

1 Used approach

For the purposes of training the NER system, the Learning Based Approach is utilized, using a supervised method. In particular, a Conditional Random Fields (CRF) based model has been used, which in this case is the CRFpp (CRF++) machine learning framework since for the purposes of implementing the whole system I should use distant learning / supervision methods.

2 NER system Training Analysis

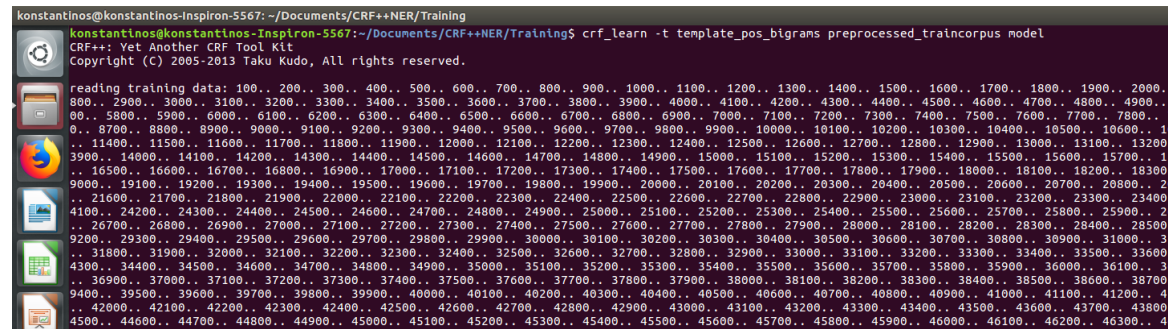
The system consists of two parts, the training of the Named Entity Recognition (NER) system and its testing.

To initiate the training, initially wp2 (aij-wikiner-en-wp2) automatically annotated corpus has been downloaded, which is also almost in IOB format ("traincorpus"). In order to change and bring the training data and its features in accordance with the CRF++ machine learning framework to the required format, it was pre-processed by using a Perl coded script (preprocessing_traincorpus.pl), using also the terminal of the operating Linux system (Ubuntu) for the purposes of its execution and its output is the preprocessed_traincorpus.txt file.

As it was mentioned above, when it comes to IOB format, every token of a sentence in a data set is being labeled with a chunk label, one of the three following exceptional chunking tags, which more specifically are I (Inside), O (Outside), or B (Beginning). In the event that a token denotes the beginning of a chunk, it is labeled as "B". The following tokens in that particular chunk are labeled as "I". Every single other token in the chunk is labeled as "O". Moreover, the labels "I" and "B" are followed by the type of the chunk, as for instance I-NP and B-NP [5].

Moreover, before the beginning of the training phase, CRF++-0.58 toolkit has been downloaded and installed by following a particular procedure. The terminal of the operating Linux system (Ubuntu) has been utilised for the purposes of the aforementioned procedure and to execute the appropriate commands, as it can be seen in the following 4 figures.

Afterwards, the preprocessed_traincorpus data set was used together with the CRF ++ template file, which is the “template_pos_bigrams” file, to initiate the training of the NER system by typing in the following command as it can be seen in the figure 1 below, using the terminal of the operating Linux system (Ubuntu) to execute the command.



```

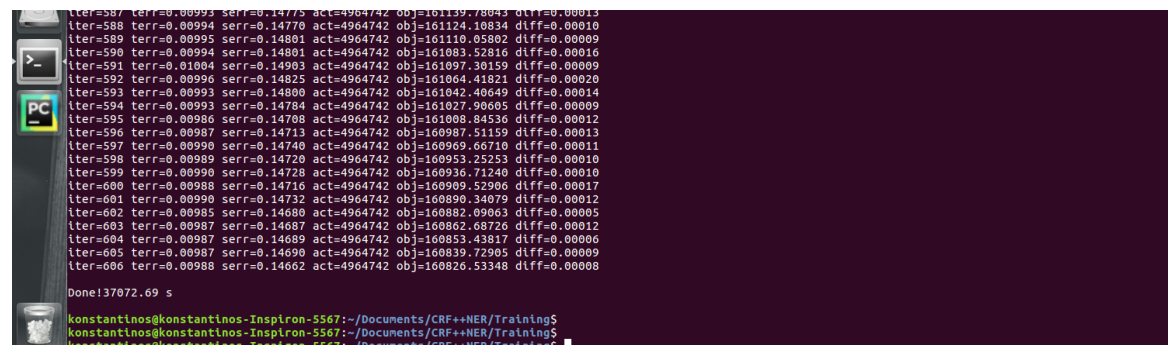
konstantinos@konstantinos-Inspiron-5567:~/Documents/CRF++NER/Training$ konstantinos@konstantinos-Inspiron-5567:~/Documents/CRF++NER/Training$ crf_learn -t template_pos_bigrams preprocessed_traincorpus model
CRF++: Yet Another CRF Tool Kit
Copyright (C) 2005-2013 Taku Kudo, All rights reserved.

reading training data: 100.. 200.. 300.. 400.. 500.. 600.. 700.. 800.. 900.. 1000.. 1100.. 1200.. 1300.. 1400.. 1500.. 1600.. 1700.. 1800.. 1900.. 2000..
800.. 2900.. 3000.. 3100.. 3200.. 3300.. 3400.. 3500.. 3600.. 3700.. 3800.. 3900.. 4000.. 4100.. 4200.. 4300.. 4400.. 4500.. 4600.. 4700.. 4800.. 4900..
600.. 5800.. 5900.. 6000.. 6100.. 6200.. 6300.. 6400.. 6500.. 6600.. 6700.. 6800.. 6900.. 7000.. 7100.. 7200.. 7300.. 7400.. 7500.. 7600.. 7700.. 7800.. 7
0.. 8700.. 8800.. 8900.. 9000.. 9100.. 9200.. 9300.. 9400.. 9500.. 9600.. 9700.. 9800.. 9900.. 10000.. 10100.. 10200.. 10300.. 10400.. 10500.. 10600.. 10
.. 11400.. 11500.. 11600.. 11700.. 11800.. 11900.. 12000.. 12100.. 12200.. 12300.. 12400.. 12500.. 12600.. 12700.. 12800.. 12900.. 13000.. 13100.. 13200..
3900.. 14000.. 14100.. 14200.. 14300.. 14400.. 14500.. 14600.. 14700.. 14800.. 14900.. 15000.. 15100.. 15200.. 15300.. 15400.. 15500.. 15600.. 15700.. 15
.. 16500.. 16600.. 16700.. 16800.. 16900.. 17000.. 17100.. 17200.. 17300.. 17400.. 17500.. 17600.. 17700.. 17800.. 17900.. 18000.. 18100.. 18200.. 18300..
9000.. 19100.. 19200.. 19300.. 19400.. 19500.. 19600.. 19700.. 19800.. 19900.. 20000.. 20100.. 20200.. 20300.. 20400.. 20500.. 20600.. 20700.. 20800.. 20
.. 21600.. 21700.. 21800.. 21900.. 22000.. 22100.. 22200.. 22300.. 22400.. 22500.. 22600.. 22700.. 22800.. 22900.. 23000.. 23100.. 23200.. 23300.. 23400..
4100.. 24200.. 24300.. 24400.. 24500.. 24600.. 24700.. 24800.. 24900.. 25000.. 25100.. 25200.. 25300.. 25400.. 25500.. 25600.. 25700.. 25800.. 25900.. 26
.. 26700.. 26800.. 26900.. 27000.. 27100.. 27200.. 27300.. 27400.. 27500.. 27600.. 27700.. 27800.. 27900.. 28000.. 28100.. 28200.. 28300.. 28400.. 28500..
9200.. 29300.. 29400.. 29500.. 29600.. 29700.. 29800.. 29900.. 30000.. 30100.. 30200.. 30300.. 30400.. 30500.. 30600.. 30700.. 30800.. 30900.. 31000.. 31
.. 31800.. 31900.. 32000.. 32100.. 32200.. 32300.. 32400.. 32500.. 32600.. 32700.. 32800.. 32900.. 33000.. 33100.. 33200.. 33300.. 33400.. 33500.. 33600..
4300.. 34400.. 34500.. 34600.. 34700.. 34800.. 34900.. 35000.. 35100.. 35200.. 35300.. 35400.. 35500.. 35600.. 35700.. 35800.. 35900.. 36000.. 36100.. 36
.. 36900.. 37000.. 37100.. 37200.. 37300.. 37400.. 37500.. 37600.. 37700.. 37800.. 37900.. 38000.. 38100.. 38200.. 38300.. 38400.. 38500.. 38600.. 38700..
9400.. 39500.. 39600.. 39700.. 39800.. 39900.. 40000.. 40100.. 40200.. 40300.. 40400.. 40500.. 40600.. 40700.. 40800.. 40900.. 41000.. 41100.. 41200.. 41
.. 42000.. 42100.. 42200.. 42300.. 42400.. 42500.. 42600.. 42700.. 42800.. 42900.. 43000.. 43100.. 43200.. 43300.. 43400.. 43500.. 43600.. 43700.. 43800..
4500.. 44600.. 44700.. 44800.. 44900.. 45000.. 45100.. 45200.. 45300.. 45400.. 45500.. 45600.. 45700.. 45800.. 45900.. 46000.. 46100.. 46200.. 46300.. 46
.. 47100.. 47200.. 47300.. 47400.. 47500.. 47600.. 47700.. 47800.. 47900.. 48000.. 48100.. 48200.. 48300.. 48400.. 48500.. 48600.. 48700.. 48800.. 48900.. 49000..

```

Figure 1: CRF++ training phase

After typing in this command, the CRF++ training model is being extracted.



```

lter=587 terr=0.00993 serr=0.14775 act=4964742 obj=161129.78043 dlf=0.00013
lter=588 terr=0.00994 serr=0.14770 act=4964742 obj=161124.10834 dlf=0.00010
lter=589 terr=0.00995 serr=0.14801 act=4964742 obj=161110.05802 dlf=0.00009
lter=590 terr=0.00994 serr=0.14801 act=4964742 obj=161083.52816 dlf=0.00016
lter=591 terr=0.01004 serr=0.14903 act=4964742 obj=161097.30159 dlf=0.00009
lter=592 terr=0.00996 serr=0.14825 act=4964742 obj=161004.41821 dlf=0.00020
lter=593 terr=0.00993 serr=0.14800 act=4964742 obj=161042.40649 dlf=0.00014
lter=594 terr=0.00993 serr=0.14784 act=4964742 obj=161027.90605 dlf=0.00009
lter=595 terr=0.00986 serr=0.14708 act=4964742 obj=161008.84536 dlf=0.00012
lter=596 terr=0.00987 serr=0.14713 act=4964742 obj=160987.51159 dlf=0.00013
lter=597 terr=0.00990 serr=0.14740 act=4964742 obj=160969.66710 dlf=0.00011
lter=598 terr=0.00989 serr=0.14720 act=4964742 obj=160953.25253 dlf=0.00010
lter=599 terr=0.00990 serr=0.14720 act=4964742 obj=160936.71240 dlf=0.00010
lter=600 terr=0.00988 serr=0.14716 act=4964742 obj=160909.52906 dlf=0.00017
lter=601 terr=0.00990 serr=0.14732 act=4964742 obj=160890.34079 dlf=0.00012
lter=602 terr=0.00985 serr=0.14680 act=4964742 obj=160882.09063 dlf=0.00005
lter=603 terr=0.00987 serr=0.14687 act=4964742 obj=160862.68726 dlf=0.00012
lter=604 terr=0.00987 serr=0.14689 act=4964742 obj=160853.43817 dlf=0.00006
lter=605 terr=0.00987 serr=0.14690 act=4964742 obj=160839.72905 dlf=0.00009
lter=606 terr=0.00988 serr=0.14662 act=4964742 obj=160826.53348 dlf=0.00008

Done!37872.69 s
konstantinos@konstantinos-Inspiron-5567:~/Documents/CRF++NER/Training$
konstantinos@konstantinos-Inspiron-5567:~/Documents/CRF++NER/Training$

```

Figure 2: CRF++ training phase finished (606 iterations and approximately 10 hours to finish)

What is more, there is an option of reading the extracted model by typing in the terminal, the command `ls model*` to display in the terminal the extracted model files and then the command `vim model.txt` in order to open and read the extracted model after the training phase of the Named Entity Recognition system. Nevertheless, it is exactly the same as locating manually and displaying the file from the path that the extracted model has been stored, in the system’s folder. Also, the extracted model’s txt file has to be stored in the testing folder of the whole system as well, for the purposes of the testing phase, which will be analysed below.

With regards to the template file, it consists of several lines where each one of them signifies one particular template. Every template of the whole template, will have a determined token from the input information (data), which is accomplished by using the special macro `%x[row,col]`, where "row" is for appointing the relative position from the present token which is being focused, as well as "col" indicates the outright column position. However, each line that is not empty, is a training data point in the training data for the CRF ++ train model. It is important to emphasize that only those features that are in the template file and has been stated exactly, will be utilized for

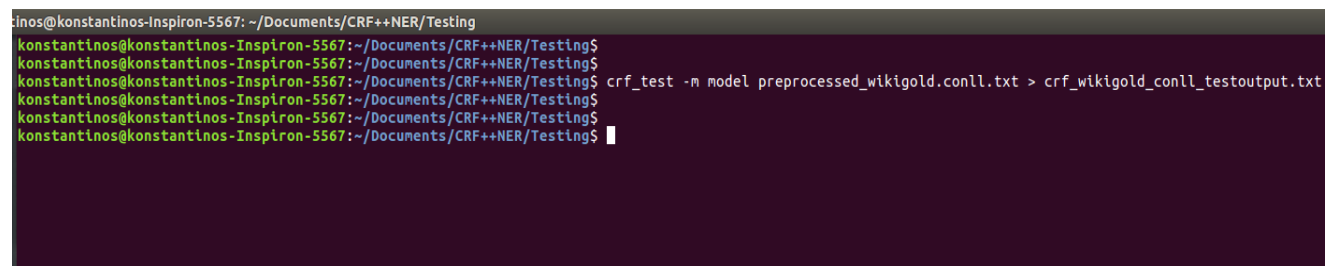
training our model. Also, bigram features are being defined by this template file. More specifically, a blend of the present yield token and past yield token (bigram) is consequently produced by using this template.

As it can be seen in the figure above, "iter" is the number of processed iterations, "terr" is the error rate concerning labels (number of error tags divided by the number of all tags), "serr" is sentence error rate (number of error sentences divided by the number of all sentences). Additionally, the "obj" stands for current object value. At the point when this esteem unites to a settled point, the iteration is being stopped by CRF++. Lastly, "diff" is the relative difference from the past object value.

3 Testing Analysis

Subsequently, in order to start the testing phase of the whole NER system, the gold standard data set wikigold.conll.txt has been collected and aggregated in order to evaluate the whole system. In order to change and bring the training data and its features in accordance with the CRF++ machine learning framework to the required format, it was pre-processed by utilizing a Python coded script (Preprocess_testing_corpus.py), using also the terminal of the operating Linux system (Ubuntu) for the purposes of its execution, as its output is the “preprocessed_wikigold.conll” txt file.

Then, the “preprocessed_wikigold.conll” data set was utilized along with an empty txt file (crf_wikigold_conll_testoutput.txt) and the extracted training model from the training phase, in order to test the whole system and create an output “crf_wikigold_conll_testoutput.txt”, which contains the predicted NER features of the “preprocessed_wikigold.conll” data set, as it can be seen in the figure 3 below. The terminal of the operating Linux system (Ubuntu) was utilized to execute the command.

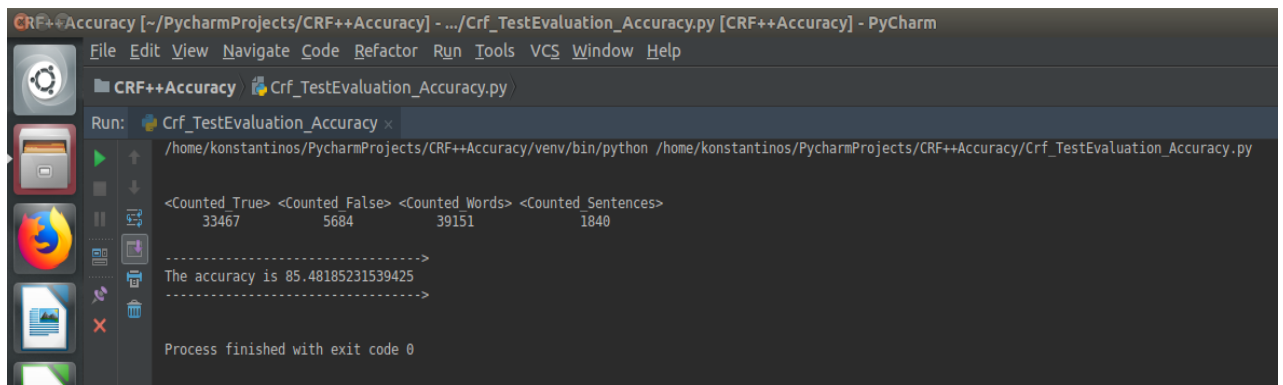
A terminal window with a dark background and light green text. The prompt is 'konstantinos@konstantinos-Inspiron-5567: ~/Documents/CRF++NER/Testing'. The command entered is 'crf_test -n model preprocessed_wikigold.conll.txt > crf_wikigold_conll_testoutput.txt'. The cursor is at the end of the command line.

```
konstantinos@konstantinos-Inspiron-5567: ~/Documents/CRF++NER/Testing$  
konstantinos@konstantinos-Inspiron-5567:~/Documents/CRF++NER/Testing$  
konstantinos@konstantinos-Inspiron-5567:~/Documents/CRF++NER/Testing$ crf_test -n model preprocessed_wikigold.conll.txt > crf_wikigold_conll_testoutput.txt  
konstantinos@konstantinos-Inspiron-5567:~/Documents/CRF++NER/Testing$  
konstantinos@konstantinos-Inspiron-5567:~/Documents/CRF++NER/Testing$
```

Figure 3: CRF++ testing using the wikigold.conll corpus

4 Evaluation

For the purposes of the evaluation of the extracted model and, our whole Named Entity Recognition (NER) system, a python coded script (Crf_TestEvaluation_Accuracy.py) has been utilized with the crf_wikigold_conll_testoutput.txt file as an input. In this script the main task is to check if the last two words of each row are the same (i.e. compares the actual label and the predicted label). As it can be seen in the figure 4 below, the accuracy score of the system performed is approximately 85.48%.



```
CRF++Accuracy [~/PycharmProjects/CRF++Accuracy] - .../Crf_TestEvaluation_Accuracy.py [CRF++Accuracy] - PyCharm
File Edit View Navigate Code Refactor Run Tools VCS Window Help
CRF++Accuracy Crf_TestEvaluation_Accuracy.py
Run: Crf_TestEvaluation_Accuracy
/home/konstantinos/PycharmProjects/CRF++Accuracy/venv/bin/python /home/konstantinos/PycharmProjects/CRF++Accuracy/Crf_TestEvaluation_Accuracy.py
<Counted True> <Counted False> <Counted Words> <Counted Sentences>
33467 5684 39151 1840
----->
The accuracy is 85.48185231539425
----->
Process finished with exit code 0
```

Figure 4: Evaluation of the NER implemented system

5 Comparison with the state of the art

By comparing with state-of-the-art models, we can observe that the performance of the implemented system is quite good and comparable. More specifically, researchers have used Distant Supervision (DS) from Wikipedia (in Nguyen and Moschitti, 2011), a supervised modeled approach for Sentence-level Relation Extraction (SLRE), where they were based on kernel methods (KM) and Support Vector Machines (SVMs) where the best in class models were in light of. In that particular research, the tests demonstrated that the particular approach is vigorous to Web records, as well as high accuracy has accomplished and more specifically, 74.29% for f1 score measure utilizing 52 YAGO relations in [1]. As it can be seen, our accuracy is higher than that and thus our approach and model is robust and better than the aforementioned one.

However, our NER system's accuracy is lower than the robust accuracy of 95% achieved after that system's evaluation, by the researchers who utilized the knowledge base and the YAGO ontology (version 2008-w40-2) for 99 relations in [1]. It can be

observed that on average and compared to state-of-the-art models / approaches, the performance of the implemented model is considered quite good.

References

- [1] Truc-Vien T. Nguyen and Alessandro Moschitti: Joint Distant and Direct Supervision for Relation Extraction. In Proceedings of the 5th International Joint Conference on Natural Language Processing. Chiang Mai, Thailand, November 8 – 13, 2011 pages 732–740.
- [2] Miao Fan, Qiang Zhou and Thomas Fang Zheng: Distant Supervision for Entity Linking. In Proceedings of the 29th Pacific Asia Conference on Language, Information and Computation. Shanghai, China, October 30 - November 1, 2015, pages 732–740.
- [3] Mike Mintz, Steven Bills, Rion Snow, Dan Jurafsky: Distant Supervision for Entity Linking. In Proceedings of the 47th Annual Meeting of the ACL and the 4th IJCNLP of the AFNLP. Suntec, Singapore, 2-7 August 2009, pages 1003–1011.
- [4] Mihai Surdeanuy, Julie Tibshiraniy, Ramesh Nallapati, Christopher D. Manning. “Multi-instance Multi-label Learning for Relation Extraction”.
- [5] NLTK, " Extracting Information from Text" [Online]. Available at: <http://www.nltk.org/book/ch07.html/> Last accessed: 22 April 2018.