

Token-level noise in large Web corpora and non-destructive normalization for linguistic applications

Felix Bildhauer*
FU Berlin/SFB 632
felix.bildhauer
@fu-berlin.de

Roland Schäfer*
FU Berlin
roland.schaefer
@fu-berlin.de

1 Introduction

We discuss token-level noise in Web corpora, using our own COW2012 corpora (Schäfer and Bildhauer 2012), which are up to 9 billion tokens large, as examples. Web corpora pose unique problems by being (at least potentially) very large and noisy at the same time. We show that it is partly because of this noise that they open up completely new roads in linguistic research. What can be called “noise“ from a POS tagging perspective, might be considered valuable evidence in research on non-standard writing. Normalization therefore has to be implemented very carefully or non-destructively. To illustrate, we describe two types of normalization tools which we have developed as well as our improved architecture to make huge non-destructively normalized corpora available to linguists.

2 Token-level noise in DECOW2012

As an example of a huge and noisy corpus, we use our DECOW2012 corpus of German (9.1 billion tokens) which was compiled and processed from Web data in 2012.¹ Similar to the findings in Liu and Curran (2006) for other Web corpora, we have 63,569,767 different types (after tokenization), of which 39,988,127 are hapax legomena. Such figures are implausible for any single natural language, and there must be a huge number of “noisy tokens” from diverse sources in the corpus. We have assessed the sources and types of noisy tokens by extracting a sample of tokens which were unknown to the standard German model for TreeTagger (Schmid 1995). We found the distribution shown in Table 1. It is known that the usual proportion of over 50% of truly noisy tokens cause further damage (i.e., degraded accuracy) in the whole linguistic post-processing chain. A good example is the significantly lower POS tagger accuracy on Web data (e.g., Giesbrecht and Evert 2009, summary in Schäfer and Bildhauer 2013). While it is possible to

extend a tagger's lexicon (or train a whole new tagger model) to improve the accuracy on the 46.8% of rare but real words and the 1.2% of numbers, dealing with the true noisy tokens involves a whole series of challenges and, above all, design decisions. We discuss two components of our normalization chain: one which applies destructive normalization (dehyphenation) and one which applies non-destructive normalization (spelling correction).

Source	%	±% CI
misspelling	20.0	5.0
tokenizer error	17.6	4.7
non-word	7.6	3.3
foreign language	6.8	3.1
real rare word	46.8	6.2
number	1.2	1.3

Table1: Sources of tokens which are unknown to TreeTagger in DECOW2012 (with 95% confidence interval for the estimate); true noisy tokens are above the line; below the line are those for which the POS tagger should (in principle) have a solution

3 Two examples of Web corpus noise

Hard-coded hyphenation is not only found in OCR'ed documents, but also in Web corpora. Sources include, e.g., texts from word processors or PDFs pasted into content management systems. The number of noisy tokens from hyphenation in Web corpora is relatively small. In Table 1, hyphenated words are included in the 7.6% of non-words, which also include HTML markup and similar material. In a language like German, automatic normalization of hyphenated words must distinguish between (1) ordinary hyphenation words like *Seitenstreifen* (“hard shoulder”) which have to be concatenated while dropping the hyphen (→ *Seitenstreifen*), (2) compounds which are actually written with a hyphen like *Philipps-Lagerverkauf* (“direct sale by Philips”) and were just accidentally ripped in halves at the hyphen position (→ *Philipps-Lagerverkauf*), and (3) abbreviated and coordinated compounds, which must never be normalized since they represent standard orthography: *weder TV- noch Radiosender* (“neither TV nor radio stations”).

First, it must be noted that simple approaches which rely on line endings being marked in the corpus (e.g., Grefenstette and Tapanainen 1994) will not work because the original line endings at which hyphenation occurred are usually absent in Web corpora. Our solution is HyDRA, an efficient (compiled) tool/library. It works in an unsupervised manner by generating frequency lists of unigrams and bigrams from a corpus in the training phase. In the production phase, it uses the database of ngrams and their frequencies as a primitive language model

* Both authors have made equal contributions.

¹ <http://www.corporafromtheweb.org/>

final product. We use the IMS Open Corpus Workbench (CWB)³, because it allows us to index our corpora in quite large chunks (of roughly 1.5 billion tokens) without major performance loss. For corpora of the size under discussion in this paper, we know no alternative. CWB has drawbacks, however, when it comes to non-destructive normalization. For example, multi-token units (MTU) and multi-unit tokens (MUT) cannot be represented adequately in CWB. To dehyphenate non-destructively, MUT capabilities would be required, and to render tokens like *undn* (cf. Section 4) in a POS tagger-friendly fashion, MTU capabilities would be required. In other words, transparently mapping *Seiten- streifen* to *Seitenstreifen* or *undn* to *und n* (with unequivocal tags for each token) in CWB is virtually impossible. Other architectures, like the ANNIS tool (Zeldes et al. 2009), which are suitable to deal with this kind of annotation, fail to perform well (or rather: at all) with corpora in the giga-token region.

We decided to use destructive normalization for dehyphenation with HyDRA because of the high accuracy of the tool and the low frequency of the phenomenon. Also, hard-coded hyphenation (mostly created by software) seems to us to be true noise, also from the linguistic perspective. On the problem with *undn*, cf. Section 6.

The increased size caused by normalization layers in the new version of our corpora (e.g., DECOW2013, scheduled for release in July 2013) makes it necessary to split the corpora in even smaller slices (roughly 20 for DECOW2013). This means that querying the whole corpus becomes quite cumbersome for linguists who cannot write their own scripts and cannot use parallelization. Around 20 single queries would have to be performed and merged manually for each lookup. Therefore, we are currently developing a simplistic Map-Reduce abstraction layer for CQP which executes CQP child processes on multiple machines. The tool does not require setting up a whole complicated cluster infrastructure (like Apache Hadoop). Processes across any number of machines communicate via a simple SSH connection and common NFS file systems. A Reduce is currently available for concordances, and we are working on a Reduce for CWB group results (combinatorial frequency tables).

6 Open problems and outlook

Some problems remain unsolved for the time being. To name just one example, POS tagging non-standard cliticized words like the aforementioned *undn* (“*and a*”) is impossible with most available POS tagger models. We are convinced that new tag sets and tagger models need to be developed which

accomplish the task of assigning some reasonable single POS tag to such forms and map them to the lemma of the non-clitic element (in the example, this is *und*). We hope to present first evaluations of our work on POS tagger improvement at the workshop, at least for English and German.

We also hope to present an evaluation of HyDRA and the spelling correction for our other Web-derived corpora (currently Danish, Dutch, English, French, Swedish). If there is significant interest, HyDRA can be released as a cross-platform library with a C ABI and C header files.

References

- Giesbrecht, E. and Evert, S. 2009. “Part-of-Speech (POS) Tagging -- a Solved Task? An Evaluation of POS Taggers for the German Web as Corpus”. In: Alegria, I. and Leturia, I. and Sharoff, S. (eds.) *Proceedings of the Fifth Web as Corpus Workshop (WAC5)*: 27-35
- Grefenstette, G. and Tapanainen, P. 1994. “What is a word? What is a sentence?” In: *Proceedings of 3rd Conference on Computational Lexicography and Text Research*.
- Liu, V. and Curran, J. R. 2006. “Web text corpus for natural language processing”. In: *Proceedings of the 11th Conference of the European Chapter of the Association for Computational Linguistics*. EACL: 233-240.
- Schäfer, R. and Bildhauer, F. 2012. “Building Large Corpora from the Web Using a New Efficient Tool Chain”. In Calzolari, N. and Choukri, K. and Declerck, T. and Doğan, M. U. and Maegaard, B. and Mariani, J. and Odijk, J. and Piperidis, S. (eds.) *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*. ELRA: 486-493.
- Schäfer, R. and Bildhauer, F. 2013, to appear. *Web Corpus Construction*. San Francisco: Morgan and Claypool.
- Schäfer, R. and Sayatz, U. 2013, submitted. “Die Kurzformen des Indefinitartikels im Deutschen”.
- Schmid, H. 1995. “Improvements in Part-of-Speech Tagging with an Application to German”. In: *Proceedings of the EACL SIGDAT-Workshop*, Dublin, Ireland.
- Subramaniam, L. V. and Roy, S. and Faruque, T. A. and Negi, S. 2009. “A survey of types of text noise and techniques to handle noisy text”. In: *Proceedings of The Third Workshop on Analytics for Noisy Unstructured Text Data*, ACM: 115-122.
- Zeldes, A. and Ritz, J. and Lüdeling, A. and Chiarcos, Chr. 2009. ANNIS: “A Search Tool for Multi-Layer Annotated Corpora”. In: *Proceedings of Corpus Linguistics 2009*.

³ <http://cwb.sourceforge.net/>