions.html#cross-entropy\)](http://ml-cheatsheet.readthedocs.io/en/latest/loss_functions.html#cross-entropy), efficiently implemented in TF as [cross entropy with logits \(https://www.tensorflow.org/api_docs/python/tf/nn/softmax_cross_entropy_with_logits_v2\)](https://www.tensorflow.org/api_docs/python/tf/nn/softmax_cross_entropy_with_logits_v2). Note that it should be applied to logits of the model (not to softmax probabilities!). Also note, that we do not want to take into account loss terms coming from <PAD> tokens. So we need to mask them out, before computing [mean \(https://www.tensorflow.org/api_docs/python/tf/reduce_mean\)](https://www.tensorflow.org/api_docs/python/tf/reduce_mean).

```
In [20]: 1 def compute_loss(self, n_tags, PAD_index):
2         """Computes masked cross-entropy loss with logits."""
3
4         # Create cross entropy function (tf.nn.softmax_cross_entropy_with_logits_v2)
5         ground_truth_tags_one_hot = tf.one_hot(self.ground_truth_tags, n_tags)
6         ##### YOUR CODE HERE #####
7         loss_tensor = tf.nn.softmax_cross_entropy_with_logits_v2(labels=ground_truth_tags_one_hot, logits=self.logits)
8
9         mask = tf.cast(tf.not_equal(self.input_batch, PAD_index), tf.float32)
10        # Create Loss function which doesn't operate with <PAD> tokens (tf.reduce_mean)
11        # Be careful that the argument of tf.reduce_mean should be
12        # multiplication of mask and loss_tensor.
13        ##### YOUR CODE HERE #####
14        self.loss = tf.reduce_mean(tf.multiply(mask, loss_tensor))
```

```
In [21]: 1 BiLSTMModel.__compute_loss = classmethod(compute_loss)
```

The last thing to specify is how we want to optimize the loss. We suggest that you use [Adam](https://www.tensorflow.org/api_docs/python/tf/train/AdamOptimizer) (https://www.tensorflow.org/api_docs/python/tf/train/AdamOptimizer) optimizer with a learning rate from the corresponding placeholder. You will also need to apply clipping to eliminate exploding gradients. It can be easily done with [clip_by_norm](https://www.tensorflow.org/api_docs/python/tf/clip_by_norm) (https://www.tensorflow.org/api_docs/python/tf/clip_by_norm) function.

```
In [22]: 1 def perform_optimization(self):
2         """Specifies the optimizer and train_op for the model."""
3
4         # Create an optimizer (tf.train.AdamOptimizer)
5         ##### YOUR CODE HERE #####
6         self.optimizer = tf.train.AdamOptimizer(self.learning_rate_ph)
7         self.grads_and_vars = self.optimizer.compute_gradients(self.loss)
8
9         # Gradient clipping (tf.clip_by_norm) for self.grads_and_vars
10        # Pay attention that you need to apply this operation only for gradients
11        # because self.grads_and_vars also contains variables.
12        # List comprehension might be useful in this case.
13        clip_norm = tf.cast(1.0, tf.float32)
14        ##### YOUR CODE HERE #####
15        self.grads_and_vars = [(grad if grad is None else tf.clip_by_norm(grad, clip_norm), var)
16                               for grad, var in self.grads_and_vars]
17
18        self.train_op = self.optimizer.apply_gradients(self.grads_and_vars)
```

```
In [23]: 1 BiLSTMModel.__perform_optimization = classmethod(perform_optimization)
```

Congratulations! You have specified all the parts of your network. You may have noticed, that we didn't deal with any real data yet, so what you have written is just recipes on how the network should function. Now we will put them to the constructor of our Bi-LSTM class to use it in the next section.

```
In [24]: 1 def init_model(self, vocabulary_size, n_tags, embedding_dim, n_hidden_rnn, PAD_index):
2         self.__declare_placeholders()
3         self.__build_layers(vocabulary_size, embedding_dim, n_hidden_rnn, n_tags)
4         self.__compute_predictions()
5         self.__compute_loss(n_tags, PAD_index)
6         self.__perform_optimization()
```

```
In [25]: 1 BiLSTMModel.__init__ = classmethod(init_model)
```

Train the network and predict tags

[Session.run](https://www.tensorflow.org/api_docs/python/tf/Session#run) (https://www.tensorflow.org/api_docs/python/tf/Session#run) is a point which initiates computations in the graph that we have defined. To train the network, we need to compute *self.train_op*, which was declared in *perform_optimization*. To predict tags, we just need to compute *self.predictions*. Anyway, we need to feed actual data through the placeholders that we defined before.

```
In [26]: 1 def train_on_batch(self, session, x_batch, y_batch, lengths, learning_rate, dropout_keep_probability):
2         feed_dict = {self.input_batch: x_batch,
3                     self.ground_truth_tags: y_batch,
4                     self.learning_rate_ph: learning_rate,
5                     self.dropout_ph: dropout_keep_probability,
6                     self.lengths: lengths}
7
8         session.run(self.train_op, feed_dict=feed_dict)
```

```
In [27]: 1 BiLSTMModel.train_on_batch = classmethod(train_on_batch)
```

Implement the function *predict_for_batch* by initializing *feed_dict* with input *x_batch* and *lengths* and running the *session* for *self.predictions*.

```
In [28]: 1 def predict_for_batch(self, session, x_batch, lengths):
2         #####
3         ##### YOUR CODE HERE #####
4         #####
5         feed_dict = {self.input_batch:x_batch,
6                       self.lengths: lengths}
7
8         predictions = session.run(self.predictions, feed_dict=feed_dict)
9         return predictions
```

```
In [29]: 1 BiLSTMModel.predict_for_batch = classmethod(predict_for_batch)
```

We finished with necessary methods of our BiLSTMModel model and almost ready to start experimenting.

Evaluation

To simplify the evaluation process we provide two functions for you:

- *predict_tags*: uses a model to get predictions and transforms indices to tokens and tags;
- *eval_conll*: calculates precision, recall and F1 for the results.

```
In [30]: 1 from evaluation import precision_recall_f1
```

```
In [31]: 1 def predict_tags(model, session, token_idxs_batch, lengths):
2         """Performs predictions and transforms indices to tokens and tags."""
3
4         tag_idxs_batch = model.predict_for_batch(session, token_idxs_batch, lengths)
5
6         tags_batch, tokens_batch = [], []
7         for tag_idxs, token_idxs in zip(tag_idxs_batch, token_idxs_batch):
8             tags, tokens = [], []
9             for tag_idx, token_idx in zip(tag_idxs, token_idxs):
10                tags.append(idx2tag[tag_idx])
11                tokens.append(idx2token[token_idx])
12            tags_batch.append(tags)
13            tokens_batch.append(tokens)
14        return tags_batch, tokens_batch
15
16
17 def eval_conll(model, session, tokens, tags, short_report=True):
18     """Computes NER quality measures using CONLL shared task script."""
19
20     y_true, y_pred = [], []
21     for x_batch, y_batch, lengths in batches_generator(1, tokens, tags):
22         tags_batch, tokens_batch = predict_tags(model, session, x_batch, lengths)
23         if len(x_batch[0]) != len(tags_batch[0]):
24             raise Exception("Incorrect length of prediction for the input, "
25                             "expected length: %i, got: %i" % (len(x_batch[0]), len(tags_batch[0])))
26         predicted_tags = []
27         ground_truth_tags = []
28         for gt_tag_idx, pred_tag, token in zip(y_batch[0], tags_batch[0], tokens_batch[0]):
29             if token != '<PAD>':
30                 ground_truth_tags.append(idx2tag[gt_tag_idx])
31                 predicted_tags.append(pred_tag)
32
33         # We extend every prediction and ground truth sequence with 'O' tag
34         # to indicate a possible end of entity.
35         y_true.extend(ground_truth_tags + ['O'])
36         y_pred.extend(predicted_tags + ['O'])
37
38     results = precision_recall_f1(y_true, y_pred, print_results=True, short_report=short_report)
39     return results
```

Run your experiment

Create *BiLSTMModel* model with the following parameters:

- *vocabulary_size* — number of tokens;
- *n_tags* — number of tags;
- *embedding_dim* — dimension of embeddings, recommended value: 200;
- *n_hidden_rnn* — size of hidden layers for RNN, recommended value: 200;
- *PAD_index* — an index of the padding token (<PAD>).

Set hyperparameters. You might want to start with the following recommended values:

- *batch_size*: 32;
- 4 epochs;
- starting value of *learning_rate*: 0.005

- *learning_rate_decay*: a square root of 2;
- *dropout_keep_probability*: try several values: 0.1, 0.5, 0.9.

However, feel free to conduct more experiments to tune hyperparameters and earn extra points for the assignment.

```
In [32]: 1 tf.reset_default_graph()
2
3 ##### YOUR CODE HERE #####
4 model = BiLSTMModel(vocabulary_size=len(token2idx),
5                     n_tags=len(tag2idx),
6                     embedding_dim=200,
7                     n_hidden_rnn=200,
8                     PAD_index=token2idx['<PAD>']
9                     )
10
11 batch_size = 32
12 n_epochs = 6
13 learning_rate = 0.005
14 learning_rate_decay = np.sqrt(2)
15 dropout_keep_probability = 0.9
```

WARNING:tensorflow:From <ipython-input-16-120aaaafea26>:13: BasicLSTMCell.__init__ (from tensorflow.python.ops.rnn_cell_impl) is deprecated and will be removed in a future version.

Instructions for updating:

This class is equivalent as tf.keras.layers.LSTMCell, and will be replaced by that in Tensorflow 2.0.

WARNING:tensorflow:From <ipython-input-16-120aaaafea26>:33: bidirectional_dynamic_rnn (from tensorflow.python.ops.rnn) is deprecated and will be removed in a future version.

Instructions for updating:

Please use `keras.layers.Bidirectional(keras.layers.RNN(cell))`, which is equivalent to this API

WARNING:tensorflow:From C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\ops\rnn.py:464: dynamic_rnn (from tensorflow.python.ops.rnn) is deprecated and will be removed in a future version.

Instructions for updating:

Please use `keras.layers.RNN(cell)`, which is equivalent to this API

WARNING:tensorflow:From C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\ops\init_ops.py:1251: calling VarianceScaling.__init__ (from tensorflow.python.ops.init_ops) with dtype is deprecated and will be removed in a future version.

Instructions for updating:

Call initializer instance with the dtype argument instead of passing it to the constructor

WARNING:tensorflow:From C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\ops\rnn_cell_impl.py:738: calling Zeros.__init__ (from tensorflow.python.ops.init_ops) with dtype is deprecated and will be removed in a future version.

If you got an error `Tensor conversion requested dtype float64 for Tensor with dtype float32` in this point, check if there are variables without dtype initialised. Set the value of dtype equals to `tf.float32` for such variables.

Finally, we are ready to run the training!

```
In [33]: 1 %%time
2 sess = tf.Session()
3 sess.run(tf.global_variables_initializer())
4
5 print('Start training... \n')
6 for epoch in range(n_epochs):
7     # For each epoch evaluate the model on train and validation data
8     print('-' * 20 + ' Epoch {}'.format(epoch+1) + 'of {}'.format(n_epochs) + '-' * 20)
9     print('Train data evaluation:')
10    eval_conll(model, sess, train_tokens, train_tags, short_report=True)
11    print('Validation data evaluation:')
12    eval_conll(model, sess, validation_tokens, validation_tags, short_report=True)
13
14    # Train the model
15    for x_batch, y_batch, lengths in batches_generator(batch_size, train_tokens, train_tags):
16        model.train_on_batch(sess, x_batch, y_batch, lengths, learning_rate, dropout_keep_probability)
17
18    # Decaying the learning rate
19    learning_rate = learning_rate / learning_rate_decay
20
21    print('...training finished.')
```

Start training...

----- Epoch 1 of 6 -----

Train data evaluation:

processed 105778 tokens with 4489 phrases; found: 76264 phrases; correct: 149.

precision: 0.20%; recall: 3.32%; F1: 0.37

Validation data evaluation:

processed 12836 tokens with 537 phrases; found: 9296 phrases; correct: 22.

precision: 0.24%; recall: 4.10%; F1: 0.45

----- Epoch 2 of 6 -----

Train data evaluation:

processed 105778 tokens with 4489 phrases; found: 3981 phrases; correct: 1121.

precision: 28.16%; recall: 24.97%; F1: 26.47

..

Now let us see full quality reports for the final model on train, validation, and test sets. To give you a hint whether you have implemented everything correctly, you might expect F-score about 40% on the validation set.

The output of the cell below (as well as the output of all the other cells) should be present in the notebook for peer2peer review!

```
In [34]: 1 %%time
2 print('-' * 20 + ' Train set quality: ' + '-' * 20)
3 train_results = eval_conll(model, sess, train_tokens, train_tags, short_report=False)
4
5 print('-' * 20 + ' Validation set quality: ' + '-' * 20)
6 ##### YOUR CODE HERE #####
7 validation_results = eval_conll(model, sess, validation_tokens, validation_tags, short_report=False)
8
9 print('-' * 20 + ' Test set quality: ' + '-' * 20)
10 ##### YOUR CODE HERE #####
11 test_results = eval_conll(model, sess, test_tokens, test_tags, short_report=False)
```

----- Train set quality: -----
 processed 105778 tokens with 4489 phrases; found: 4507 phrases; correct: 4409.

precision: 97.83%; recall: 98.22%; F1: 98.02

company:	precision:	98.44%;	recall:	98.29%;	F1:	98.37;	predicted:	642
facility:	precision:	95.27%;	recall:	96.18%;	F1:	95.72;	predicted:	317
geo-loc:	precision:	99.30%;	recall:	99.50%;	F1:	99.40;	predicted:	998
movie:	precision:	90.41%;	recall:	97.06%;	F1:	93.62;	predicted:	73
musicartist:	precision:	96.58%;	recall:	97.41%;	F1:	97.00;	predicted:	234
other:	precision:	97.13%;	recall:	98.28%;	F1:	97.70;	predicted:	766
person:	precision:	98.98%;	recall:	98.98%;	F1:	98.98;	predicted:	886
product:	precision:	98.40%;	recall:	96.86%;	F1:	97.62;	predicted:	313
sportsteam:	precision:	98.14%;	recall:	97.24%;	F1:	97.69;	predicted:	215
tvshow:	precision:	82.54%;	recall:	89.66%;	F1:	85.95;	predicted:	63

----- Validation set quality: -----
 processed 12836 tokens with 537 phrases; found: 385 phrases; correct: 195.

precision: 50.65%; recall: 36.31%; F1: 42.30

company:	precision:	63.22%;	recall:	52.88%;	F1:	57.59;	predicted:	87
facility:	precision:	44.83%;	recall:	38.24%;	F1:	41.27;	predicted:	29
geo-loc:	precision:	70.37%;	recall:	50.44%;	F1:	58.76;	predicted:	81
movie:	precision:	0.00%;	recall:	0.00%;	F1:	0.00;	predicted:	16
musicartist:	precision:	62.50%;	recall:	17.86%;	F1:	27.78;	predicted:	8
other:	precision:	38.81%;	recall:	32.10%;	F1:	35.14;	predicted:	67
person:	precision:	55.56%;	recall:	26.79%;	F1:	36.14;	predicted:	54
product:	precision:	15.79%;	recall:	8.82%;	F1:	11.32;	predicted:	19
sportsteam:	precision:	31.58%;	recall:	30.00%;	F1:	30.77;	predicted:	19
tvshow:	precision:	0.00%;	recall:	0.00%;	F1:	0.00;	predicted:	5

----- Test set quality: -----
 processed 13258 tokens with 604 phrases; found: 453 phrases; correct: 228.

precision: 50.33%; recall: 37.75%; F1: 43.14

company:	precision:	61.40%;	recall:	41.67%;	F1:	49.65;	predicted:	57
facility:	precision:	51.28%;	recall:	42.55%;	F1:	46.51;	predicted:	39
geo-loc:	precision:	67.74%;	recall:	50.91%;	F1:	58.13;	predicted:	124
movie:	precision:	0.00%;	recall:	0.00%;	F1:	0.00;	predicted:	6
musicartist:	precision:	6.67%;	recall:	3.70%;	F1:	4.76;	predicted:	15
other:	precision:	38.04%;	recall:	33.98%;	F1:	35.90;	predicted:	92
person:	precision:	55.84%;	recall:	41.35%;	F1:	47.51;	predicted:	77
product:	precision:	13.64%;	recall:	10.71%;	F1:	12.00;	predicted:	22
sportsteam:	precision:	35.00%;	recall:	22.58%;	F1:	27.45;	predicted:	20
tvshow:	precision:	0.00%;	recall:	0.00%;	F1:	0.00;	predicted:	1

Wall time: 2min 26s

Conclusions

Could we say that our model is state of the art and the results are acceptable for the task? Definately, we can say so. Nowadays, Bi-LSTM is one of the state of the art approaches for solving NER problem and it outperforms other classical methods. Despite the fact that we used small training corpora (in comparison with usual sizes of corpora in Deep Learning), our results are quite good. In addition, in this task there are many possible named entities and for some of them we have only several dozens of trainig examples, which is definately small. However, the implemented model outperforms classical CRFs for this task. Even better results could be obtained by some combinations of several types of methods, e.g. see [this \(https://arxiv.org/abs/1603.01354\)](https://arxiv.org/abs/1603.01354) paper if you are interested.