

Natural Language Processing and the Web



Machine Learning Project - Documentation

Benedikt Lins (1799381) und Stefan Thaut (1800351)

Department 20 - Computer Science

January 8, 2019

1 Foundations

In this Machine Learning Project we want to develop a named entity recognizer based on a machine learning approach. A named entity is a set of tokens, which form a name. Examples are "New York" or "Angela Merkel". We use a Conditional Random Field as the machine learning model.

Our training set contains a set of tokens with a chunk-annotation and the associated named entity class relating to the IOB-notation. We also have a german training set, that additionally contains a lemma for the given token.

An analysis for the english training set shows, that about 83 % of all tokens are no named entities, as we can see in figure 1

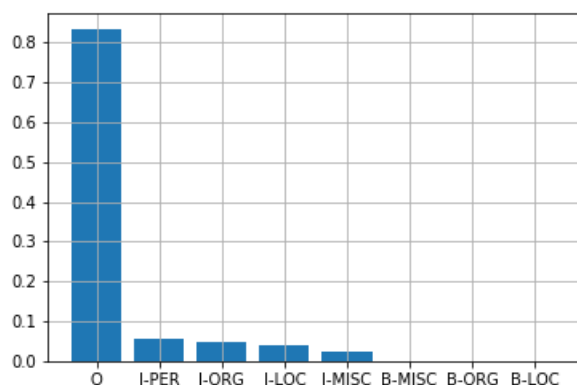


Figure 1: Class distribution of the english training set

Our CRF-model is based on the CRF-model we heard about in the fifth exercise of the lecture. So we already used the provided features:

- Is the first letter of a token capitalized? (Cap)
- Number of characters (NoOfChar)

In table 1 we can see the improvement of the measurements with this features in comparison to the baseline. The Micro F1 is increased by about just 0.035, but the Macro F1 is increased by about 0.15.

1.1 Evaluation

To have a baseline for an evaluation, with which we can compare our results, we added just a dummy feature to the CRF-model and did a testrun with the english trainingset and the english validationset as testset. The evaluationmetrics are shown in table 1.

Table 1: Evaluation results for different feature combinations

	Micro F1	Macro F1
Baseline	0.8253	0.1130
Cap, NoOfChar	0.8602	0.2693
ltcForO	0.8253	0.1130
Cap, NoOfChar, ltcForO	0.8596	0.2691

2 Additional features

2.1 Last two characters of a token

At first we have tried to use the last two characters (ltc) of a token as a feature. For the meaningfulness of this feature, we analysed the most frequent ltc for the class *O* and for the tokens, that are not in the class *O*. The ten most frequent ltc are presented in figure 2.

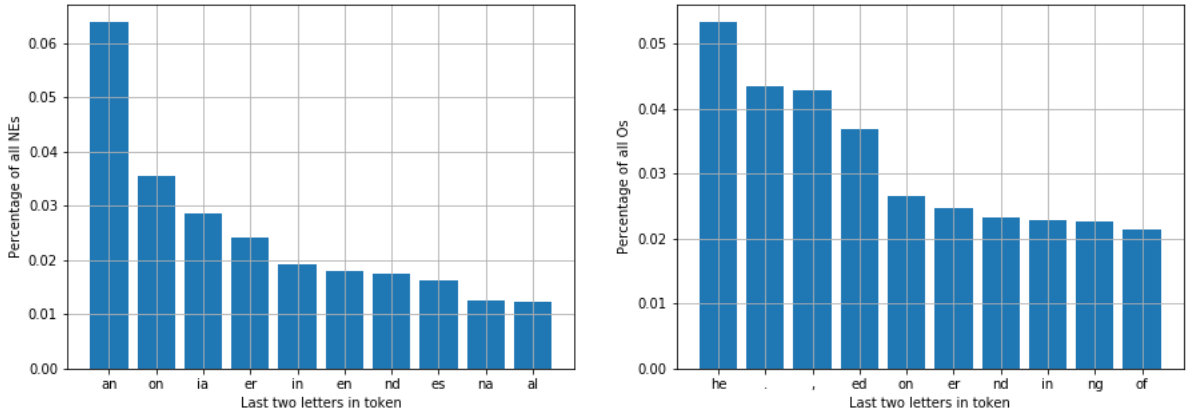


Figure 2: This figures show the ten most frequent ltc. The right side shows the ltc for the class *O* and the left side shows the ltc for all tokens, that are not in class *O*

As we can see, about six percent of all named entities end with "an" and about five percent of all non named entities end with "he". In the next step it is important to find out, if some endings are typical for named entities or for non named entities. So we examined the appearances of the most frequent ltc of named entities in tokens that are non named entities and vice versa. The results for the ten most frequent ltc are shown in figure 3. The left side shows the comparison for the most frequent ltc of named entities. The blue curve represents the occurrences of the ltc in named entities and the orange curve represents the occurrences of the ltc in non named entities. So we can imply, that the greater the distance between the both curves is, the more specific the ltc is for the class. We identified the ltc "an", "ia", and "na" as typical for named entities and analogously the ltc "he", ".", ",", "ed" and "of" as typical for non named entities.

In a first testrun with this observations, we used only the ltc, that are typical for non named entities (ltcForO). As we can see in table 1, we do not get an improvement in comparison to the baseline with this new feature. When we add this features to the capital- and the number-of-char-feature, we even decrease the evaluation in comparison to the use of both of the features.

We assumed the insufficient classification as the reason for the lower evaluation. We determined just two classes: named entities and non named entities. But there are finer classes for the named entities. So in a second step we searched for the five most frequent ltc for each class and plotted it against the frequencies of the ltc in the other classes (see figure 4). The cell for a specific ltc in a specific class is brighter, the more typical the ltc is for this class (i.e. the more token in the class end with the ltc). Thus we can use a ltc as a feature, if a cell for this ltc is nearly white or if all other cells in the column of a ltc are black (i.e. this ltc does not occur in any other class).

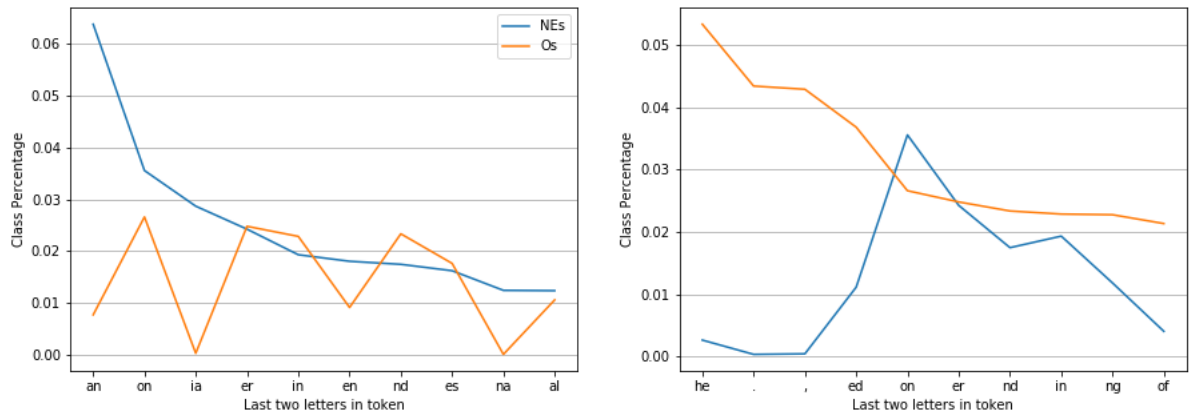


Figure 3: This figures show the percentages of the ten most frequent ltc in comparison of the classes named entity and non named entity

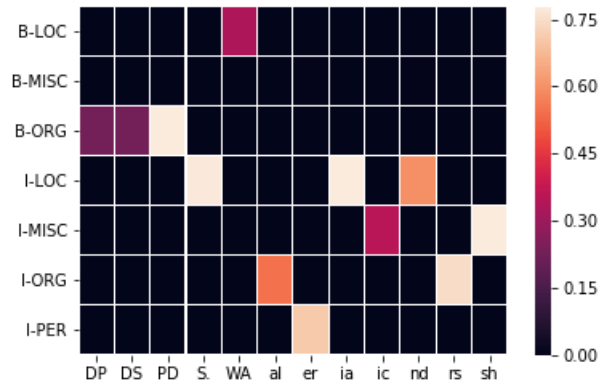


Figure 4: This figure show the percentages of the five most frequent ltc per classes

2.2 POS-Tags

With the same strategy as in figure 4 we inspected the POS tags per each class. The relevant POS tags are plotted in figure 5. All other POS tags are not on a significant level except for the class "O".

2.3 Chunks

As a last feature, we added the chunk annotation provided by the dataset. The relevant chunks are shown in figure 6. We can see, that the most relevant chunk is "I-INTJ". But we have to consider, that the relative frequency of this chunk is only by 25 % in all classes. So this chunk rather represents the class "O".

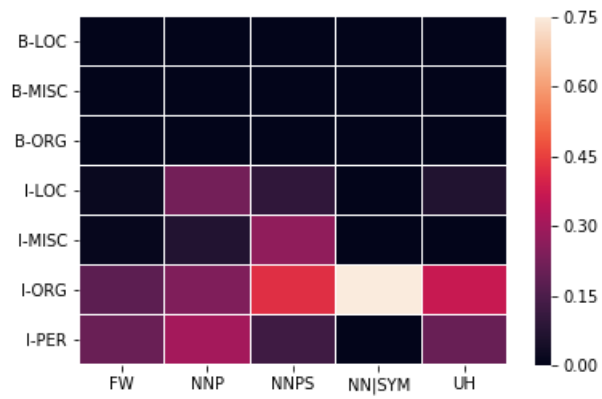


Figure 5: Relevant POS tags for each class

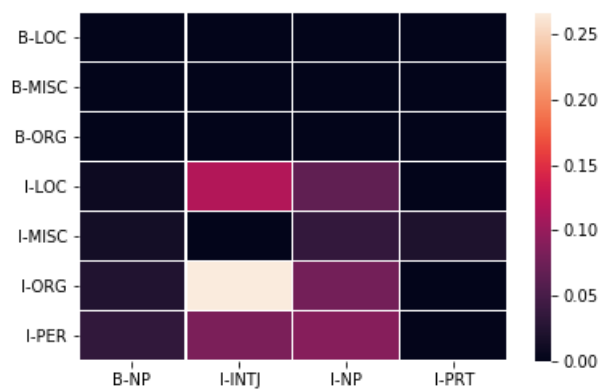


Figure 6: Relevant chunks for each class