

Doubling DOP*

Data Oriented Parsing based on Double-DOP and DOP*

Benno Kruit (10576223)
Sara Veldhoen (10545298)

January 14, 2014

1 Introduction

2 Existing Frameworks: Double-DOP and DOP*

A PTSG consists of the symbolic grammar, i.e. a set of fragments, and the corresponding weights. In the general case of DOP, all fragments are extracted from all the trees in the treebank. The number of fragments is exponential in the length of the sentences, thus the total number of fragments extracted would be far too large for efficient computation. Later models have therefore restricted the set of fragments in the grammar, thus improving computational efficiency. However, determining a proper subset of fragments to use is not trivial and this choice may negatively influence the performance of the grammar.

In this section, we outline two approaches to constrain the extraction of fragments: Double-DOP and DOP*. Furthermore, we discuss the similarities and dissimilarities for these two approaches.

2.1 Double-DOP

In the following, we discuss Double-DOP as it was presented in [1]. In this model, no unique fragments are extracted from the dataset: if a construction occurs in one tree only, it is probably not representative for the language. This is carried out by a dynamic algorithm using tree-kernels. It iterates over pairs of trees in the treebank, looking for fragments they have in common. In addition, only the largest shared fragment is stored.

The symbolic grammar that is the output of this algorithm is not guaranteed to derive each tree in the training corpus. Therefore all one-level fragments, constituting the set of PCFG-productions, are also added.

After the extraction of the symbolic grammar, the weights are obtained. This is done in a second pass over the treebank, assessing the relative frequencies.

The Double-DOP model has its main focus on determining the symbolic grammar. However, it was implemented with different estimators and maximizing objectives. Empirical results show that

2.2 DOP*

In DOP*[2], a rather different approach is taken called held-out estimation. The treebank is split in two parts, the extraction corpus (*EC*) and a held-out corpus (*HC*). An initial set of

fragments is extracted from the *EC*, containing all the fragments from its trees. The weights are then determined so as to maximize the likelihood of *HC*, under the assumption that this is equivalent to maximizing the joint probability of the *shortest derivations* of the trees in *HC*. All fragments that do not occur in such a derivation are removed from the symbolic grammar. Note that some trees in *HC* may not be derivable at all.

Consistency and bias DOP* is claimed to be the first consistent (non-trivial) DOP-estimator, [2] provides a consistency proof. On the other hand DOP* is biased, but Zollmann shows how bias actually arises from generalization: no non-overfitting DOP estimator could be unbiased. Bias is therefore not prohibited but on the contrary a desirable property of an estimator.

In [3] it is argued that there is a problem with the consistency proof given for DOP*, as well as the non-consistency proof for other DOP-estimators by [4]. Zuidema points out that these proofs use a frequency-distribution test, whereas for DOP a weight-distribution test would be more appropriate.

2.3 Comparison

Both DOP* and Double-DOP restrict the symbolic grammar based on some notion of reoccurrence of fragments.

In Double-DOP this is evident, it is explicit in the algorithm how (largest) reoccurring fragments are added to the grammar. From a computational point of view, this approach is very intuitive and can be implemented rather efficiently. However, the threshold (two) on the number of reoccurrences might seem rather trivial. Indeed, [1] reports experiments varying this threshold that show how performance drops for higher thresholds, with on the other hand a great reduction of the size of the grammar. A proper setting might well depend on the size and nature of the treebank used, and the computational costs one is willing to pay. In short, Double-DOP is computationally attractive but its theoretical foundation is not convincing us.

Theoretically, DOP* is more appealing: we decide on the symbolic grammar by assessing which fragments are actually used in derivations. The main assumption is that shortest derivations are preferred. The splitting of the treebank makes this possible (otherwise the grammar would be terribly overfitted).

2.4 Preview

DOP* is proven to be consistent and theoretically appealing. The computational problem of DOP* is that, at first, the entire set of possible fragments is extracted from *EC* and this set is reduced in a later phase. This comes with the need for huge storage and computation. In the next section, we investigate the possibility of reformulating the DOP* approach with insights from the Double-DOP model to maintain the best of both worlds.

3 A new implementation

We propose the following algorithm to perform a DOP*-like estimation using the efficient tree-kernel approach from Double-DOP.

Data: non-overlapping treebanks EC and HC

Result: a PTSG, i.e. a map of tree fragments and corresponding weights

```
1: Initialize  $M$ , a map from fragments to weights, empty
2: for  $t_H \in HC$  do
3:    $Q \leftarrow frag(t_H)$ 
4:    $F \leftarrow \emptyset$ 
5:   while  $Q$  not empty do
6:     for  $t_E \in EC$  do
7:       for  $f \in Q \cap frag(t_E)$  do
8:          $Q \leftarrow Q / \{f\}$ 
9:          $F \leftarrow Q \cup \{f\}$ 
10:      end for
11:    end for
12:  end while
13:   $D \leftarrow$  the shortest derivation of  $t_H$  using fragments in  $F$ 
14:  for  $f \in D$  do
15:     $M[f] \leftarrow M[f] + 1$ 
16:  end for
17: end for
18: for  $n \in V_N$  do
19:   Add the counts of all fragments rooted at  $n$ 
20:   Normalize the weights of these fragments
21: end for
```

In words: Iterate over the trees in HC (lines 2-17): create a list F of all of its fragments that occur in EC (lines 3-12). Create the shortest derivation D from fragments in F (NB: this is not always possible, in which case $D = \emptyset$). Increment the counts of these fragments in the map M

In Double-DOP, we would iterate over pairs of trees. With treebank size n , that would require $n!$ comparisons. In our proposed algorithm, we need to iterate over each combination of a tree in HC with (worst case¹) all trees in EC . This requires $n/2 \times n/2$ comparisons. Furthermore, we immediately store the counts, so we don't need to iterate over the data again to assess the weights. We do however need to normalize the weights, such that the sum of weights of all fragment sharing the same root equals 1.

4 Results

5 Conclusion

References

- [1] F. Sangati and W. Zuidema, "Accurate parsing with compact tree-substitution grammars: Double-dop," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 84–95, Association for Computational Linguistics, 2011.

¹The condition in line 5 avoids unnecessary search for fragments. We might be able to pose a smarter constraint here, such that we only continue if we can possibly find fragments that lead to a shorter derivation

- [2] A. Zollmann and K. Sima'an, "A consistent and efficient estimator for data-oriented parsing," *Journal of Automata, Languages and Combinatorics*, vol. 10, no. 2/3, pp. 367–388, 2005.
- [3] W. Zuidema, "Theoretical evaluation of estimation methods for data-oriented parsing," in *Proceedings of the Eleventh Conference of the European Chapter of the Association for Computational Linguistics: Posters & Demonstrations*, pp. 183–186, Association for Computational Linguistics, 2006.
- [4] M. Johnson, "The dop estimation method is biased and inconsistent," *Computational Linguistics*, vol. 28, no. 1, pp. 71–76, 2002.