

Named Entity Recognition

Lodi Dodevksa, Viktor Petreski

University of Ljubljana

Faculty for Computer and Information Science

Večna pot 113, SI-1000 Ljubljana

lodi.dodevska@gmail.com, vpetreski96@gmail.com

1 Introduction

After we did some research on classic and state-of-the-art Named entity recognition models, we chose to implement Conditional Random Fields (CRFs), which will represent the baseline for our work. Next, we decided to continue with a model of Bi-LSTM with a CRF layer. We chose to use rather simple embedding for the data, because first we wanted to get familiar with the Bi-LSTM models, explore the influence of the parameters on the performance and then use BERT and/or ELMO or similar embeddings for the final part of the assignment. We are using CoNLL-2003 and GMB datasets, which are popular for experimenting in NER-related tasks. We did our experiments on both of the obtained datasets at the same time.

2 Related work

Named entity recognition is a very developed branch of natural language processing. It can be split in two groups: generic (tries to localize names of people, organizations, locations etc.) and domain-specific (mostly used in the pharmaceutical, biological or medical research area, for extracting proteins, genes etc.). In this assignment we will focus on the first category. We checked a few papers in order to familiarize with the existing methodologies. The algorithms that are used can be roughly divided in four groups:

- Rule-based approach: follows a set of predefined rules to extract the named entities from a text. The disadvantage of these methods is that the rules must be created manually, based on the domain we are researching (Kim and Woodland, 2000).
- Supervised learning models: such as Hidden Markov Models, Support Vector Machines (Isozaki and Kazawa, 2002), Decision Trees or Conditional Random Fields. An in-

teresting example of CRF is (Finkel et al., 2005) which explains the baseline for Stanford Named Entity Recognizer. These models are trained on an annotated corpus and they obtain fairly good results when used in the right domain, but their drawback is that they are corpus-dependent and there may not be available corpus for each domain we would like to explore.

- Unsupervised learning models: do not depend on an annotated corpus. They aim to find hidden principles or patterns that represent the data appropriately. An example of an unsupervised model is Google's BERT, which was pretrained on a large text corpus from Wikipedia.
- Deep learning approach: most of the successful state-of-the-art NER models are focused on this area. (Yadav and Bethard, 2019) gives a nice overview of deep neural network architectures and highlights some recent improvements.

3 Dataset

For initial analysis we chose Groningen Meaning Bank (GMB) dataset, which is a freely available annotated corpus of texts mostly used for named entity recognition POS tagging tasks.

The dataset is given in csv format. We used Pandas DataFrame from python for reading and analyzing it. The dataset contains 25 columns, but not all of them are going to be useful for our work in the future. Out of 25 columns, we will probably keep only *sentence_idx*, *pos*, *word*, *lemma* and *tag* for further usage. One row in the dataset corresponds to one word from the text.

First we had to deal with NaN values, or at least we thought. Since the dataset is preprocessed to some extent, there was only one erroneous line which contained only NaN values, which was a

surprise to us. After we removed the line, there were 1050794 rows left. Then we removed the duplicate sentences. After that there are 768956 rows (sentences) left and 30172 total unique words.

Next, we explored the tag column, which contains the appropriate named entity tag for each token. The following types of entities are covered:

1. *geo* = Geographical Entity
2. *org* = Organization
3. *per* = Person
4. *gpe* = Geopolitical Entity
5. *tim* = Time indicator
6. *art* = Artifact
7. *eve* = Event
8. *nat* = Natural Phenomenon

The entities are tagged using IOB format, which means that there are prefixes I and B before a tag. B indicates that the tag is at the beginning of a chunk, I denotes that the tag is inside of a chunk. If the token does not belong to any of the aforementioned entities, its corresponding tag is O. The total values from each type of entity are shown in Table 1.

Entity	Total
geo	32969
org	27136
per	24991
tim	19517
gpe	11989
art	478
eve	433
nat	205

Table 1: Number of appearances of each entity type

We can see that the tags are not uniformly distributed. We will have to take that into account when we will analyze the performance of the NER model that we will use in the future and identify the number of false positives carefully.

The CoNLL-2003 dataset is the second dataset that we are going to use. It is already completely preprocessed, so there is not much work to do in this part. There is a German and an English version, but we are using only the English one. The

dataset contains 4 columns: word, POS tag, chunk tag and NER tag. We are not going to use chunk tag for this task. Something that is very important is that this dataset contains less tags than GMB. There are *per* = Person, *org* = Organization, *loc* = Location and *misc* = Miscellaneous, which is equivalent to the *O* tag in the GMB.

4 Data preparation

Before we started with training the models we had to prepare appropriate train and test sets. CoNLL-2003 dataset already contains three subsets: the largest set is the training data and the two smaller sets (*testA* and *testB*) are for testing. On the other hand, the GMB dataset did not come split into subsets for training and testing, so we split the dataset into training and test subset with ratio 70:30. In the Bi-LSTM models, we used 10% of each training set as a validation set.

The input sequences for CRF and Bi-LSTM-CRF are different. CRF uses the words and the corresponding POS tags for creating the input features. Bi-LSTM-CRF uses only the tokens of the sentences. Bi-LSTM-CRF expects a numerical input, so we convert each token and label (tag) of the vocabulary to a corresponding integer index. The good side of this representation is that it does not require a lot of time to train. Another thing that we must take care of is the length of the sequences. In order to hand over the data as an input to Keras model, all the sequences need to have the same length. So, we found the longest sentence and set its length to be the maximum length of the input sequence. The sentences that are shorter than this are padded to get to the required length.

5 Models

As mentioned before, we are using two different models for Named entity recognition. The first model (CRF) is supposed to represent the baseline for our work and the second model (Bi-LSTM-CRF) is the one we will continue to work on and try to improve by the end of the third assignment.

5.1 CRF

Conditional Random Field is a statistical model which is used very often in the NER field. Up to some time ago the most successful NER models were based on CRF. We provide a list of features for each token in the sentence as an input to this model. Based on these features, the model tries to

predict a label for each word in the sentence, or, better said, an optimal sequence of labels for the entire sentence. The hard part here is finding the types of features that are useful and can improve the predictions. Luckily, CRFs are widely popular and there is a lot of research done. We will use the feature dictionary suggested by *sklearn-crfsuite*. Some of the token's properties that are taken into account are whether it starts with uppercase or lowercase, what is its POS tag, is it a digit or not, etc.

5.2 Bi-LSTM with CRF layer

Bi-LSTM networks are used when we need to process sequences where the past and future features are equally important, just like in our case. The first layer in our network is the Input layer, which carries only the shape of the input data. Next comes the Embedding layer. It takes three parameters as an input: size of the vocabulary of the dataset, dimensions of the embedding (size of the vector space in which the words will be embedded) and the length of the input. Since all our sentences are padded, this length is actually the length of the longest sentence.

This layer is followed by Bi-LSTM layers. Bi-LSTM is comprised of two separate LSTM layers. One reads the sentence from the beginning to the end (left to right) and the other reads it in a reversed order. We are going to optimize this layer with respect to *units*, *dropout* and *recurrent_dropout*. The *return_sequence* parameter is set to True, which means that the layer will return the full sequence of the output, and not just a final value. That is why after the Bi-LSTM comes the TimeDistributed layer. This layer allows to apply the next layer to every element of the sequence independently. The outputs of the Bi-LSTM are actually scores of each label. These scores are handed over to the CRF layer and the label that has the highest score will be selected as a prediction for the token. We could easily create a model without the last CRF layer, but researchers have shown that using CRF can improve the performance of the network, because it allows to better capture the context of the sequence.

6 Results

In almost all of the cases we obtained about 97% - 99% accuracy. But, we could not use this measure to evaluate our models because it is not ap-

propriate for this type of predictions. We are expecting that most of the tokens are assigned the label **O** since it is the most common label and it outnumbers the other labels which are much more important. That is why we are using F1-score for evaluation, which is a better metric in the case of imbalanced labels.

For the CRF model we used *lbfgs* (*Gradient descent using the L-BFGS method*) algorithm for optimization, because it is supposed to achieve better solution than regular Gradient descent with less iterations. We optimized the regularization parameters *c1* and *c2*. We also combined the training dataset of CoNLL-2003 with one of its test sets, in order to see whether that will contribute to better performance. In Table 2 are the obtained results on the CoNLL-2003 dataset. We achieved better F1-score (89%) when we trained the CRF on the original training set and on the combined original set + test_B. Table 3 shows the best F1-scores for the GMB dataset. 84%-85% F1-score is achieved with all the given parameter combinations.

c1	c2	F1-score	Train set	Test set
0.0910	0.0029	89%	Train	A
0.1000	0.0010	81%	Train	B
0.5000	0.1000	81%	Train, A	B
0.0677	0.1949	82%	Train, A	B
0.3724	0.3135	88%	Train, B	A
0.1561	0.3552	89%	Train, B	A

Table 2: F1-scores for CRF models on CoNLL-2003 dataset

c1	c2	F1-score	Train set	Test set
0.0910	0.0029	84%	Train	A
0.1000	0.0010	84%	Train	B
0.5000	0.1000	85%	Train, A	B
0.0677	0.1949	85%	Train, A	B
0.3724	0.3135	84%	Train, B	A
0.1561	0.3552	85%	Train, B	A

Table 3: F1-scores for CRF models on GMB dataset

For the Bi-LSTM model, we used relatively small dropout and recurrent dropout rates, somewhere between 0.1 and 0.2. The performance did not improve with increasing these values. We also tried different embedding sizes, but the best F1-score was achieved for embedding size 45-50. Even though we started with fairly small number of epochs, eventually we realized that with more epochs we can obtain better results, but, as always, there is a point of diminishing return. After increasing the number of epochs over 15, the im-

provement was negligent.

The best results for CoNLL-2003 are presented in Table 4 and for Groningen Meaning Bank dataset in Table 5. We achieved F1-scores of 85% and 82% for CoNLL and GMB respectively. We expect to improve these values with using more complex embedding in the future.

Emb. size	Drop-out	Rec. dropout	Epochs	LSTM units	F1 score	Train set	Test set
40	0.1	0.1	10	50	83%	Train, B	A
30	0.1	0.1	8	50	62%	Train, B	A
40	0.1	0.1	8	50	75%	Train, A	B
45	0.15	0.1	12	80	85%	Train, B	A
50	0.2	0.15	15	90	85%	Train, B	A
40	0.1	0.1	8	50	73%	Train, A	B
45	0.1	0.1	30	85	85%	Train, B	A

Table 4: F1-scores for Bi-LSTM-CRF on CoNLL-2003 dataset

Emb. size	Drop-out	Rec. dropout	Epochs	LSTM units	F1 score	Train set	Test set
20	0.5	0.5	5	40	77%	GMB train	GMB test
30	0.3	0.3	5	40	77%	GMB train	GMB test
30	0.3	0.3	5	50	79%	GMB train	GMB test
30	0.3	0.3	10	50	81%	GMB train	GMB test
40	0.1	0.1	10	80	82%	GMB train	GMB test
45	0.15	0.1	12	80	82%	GMB train	GMB test

Table 5: F1-scores for Bi-LSTM-CRF on GMB dataset

7 Conclusion and future work

Now that we are familiar with Bi-LSTMs, we are willing to explore this field furthermore. The first step will be to use BERT and/or ELMO embeddings (Straková et al., 2019) and to optimize the parameters for the new models. Next, we could try to develop a different kind of model, i.e. Bi-LSTM-CNN (Chiu and Nichols, 2016), if we have enough time.

References

- Jason P.C. Chiu and Eric Nichols. 2016. [Named entity recognition with bidirectional LSTM-CNNs](#). *Transactions of the Association for Computational Linguistics*, 4:357–370.
- Jenny Rose Finkel, Trond Grenager, and Christopher Manning. 2005. [Incorporating non-local information into information extraction systems by gibbs sampling](#). In *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, ACL ’05, page 363–370, USA. Association for Computational Linguistics.

- Hideki Isozaki and Hideto Kazawa. 2002. [Efficient support vector classifiers for named entity recognition](#). In *Proceedings of the 19th International Conference on Computational Linguistics - Volume 1*, COLING ’02, page 1–7, USA. Association for Computational Linguistics.

- Ji-Hwan Kim and Philip C. Woodland. 2000. A rule-based named entity recognition system for speech input. In *INTERSPEECH*.

- Jana Straková, Milan Straka, and Jan Hajic. 2019. [Neural architectures for nested NER through linearization](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5326–5331. Association for Computational Linguistics.

- Vikas Yadav and Steven Bethard. 2019. A survey on recent advances in named entity recognition from deep learning models.