

Optimizing a PoS Tag Set for Dependency Parsing

Anonymous ACL submission

Abstract

Abstraction ensues.

1 Introduction

Part-of-speech (PoS) tagging is an important preprocessing step for many NLP tasks, such as dependency parsing (Nivre et al., 2007; Hajič et al., 2009), named entity recognition (Tjong Kim Sang and De Meulder, 2003) and sentiment analysis (Wilson et al., 2009). Whereas much effort has gone into the development of PoS-taggers – to the effect that this task is often considered more or less a solved task – considerably less effort has been devoted to the empirical evaluation of the PoS tag sets themselves. Error analysis of PoS-taggers indicate that, whereas tagging improvement through means of learning algorithm or feature engineering seems to have reached something of a plateau, linguistic and empirical assessment of the distinctions made in the PoS tag set may be a venue worth investigating further (Manning, 2011). Clearly, the utility of a PoS tag set is tightly coupled with the the downstream task for which it is performed. Even so, PoS tag sets are usually employed in a “one size fits all” fashion, regardless of the requirements posed by the task which makes use of this information.

It is well known that syntactic parsing often benefits from quite fine-grained morphological distinctions (Zhang and Nivre, 2011; Seeker and Kuhn, 2013). Morphology interacts with syntax though phenomena such as agreement and case marking, and incorporating information on morphological properties of words can therefore often improved parsing performance. However, in a realistic setting where the aim is to automatically parse raw text, the generation of morphological information will often require a separate step of morphological analysis that can be quite costly.

In this paper, we optimize a PoS tag set for the task of data-driven dependency parsing of Norwegian. We report on a set of experiments where various morphological distinctions are introduced to the PoS annotations and evaluated both intrinsically, i.e. in terms of PoS-tagging accuracy, and extrinsically, in terms of parsing accuracy. Our results show that the introduction of morphological distinctions not present in the original tag set, whilst compromising tagger accuracy, actually leads to significantly improved parsing accuracy. This optimization allows us to bypass the additional step of morphological analysis, framing the whole preprocessing problem as a simple tagging task.

The article is structured as follows. We start out by reviewing previous work on tag set evaluation in Section 2, while Section 3 details the treebank that provides the initial gold annotations used for our experiments. Section 4 describes the experimental setup used in our work and Section 5 goes on to provide the results from our tag set optimization. Finally, Section 6 summarizes our main findings and discusses some avenues for future work.

2 Previous/Related work

There has been little previous work dedicated specifically to extrinsic evaluation of the effects of PoS tag sets on downstream applications. There has, however, been some work that evaluates the effects of PoS tag set granularity on PoS tagging. This includes investigation of the effects of PoS tag sets on tagging of Swedish (Megyesi, 2001; Megyesi, 2002) and English (MacKinlay, 2005).

Megyesi (2001) and Megyesi (2002) trained and evaluated a range of PoS taggers on the Stockholm-Umeå Corpus (SUC) (Gustafson-Capková and Hartmann, 2006), annotated with a tag set based on a Swedish version of PAROLE

tags totaling 139 tags. Furthermore, they investigated the effects of tag set size on tagging by mapping the original tag set into smaller subsets designed for parsing. They argue that a tag set with complete morphological tags may not be necessary for all NLP applications, for instance syntactic parsing. They found that the smallest tag set comprising 26 tags yields the lowest tagger error rate. However, for some of the taggers, augmenting the tag set with more linguistically informative tags may actually lead to a drop in error rate. They argue that this shows that the size of the tag set as well as the type of information in the tags are crucial factors for tagger performance. However/unfortunately, they do not report results of parsing with the various PoS tag sets.

Similarly, MacKinlay (2005) investigated the effects of PoS tag sets on tagger performance in English, specifically the Wall Street Journal portion of Penn Treebank (Marcus et al., 1993). They mapped the original tag set of Penn Treebank to more fine-grained tag sets using linguistic insight to investigate whether additional linguistic information included in finer-grained tags could assist the tagger. Experimenting with both lexically and syntactically conditioned modifications such as distinguishing between count nouns and non-count nouns and between transitive and intransitive verbs, they found that more fine-grained tag sets rarely led to improvements in tagger accuracy; the most successful modification yielded an improvement in tagger accuracy of 0.05 percentage points.

Transition/overgang her ...

3 The Norwegian Dependency Treebank

We used the newly developed Norwegian Dependency Treebank (NDT) (Solberg et al., 2014), the first publicly available treebank for Norwegian. It was developed at the National Library of Norway in collaboration with the University of Oslo, and contains manual syntactic and morphological annotation. The treebank contains data from both varieties of Norwegian; 311 000 tokens of Bokmål and 303 000 tokens of Nynorsk. We will in the following only be using the Bokmål portion of the treebank. The annotated texts are mostly newspaper text, but also include government reports, parliament transcripts and excerpts from blogs. The annotation process of the treebank was supported by the rule-based Oslo-Bergen Tagger (Hagen et

Tag	Description
adj	Adjective
adv	Adverb
det	Determiner
inf-merke	Infinitive marker
interj	Interjection
konj	Conjunction
prep	Preposition
pron	Pronoun
sbu	Subordinate conjunction
subst	Noun
ukjent	Unknown (foreign word)
verb	Verb

Table 1: Overview of the PoS tag set of NDT.

al., 2000) and then manually corrected by annotators, adding syntactic dependency analysis to the morphosyntactic annotation.

Morphological Annotation The morphological annotation and PoS tag set of NDT follows the Oslo-Bergen Tagger (Hagen et al., 2000; Solberg, 2013), which in turn is largely based on the work of Faarlund et al. (1997). The tag set consists of 12 morphosyntactic PoS tags, with 7 additional tags for punctuation and symbols. The tag set is thus rather coarse-grained, with broad categories such as *subst* (noun) and *verb* (verb). The PoS tags are complemented by a large set of morphological features, providing information about morphological properties such as definiteness, number and tense. These features are used in our tag set modifications, where the coarse PoS tag of relevant tokens is concatenated with one or more of these features to include more linguistic information in the tags.

Syntactic Annotation The syntactic annotation choices in NDT are largely based on the Norwegian Reference Grammar (Faarlund et al., 1997). The annotation choices are outlined in Table 2, taken from/courtesy of Solberg et al. (2014), providing overview of the analyses of syntactic constructions that often distinguish dependency treebanks, such as coordination and the treatment of auxiliary and main verbs. The set of dependency relations comprises 29 dependency relations, including ADV (adverbial), SUBJ (subject) and KOORD (coordination).

Head	Dependent
Preposition	Prepositional complement
Finite verb	Complementizer
First conjunct	Subsequent conjuncts
Finite auxiliary	Lexical/main verb
Noun	Determiner

Table 2: Central head-dependent annotation choices in NDT.

4 Experimental Setup

In preparation to conducting our experiments with linguistically motivated tag set modifications, a concrete setup for the experiments needed to be established, which is presented in the following.

Data Set Split As there was no standardized data set split of NDT due to its very recent development, we needed to establish a data set split (training/development/test). Our data set split of the treebank follows the standard 80-10-10 (training/development/test) split and will be distributed with the treebank and proposed as the new standard(?). In creating the data set split, care has been taken to preserve contiguous texts in the various data sets while keeping the split balanced in terms of genre (and source). Our proposed data set split was used in the Norwegian contribution to the Universal Dependencies project (Øvrelid and Hohle, 2016). The split will be made available at a companion website.

Tagger For our experiments with tag set modifications, we wanted a PoS tagger both reasonably fast and accurate. There is often a trade-off between the two factors/considerations, as the most accurate taggers tend to suffer in terms of speed due to their complexity. However, a tagger that achieves both close to state-of-the-art accuracy as well as very high speed is TnT (Brants, 2000). TnT was furthermore recently used to evaluate the recently proposed universal tag set (Petrov et al., 2012). The sum of these factors led to TnT being the tagger of choice for our experiments.

Parser In choosing a syntactic parser for our experiments, we considered previous work on dependency parsing of Norwegian, specifically that of (Solberg et al., 2014). They found the Mate parser (Bohnet, 2010) to be the most successful parser for the parsing of NDT. Furthermore, recent dependency parser comparisons (Choi et al.,

2015) showed that Mate performed very well on parsing of English, beating a range of contemporary state-of-the-art parsers.

Tag Set Mapping In order to carry out the tag set modifications, we created a mapping that takes as input a set of tags and associated morphological features to be included in said tags. This mapping maps the relevant existing tags to new, more fine-grained tags including more relevant morphological features for the applicable tokens. This is done by concatenating the original tag with said features for all applicable tokens, i.e., tokens that are assigned said tag and contains the set of features in its morphological features.

Baseline It is common practice to compare the performance of PoS taggers to a pre-computed baseline for an initial point of comparison. For PoS tagging, a commonly used baseline is the Most Frequent Tag (MFT) baseline, which we use in our experiments. This involves labeling each word with the tag it was assigned most frequently in the training. All unknown words, i.e., words not seen in the training data, are assigned the tag most frequently assigned to words seen only once in the training. Unknown and infrequent words have in common that they rarely occur, and we might therefore expect them to have similar properties.

Tags & Features As we seek to quantify the effects of PoS tagging in a realistic setting, i.e., in application to raw text, we want to evaluate the parser on automatically assigned PoS tags. For the training of the parser, however, we have two options: using either gold standard or automatically assigned tags. In order to settle on a configuration, we conducted experiments with gold standard and automatically assigned tags to see how they differ with respect to performance. The results of our experiments reveal that the combination of training and testing on automatic tags is superior to training on gold standard tags and testing on automatic tags. Consequently, the parser was both trained and tested on automatically assigned tags in our experiments.

We removed any morphological features in order to simulate a realistic setting. Moreover, it is crucial that we remove these features when working with automatically assigned tags, as the features is still gold standard. For instance, if a verb token is erroneously tagged as a noun, we

could potentially have a noun token with verbal features such as tense, which markedly obfuscates the training and parsing. Another important aspect is that we want to isolate the effect of PoS tags, necessitating the exclusion of morphological features.

In the following, we report on a set of experiments which evaluate increasingly fine-grained PoS tag sets. The PoS tag sets are motivated by morphosyntactic properties of Norwegian and are evaluated both in terms of tagger and parser accuracy.

5 Tag Set Optimization/Experiments

By introducing more fine-grained linguistically motivated distinctions in a tag set, we increase the linguistic information represented in the tags, which may assist the parser in recognizing and generalizing syntactic patterns. However, the addition of more linguistic information to the tags and thus a more fine-grained tag set will most likely lead to a drop in tagger accuracy/tagging quality due to the increase in complexity. The best tagging does not necessarily lead to the best parse, and it is therefore interesting to investigate how the tag set modifications may affect the interplay between tagging and parsing. As we want to investigate how we can improve the linguistic quality of a tag set, we will not introduce tag set modifications which are not linguistically motivated. We will only consider distinctions we deem linguistically sensible, even if we expect them to impair the computational tractability; investigating the interplay between these two considerations is ultimately our goal.

5.1 Initial/Baseline Experiments

In an initial round of experiments, we concatenated the tag of each token with its set of morphological features in order to map the original tag set to a new, more fine-grained tag set (hereafter referred to as the *full* tag set). The result of this was a total of 368 tags, which is clearly very fine-grained. The two initial tag sets, i.e., the original tag set comprising 19 tags and the full tag set comprising 368 tags, thus represent two extremes in terms of granularity. To evaluate these tag sets and investigate how the tag set granularity would affect the performance of tagging and parsing, we trained and evaluated TnT and Mate on the training and development data, respectively, on the two

Tag set	MFT	Accuracy	LAS	UAS
Original	94.14%	97.47%	87.01%	90.19%
Full	85.15%	93.48%	87.15%	90.39%

Table 3: Results of tagging and parsing with the two initial tag sets.

tag sets. In Table 3, we report the results of these experiments. We see that the tagger accuracy drastically drops when going from the original to the full tag set. TnT reports an accuracy of 97.47% on the original tag set, which is reduced to 93.48% for the full tag set. These results confirm our hypothesis that the very high linguistic quality in the full, fine-grained tag set comes at the expense of drops in tagger performance. However, the additional linguistic information provided by the full tag set improves the parser performance. With the original tag set, Mate reports an LAS of 87.01% and a UAS of 90.19%, which increases to 87.15% and 90.39%, respectively, when using the full tag set. As we are looking for linguistically informed distinctions that improve the syntactic parsing, these results are promising and serve to indicate that additional morphological information assists the syntactic parsers, motivating the optimization of the existing PoS tag set.

5.2 Tag Set Experiments

We modified the tags for nouns, verbs, adjectives, determiners and pronouns in NDT by appending selected sets of morphological features to each tag in order to increase the linguistic information expressed by the tags. For each tag, we first experimented with each of the features in isolation before employing various combinations of them. We based our choices of combinations on how promising the features are and what we deem worth investigating in terms of linguistic utility, in order to see how the features might interact.

The morphological properties of the various parts-of-speech is reflected in the morphological features associated with the respective PoS tags. For instance, as nouns can take on gender, definiteness and number, additionally "being of" genitive case, the treebank operates with features for gender, definiteness, number and case with "accompanying" values. In addition to morphological properties such as definiteness, tense and number, all classes except for verbs can be distinguished "on" type, e.g., common nouns and proper nouns, while

pronouns can be reflexive, reciprocal, personal or interrogative.

5.3 Experiments

"Nearly" all tag set modifications leads to a drop in tagger accuracy, as expected. Few exceptions...

Nouns In Norwegian, nouns are assigned gender (feminine, masculine, neuter), definiteness (indefinite or definite) and number (singular or plural). There is agreement in gender, definiteness and number between nouns and their modifiers (adjectives and determiners). Additionally, NDT has a separate case feature to distinguishing nouns in genitive case.

Experimenting with nouns, we found that all tag set modifications yielded large increases in LAS and UAS. The most informative features are definiteness, which leads to an increase in LAS by 1.26 percentage points, to 88.27%, and type, yielding an LAS of 88.07%. Turning to combinations of features, we found that the combination of type and case, as well as type and definiteness, were the most promising, which led us to combine type, case and definiteness in a final experiment, resulting in LAS of 88.81% and UAS of 91.73%, constituting large increases from parsing with the original tag set, 1.80 percentage points and 1.54 percentage points, respectively.

Verbs Verbs are inflected for tense (infinitive, present, preterite, past perfect) in Norwegian and can additionally take on mood (imperative, indicative) and voice (active or passive). Note that voice and mood have only a single value, `pass` (passive) and `imp` (imperative), respectively. Verbs which are not passive are implicitly active, and verbs which are not imperative are indicative.

Introducing features in isolation, mood is the only feature leading to increases in LAS, with a reported LAS of 87.04%. Imperative clauses are fundamentally different from indicative clauses, as they lack an overt subject. Tense yields an LAS of 86.97% while voice yields an LAS of 86.96%. Combining mood and tense resulted in an LAS of 87.12% and UAS of 90.31%.

In an additional experiment, we mapped the verb tenses (mood, in the case of imperative) to finiteness. All verbs have finiteness, hence this distinction has broad coverage. This mapping is syntactically grounded as finite verbs and nonfinite verbs appear in completely different syntactic construction, and proved to greatly improve the

Feature	MFT	Accuracy	LAS	UAS
—	94.14%	97.47%	87.01%	90.19%
Definiteness	94.13%	97.49%	87.30%	90.42%
Gender	93.86%	97.28%	87.09%	90.31%
Number	94.06%	97.49%	87.04%	90.18%
Type	94.14%	97.61%	87.00%	90.11%

Table 4: Results of experiments with modified PoS tags for determiners.

parsing, as we saw the overall largest parser accuracy scores, with 87.30% for LAS and 90.43% for UAS, 0.29 and 0.24 percentage points higher than the baseline, respectively. This coincides with the observations seen for Swedish in Øvrelid (2008), where finiteness was found to be a very beneficial linguistic feature for parsing.

Adjectives Adjectives agree with the noun they modify in terms of gender, number and definiteness in Norwegian. Adjectives are also inflected for degree, either positive, comparative or superlative.

All features except for number (LAS of 86.99%) led to increases in parser accuracy scores. Degree being the most promising with a reported LAS of 87.29%, while distinguishing adjectives on definiteness yield an LAS of 87.14% and gender leads to LAS of 87.10%.

Combining the most promising features, we found that none of the combinations surpass the LAS with the distinction of degree alone. Definiteness and degree with an LAS of 87.23% and definiteness and number with an LAS of 87.27% and UAS identical to that of degree alone. The most fine-grained distinction, definiteness, degree and number, yields an LAS of 87.14%.

Determiners Like adjectives, determiners in Norwegian agree with the noun they modify in terms of gender, number and definiteness. Possessive pronouns are treated/analyzed as determiners in NDT.

Pronouns Pronouns in Norwegian include personal, reciprocal, reflexive and interrogative. They can exhibit gender, number and person. Personal pronouns have case (accusative or nominative).

Category	Feature(s)	MFT	Accuracy	LAS	UAS
Baseline	—	94.14%	97.47%	87.01%	90.19%
Noun	Type, case & definiteness	89.61%	97.05%	88.81%	91.73%
Verb	Finiteness	93.72%	97.35%	87.30%	90.43%
Adjective	Degree	94.13%	97.41%	87.29%	90.44%
Determiner	Definiteness	94.13%	97.49%	87.30%	90.42%
Pronoun	Type & case	94.12%	97.51%	87.30%	90.41%

Table 5: Results of tagging and parsing with the most successful tag set modification for each category.

6 Summary/Conclusion and Future Work

References

Tag	Description	References
adj komp	Comparative adjective	Bernd Bohnet. 2010. Very High Accuracy and Fast Dependency Parsing is not a Contradiction. In <i>Proceedings of the 23rd International Conference on Computational Linguistics</i> , pages 89–97, Beijing, China.
adj pos	Positive adjective	
adj sup	Superlative adjective	
det be	Definite determiner	Thorsten Brants. 2000. TnT - A Statistical Part-of-Speech Tagger. In <i>Proceedings of the Sixth Applied Natural Language Processing Conference</i> , Seattle, WA, USA.
det ub	Indefinite determiner	
pron pers	Personal pronoun	
pron pers akk	Personal pronoun, accusative	
pron pers nom	Personal pronoun, nominative	
pron refl	Reflexive pronoun	
pron res	Reciprocal pronoun	Jinho D. Choi, Joel Tetreault, and Amanda Stent. 2015. It Depends: Dependency Parser Comparison Using A Web-Based Evaluation Tool. In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics</i> , pages 387–396, Beijing, China.
pron sp	Interrogative pronoun	
subst appell	Common noun	
subst appell be	Common noun, definite	
subst appell be gen	Common noun, definite, genitive	
subst appell ub	Common noun, indefinite	
subst appell ub gen	Common noun, indefinite, genitive	Jan Terje Faarlund, Svein Lie, and Kjell Ivar Vannebo. 1997. <i>Norsk referansegrammatikk</i> . Universitetsforlaget, Oslo, Norway.
subst prop	Proper noun	
subst prop gen	Proper noun, genitive	
verb fin	Finite verb	Sofia Gustafson-Capková and Britt Hartmann. 2006. Manual of the Stockholm Umeå Corpus version 2.0.
verb infin	Nonfinite verb	

Table 6: Overview of the optimized tag set.

Tag set	MFT	Accuracy	LAS	UAS
Original	94.14%	97.47%	87.01%	90.19%
Full	85.15%	93.48%	87.15%	90.39%
Optimized	89.20%	96.85%	88.87%	91.78%

Table 7: Results of tagging and parsing with the optimized tag set, compared to the initial tag sets.

- Mitchell Marcus, Beatrice Santorino, and Mary Ann Marcinkiewicz. 1993. Building A Large Annotated Corpus of English: The Penn Treebank. Technical report, University of Philadelphia, Philadelphia, PA, USA.
- Beáta Megyesi. 2001. Comparing Data-Driven Learning Algorithms for PoS Tagging of Swedish. In *Proceedings of the 2001 Conference on Empirical Methods in Natural Language Processing*, pages 151–158, Pittsburgh, PA, USA.
- Beáta Megyesi. 2002. *Data-Driven Syntactic Analysis: Methods and Applications for Swedish*. Ph.D. thesis, Royal Institute of Technology, Stockholm, Sweden.
- Joakim Nivre, Johan Hall, Sandra Kübler, Ryan McDonald, Jens Nilsson, Sebastian Riedel, and Deniz Yuret. 2007. CoNLL 2007 Shared Task on Dependency Parsing. In *Proceedings of the 6th Conference on Natural Language Learning*, pages 915–932.
- Slav Petrov, Dipanjan Das, and Ryan McDonald. 2012. A Universal Part-of-Speech Tagset. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation*, pages 2089–2096, Istanbul, Turkey.
- Wolfgang Seeker and Jonas Kuhn. 2013. Morphological and syntactic case in statistical dependency parsing. *Computational Linguistics*, 39(1):23–55.
- Per Erik Solberg, Arne Skjærholt, Lilja Øvrelid, Kristin Hagen, and Janne Bondi Johannessen. 2014. The Norwegian Dependency Treebank. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation*, pages 789–795, Reykjavik, Iceland.
- Per Erik Solberg. 2013. Building Gold-Standard Treebanks for Norwegian. In *Proceedings of the 19th Nordic Conference of Computational Linguistics*, pages 459–464, Oslo, Norway.
- Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the conll-2003 shared task: Language-independent named entity recognition. In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003 - Volume 4*, pages 142–147, Stroudsburg, PA, USA.
- Theresa Wilson, Janyce Wiebe, and Paul Hoffman. 2009. Recognizing contextual polarity: An exploration of features for phrase-level sentiment analysis. *Computational Linguistics*, 35(3):399–433.
- Yue Zhang and Joakim Nivre. 2011. Transition-Based Dependency Parsing with Rich Non-Local Features. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 188–193, Portland, OR, USA.
- Lilja Øvrelid and Petter Hohle. 2016. Universal Dependencies for Norwegian. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation*, Portorož, Slovenia.
- Lilja Øvrelid. 2008. Finite Matters: Verbal Features in Data-Driven Parsing of Swedish. In *Proceedings of the Sixth International Conference on Natural Language Processing*, Gothenburg, Sweden.