

האוניברסיטה העברית בירושלים – הפקולטה למתמטיקה ולמדעי הטבע המכון
למדעי המחשב

Neutralizing Linguistically Problematic Annotations in Unsupervised Dependency Parsing Evaluation

ניטרול תיוגים בעיתיים תחבירית
בהערכה של ניתוח תלות לא מפוקח

מוגש על ידי
רועי שוורץ

עבודה זו הונחתה על ידי
פרופ' ארי רפפורט

עבודת גמר לתואר מוסמך במדעי המחשב
א' באב תשע"א ; 01 לאוגוסט 2011

Roy Schwartz: *Neutralizing Linguistically Problematic Annotations*
in Unsupervised Dependency Parsing Evaluation, M.Sc. Thesis, © April
2011, Nissan 5771

Supervisor:
Prof. Ari Rappoport

School of Computer Science and Engineering

Hebrew University of Jerusalem
Jerusalem, Israel

Preface

The work presented in this thesis is largely based on a work presented in ACL-HLT 2011 conference (Schwartz et al., 2011). It is a joint work with Omri Abend and Dr. Roi Reichart. I would like to thank Omri and Roi for their brilliant insights and tremendous help, which has taught me much. I would also like to thank my advisor, Prof. Ari Rappoport, for his superb intuitions and support in bringing this project to life.

I would like to thank Shay Cohen for his assistance with his implementation of the DMV parser and Taylor Berg-Kirkpatrick, Phil Blunsom and Jennifer Gillenwater for providing us with their data sets. Last, I would also like to thank Valentin I. Spitzkovsky for his comments and for providing us with his data sets.

ניתוח תלות (Dependency Parsing) הינו אתגר מרכזי בתחום עיבוד שפה טבעית (Natural Language Processing). משימת ניתוח תלות כוללת תיוג אוטומטי של משפטים בשפה טבעית, כך שכל מילה מתויגת כתלויה במילה אחרת במשפט. מבנה דקדוקי זה הוא חשוב הן מבחינה בלשנית (Melčuk, 1988), והן בזכות היותו שימושי באפליקציות מסדר גבוה כגון תרגום מכונה (Ding and Palmer, 2004; Xu et al., 2009).

כלי ניתוח תלות אוטומטיים משתמשים לרוב בטכניקות של למידת מכונה (Machine Learning). באלגוריתמי למידה מסורתיים, כלים אוטומטיים מקבלים דוגמאות אימון מתויגות, ומשתמשים בהן על מנת להסיק הכללות שיעזרו לתייג דוגמא (לא מתויגת) חדשה (למידה מפוקחת). לאחרונה, כלים לא מפוקחים הולכים וצוברים פופולאריות – כלים אלו מקבלים דוגמאות לא מתויגות ומנסים למצוא בהן קשרים סטטיסטיים פנימיים, ובעזרתם לתייג דוגמאות חדשות. בעבודה זו נתמקד בגרסה הלא מפוקחת של ניתוח תלות (Unsupervised Dependency Parsing), אשר בה דוגמאות האימון הם משפטים בשפה טבעית.

למידה לא מפוקחת הינה מאתגרת למדי, שכן היא דורשת מהכלים האוטומטיים למצוא את הקשרים הפנימיים המשמשים לבניית סכימת התיוג, ללא צפייה ולו בדוגמא אחת מתויגת. אולם, במקרים מסוימים ישנה יותר מדרך אחת בלשנית סבירה לתיוג. במקרים כאלה, התיוג הנבחר ידנית על ידי מתייג אנושי אינו בהכרח מייצג קשרים עמוקים הקיימים בטקסט, אלא עלול להיות תוצאה של החלטה אקראית. מקרים אלו משפיעים בצורה ישירה על הערכת הביצועים של כלים לא מפוקחים: מכיוון שדוגמאות האימון אינן מצביעות בצורה ברורה לכיוון אף אחד מהתיוגים האפשריים, הכלים הלא מפוקחים אינם יכולים לנחש מהו התיוג "הנכון" (כלומר זה שנבחר על ידי המתייג האנושי). כתוצאה מכך, בחירה בתיוג "השגוי", על אף שהוא סביר בלשנית, יכולה לגרום לפגיעה חמורה בביצועי הכלי.

בעבודה זו אנו נתמקד במבני תלות בעייתיים – שעבורם אין תיוג תלות שנמצא תחת מוסכמה בלשנית. דוגמאות למבנים כאלה כוללות למשל מבני איחוד ("דני ודנה") ומבני פעלים ("יכול לאכול"), אשר יש יותר מדרך קבילה אחת לתייגם (Kübler et al., 2009), ולמרות זאת, "תקן הזהב" (Gold Standard) – שכנגדו מבוצעת הערכת הביצועים – מגדיר רק את אחד התיוגים כנכון.

נתחיל בהצגת ניסויים עם שלושה כלי ניתוח מובילים (Klein and Manning, 2004; Cohen and Smith, 2009; Spitkovsky et al., 2010a), ונראה שבכל אחד מהם קיימת תת-קבוצה קטנה של פרמטרים, אשר שינוי הערך שלהם יכול לגרום לשיפור משמעותי בביצועי כלי הניתוח על פי שיטות ההערכה הנפוצות. פרמטרים אלו שולטים על התיוג של מבנים מקומיים, בהם אין הסכמה בלשנית לגבי התיוג הנכון, ושבהם כלי הניתוח למדו תיוג "שגוי" (כלומר, תיוג סביר בלשנית השונה מהתיוג של תקן הזהב). כתוצאה מכך, אילוץ כלי הניתוח לבחור את התיוג "הנכון" – על ידי שינוי הפרמטרים השולטים על התיוג של מבנים אלו – מביא לשיפור משמעותי בביצועים. אולם, שיפור זה אינו מבטא שיפור אמיתי באיכות כלי התיוג, שכן שתי הגרסאות של הכלים (עם הפרמטרים המקוריים ועם הפרמטרים החדשים)

מתייגות בצורה זהה את כל מבני השפה, מלבד מספר מבנים בעייתיים, אותם שתי הגרסאות מתייגות בצורה שהיא סבירה מבחינה בלשנית.

נראה גם שכיווני מחקר שונים בשני העשורים האחרונים גרמו ליצירתם בפועל של שלושה תיוגים שונים לתקן הזהב, כאשר השוני בא לידי ביטוי בעיקר בתיוג של מבנים בעייתיים. המשמעות של ממצא זה היא שעבודות שונות שחקרו את אותה המטלה מעריכות את עצמן מול תקני זהב שונים, אך משוות את תוצאות ההערכה שלהן אחת מול השנייה.

משני הממצאים לעיל נסיק ששיטות ההערכה הקיימות אינן מסוגלות לספק אינדיקציה אמיתית לאיכות כלי ניתוח תלות.

לבסוף נציג אמצעי חדש להערכת ביצועים, *Neutral Edge Direction* (NED), ונראה שהוא מבטל חלק גדול מהתופעה הבלתי רצויה, ולפי כך מספק אינדיקציה טובה הרבה יותר לאיכותם של כלי ניתוח תלות.

Dependency parsing is a central task in the field of *Natural Language Processing* (*NLP*). The task involves the automatic labeling of natural language sentences with dependency structures, such that each word is labeled as the dependent of another word in the sentence (its *syntactic head*). This formalism is important, both in the linguistic aspect (Mel'čuk, 1988) and in empirical aspects, as it is useful in higher level NLP applications such as Statistical Machine Translation (Ding and Palmer, 2004; Xu et al., 2009).

Dependency parsing tools (*dependency parsers*) usually use machine learning techniques. In traditional learning algorithms, automatic tools receive labeled training data, and attempt at inferring a generalization of this data in order to annotate new (unlabeled) data (*supervised learning*). Recently, there has been a rise in the popularity of *unsupervised* algorithms, where automatic tools are trained on plain text only (unlabeled data) and induce a labeling for a given test sentence. We focus on the unsupervised setting in this work.

Unsupervised learning is a challenging task, as it requires the automatic tools to establish the inter-connections in the data that are used to build the annotation. However, in some cases, there is more than one reasonable way to annotate the data. In such cases, the annotation does not reflect profound data correlations, but may be the result of some arbitrary decision which the manual annotator is forced to make. These cases affect unsupervised tools drastically. As the data does not always clearly point out towards one specific annotation, unsupervised tools have no way of guessing the “correct” annotation (as they do not see any labeled data during the training phase). As a result, making the wrong (plausible) choice could result in a serious performance decrease.

In this work we consider an example of such a phenomenon, and focus on *problematic* dependency structures, for which the dependency annotation is not under linguistic consensus. Examples include coordination structures (e.g., “John and Mary”) and infinitive verbs (e.g., “to eat”), which both have more than one linguistically plausible dependency annotation (Kübler et al., 2009). Although there is more than one way to annotate the data, the gold standards against which dependency parsers are evaluated are forced to choose one of these annotations.

We start by experimenting with three leading unsupervised parsers (Klein and Manning, 2004; Cohen and Smith, 2009; Spitkovsky et al., 2010a), and show that

for each of them, a small set of parameters can be found whose modification yields a significant improvement in the standard evaluation measures. These parameters correspond to local cases where no linguistic consensus exists as to the correct annotation, and in which the parsers learned a wrong annotation (that is, a plausible annotation which is not the same as the gold standard). As a result, forcing the parsers to choose the “right” annotation, by modifying the parameters controlling the annotations of these structures, yields a substantial improvement.

Nevertheless, both versions of the parsers (the one with the original parameters and the one with the modified ones) tag every language structure the same, except for problematic structures, where both take a linguistically plausible approach. As a result, this performance improvement does not represent true quality improvement.

We also show that three substantially different gold standards annotations, which differ from one another mainly in problematic structures, are currently in use for this task. That is, different works on the same task evaluate themselves against different gold standard annotations.

We conclude that the standard unsupervised dependency parsing evaluation does not always provide a true indication of algorithm quality.

Finally, we present a novel measure, *Neutral Edge Direction* (NED), and show that it greatly reduces these undesired phenomena, thus providing a better indication of algorithm quality.

Contents

1	Introduction	i
2	Related Work	iii
2.1	Dependency Model with Valence (DMV)	v
3	Significant Effects of Edge Flipping	vii
3.1	Parameter Changes	ix
3.2	Experimenting with Edge Flipping	xi
4	Linguistically Problematic Annotations	xiii
4.1	Evaluation Inconsistency Across Papers	xiv
5	The Neutral Edge Direction (NED) Measure	xvi
5.1	Undirected Evaluation	xvi
5.2	Neutral Edge Direction (NED)	xvii
5.3	Discussion	xviii
6	Experimenting with NED	xxi
7	Discussion	xxiv
8	Conclusion	xxviii

Introduction

Unsupervised induction of dependency parse trees is a major NLP task that attracts a substantial amount of research (Klein and Manning, 2004; Cohen et al., 2008; Headden III et al., 2009; Spitkovsky et al., 2010a; Gillenwater et al., 2010; Berg-Kirkpatrick et al., 2010; Blunsom and Cohn, 2010, *inter alia*). Parser quality is usually evaluated by comparing its output to a gold standard whose annotations are linguistically motivated. However, there are cases in which there is no linguistic consensus as to what the correct annotation is (Kübler et al., 2009). Examples include which verb is the head in a verb group structure (e.g., “can” or “eat” in “can eat”), and which noun is the head in a sequence of proper nouns (e.g., “John” or “Doe” in “John Doe”). We refer to such annotations as (*linguistically*) *problematic*. For such cases, evaluation measures should not punish the algorithm for deviating from the gold standard.

In this work we show that the evaluation measures reported in current works are highly sensitive to the annotation in problematic cases, and propose a simple new measure that greatly neutralizes the problem.

We start from the following observation: for three leading algorithms (Klein and Manning, 2004; Cohen and Smith, 2009; Spitkovsky et al., 2010a), a small set (at most 18 out of a few thousands) of parameters can be found

whose modification dramatically improves the standard evaluation measures (the attachment score measure by 9.3-15.1%, and the undirected measure by a smaller but still significant 1.3-7.7%). The phenomenon is implementation independent, occurring with several algorithms based on a fundamental probabilistic dependency model.

We show that these parameter changes can be mapped to edge direction changes in local structures in the dependency graph, and that these correspond to problematic annotations. Thus, the standard evaluation measures do not reflect the true quality of the evaluated algorithm.

We explain why the standard undirected evaluation measure is in fact sensitive to such edge direction changes, and present a new evaluation measure, *Neutral Edge Direction* (NED), which greatly alleviates the problem by ignoring the edge direction in local structures. Using NED, manual modifications of model parameters always yields small performance differences. Moreover, NED sometimes punishes such manual parameter tweaking by yielding worse results. We explain this behavior using an experiment revealing that NED always prefers the structures that are more consistent with the modeling assumptions lying in the basis of the algorithm. When manual parameter modification is done against this preference, NED results decrease.

The contributions of this work are as follows. First, we show the impact of a small number of annotation decisions on the performance of unsupervised dependency parsers. Second, we observe that often these decisions are linguistically controversial and therefore this impact is misleading. This reveals a problem in the common evaluation of unsupervised dependency parsing. This is further demonstrated by noting that recent papers evaluate the task using three gold standards which differ in such decisions and which yield substantially different results. Third, we present the NED measure, which is agnostic to errors arising from choosing the non-gold direction in such cases.

Chapter 2

Related Work

Grammar induction received considerable attention over the years (see (Clark, 2001; Klein, 2005) for reviews). Dependency parsing is the task of inducing a dependency tree from a plain text sentence (see (Mel’čuk, 1988) for a theoretical background of this task). Figure 2.1 presents the dependency parse of the sentence “I want to eat”. As the figure shows, the word “want” is the root of this sentence. It has two dependents: “I” and “eat”, which also has a dependent – “to”.

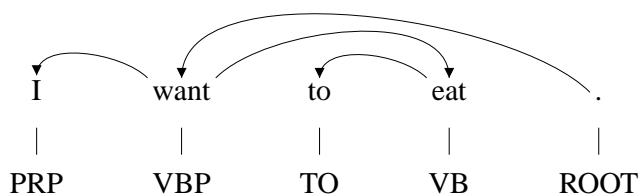


Figure 2.1: A dependency parse tree of the sentence “I want to eat”. As the figure shows, the word “want” is the root of this sentence. It has two dependents: “I” and “eat”, which also has a dependent – “to”.

For unsupervised dependency parsing, the *Dependency Model with Valence (DMV)* (Klein and Manning, 2004) was the first to beat the simple right-branching baseline. A technical description of DMV is given in Sec-

tion 2.1.

The great majority of recent works, including those experimented with in this work, are elaborations of DMV. Smith and Eisner (2005) improved the DMV results by generalizing the function maximized by DMV’s EM training algorithm. Smith and Eisner (2006) used a structural locality bias, experimenting on five languages. Cohen et al. (2008) extended DMV by using a variational EM training algorithm and adding logistic normal priors. Cohen and Smith (2009, 2010) further extended it by using a *shared* logistic normal prior which provided a new way to encode the knowledge that some POS tags are more similar than others. A bilingual joint learning further improved their performance.

Headden III et al. (2009) obtained the best reported results on WSJ10 by using a lexical extension of DMV. Gillenwater et al. (2010) used posterior regularization to bias the training towards a small number of parent-child combinations. Berg-Kirkpatrick et al. (2010) added new features to the M step of the DMV EM procedure. Berg-Kirkpatrick and Klein (2010) used a phylogenetic tree to model parameter drift between different languages. Spitkovsky et al. (2010a) explored several training protocols for DMV. Spitkovsky et al. (2010c) showed the benefits of Viterbi (“hard”) EM to DMV training. Spitkovsky et al. (2010b) presented a novel *lightly-supervised* approach that used hyper-text mark-up annotation of web-pages to train DMV.

A few non-DMV-based works were recently presented. Daumé III (2009) used shift-reduce techniques. Blunsom and Cohn (2010) used tree substitution grammar to achieve best results on WSJ^∞ .

Druck et al. (2009) took a semi-supervised approach, using a set of rules such as “A noun is usually the parent of a determiner which is to its left”, experimenting on several languages. Naseem et al. (2010) further extended this idea by using a single set of rules which globally applies to six different

languages. The latter used a model similar to DMV.

The controversial nature of some dependency structures was discussed in (Nivre, 2006; Kübler et al., 2009). Klein (2005) discussed controversial constituency structures and the evaluation problems stemming from them, stressing the importance of a consistent standard of evaluation.

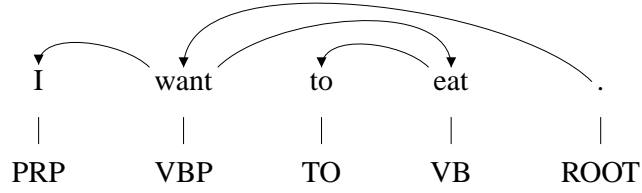
A few works explored the effects of annotation conventions on parsing performance. Nilsson et al. (2006) transformed the dependency annotations of coordinations and verb groups in the Prague TreeBank. They trained the supervised MaltParser (Nivre et al., 2006) on the transformed data, parsed the test data and re-transformed the resulting parse, thus improving performance. Klein and Manning (2004) observed that a large portion of their errors is caused by predicting the wrong direction of the edge between a noun and its determiner. Kübler (2005) compared two different conversion schemes in German supervised constituency parsing and found one to have positive influence on parsing quality.

2.1 Dependency Model with Valence (DMV)

DMV (Klein and Manning, 2004) defines a probabilistic grammar for unlabeled dependency structures. It is defined as follows: the root of the sentence is first generated, and then each head recursively generates its right and left dependents. The parameters of the model are of two types: P_{STOP} and P_{ATTACH} . $P_{STOP}(dir, h, adj)$ determines the probability to stop generating arguments, and is conditioned on 3 arguments: the head h , the direction dir ((*L*)*e**f**t* or (*R*)*i**g**h**t*) and adjacency adj (whether the head already has dependents ((*Y*)*e**s*) in direction dir or not ((*N*)*o*)). $P_{ATTACH}(arg|h, dir)$ determines the probability to generate arg as head h 's dependent in direction dir .

Figure 2.2 presents an example of the PTB gold standard parse of the sentence “I want to eat” and the DMV model probability for this parse.

vi



$$\begin{aligned}
 P = & (1 - P_{STOP}(ROOT, L, N)) \times P_{ATTACH}(VBP|ROOT, L) \\
 \times & (1 - P_{STOP}(VBP, L, N)) \times P_{ATTACH}(PRP|VBP, L) \\
 \times & (1 - P_{STOP}(VBP, R, N)) \times P_{ATTACH}(VB|VBP, R) \\
 \times & (1 - P_{STOP}(VB, L, N)) \times P_{ATTACH}(TO|VB, L) \\
 \times & P_{STOP}(VBP, L, Y) \times P_{STOP}(VBP, R, Y) \\
 \times & P_{STOP}(PRP, L, N) \times P_{STOP}(PRP, R, N) \\
 \times & P_{STOP}(VB, L, Y) \times P_{STOP}(VB, R, N) \\
 \times & P_{STOP}(TO, L, N) \times P_{STOP}(TO, R, N)
 \end{aligned}$$

Figure 2.2: The gold standard dependency parse of the sentence “I want to eat” and the DMV probability for this sentence.

Chapter 3

Significant Effects of Edge Flipping

In this chapter we present recurring error patterns in some of the leading unsupervised dependency parsers. These patterns are all local, confined to a sequence of up to three words (but mainly of just two consecutive words). They can often be mended by changing the directions of a few types of edges.

The modified parameters described in this chapter were handpicked to improve performance: we examined the local parser errors occurring the largest number of times, and found the corresponding parameters. Note that this is a valid methodology, since our goal is not to design a new algorithm but to demonstrate that modifying a small set of parameters can yield a major performance boost and eventually discover problems with evaluation methods or algorithms.

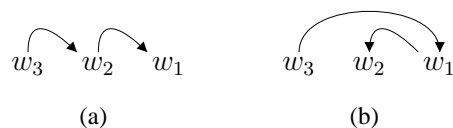


Figure 3.1: A dependency structure on the words w_1, w_2, w_3 before (Figure 3.1(a)) and after (Figure 3.1(b)) an *edge-flip* of $w_2 \rightarrow w_1$.

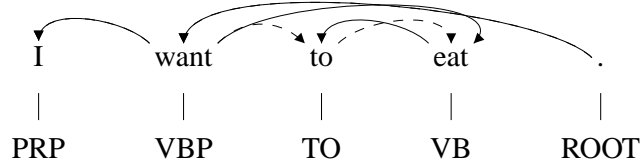


Figure 3.2: A parse of the sentence “I want to eat”, before (straight line) and after (dashed line) an *edge-flip* of the edge “to” \leftarrow “eat”.

We start with a few definitions. Consider Figure 3.1(a) that shows a dependency structure on the words w_1, w_2, w_3 . Edge flipping (henceforth, *edge-flip*) the edge $w_2 \rightarrow w_1$ is the following modification of a parse tree: (1) setting w_2 ’s parent as w_1 (instead of the other way around), and (2) setting w_1 ’s parent as w_3 (instead of the edge $w_3 \rightarrow w_2$). Figure 3.1(b) shows the dependency structure after the *edge-flip*.

Note that (1) imposes setting a new parent to w_2 , as otherwise it would have had no parent. Setting this parent to be w_3 is the minimal modification of the original parse, since it does not change the attachment of the structure $[w_2, w_1]$ to the rest of the sentence, but only the direction of the internal edge.

Figure 3.2 presents a parse of the sentence “I want to eat”, before and after an *edge-flip* of the edge “to” \leftarrow “eat”.

Since unsupervised dependency parsers are generally structure prediction models, the predictions of the parse edges are not independent. Therefore, there is no single parameter which completely controls the edge direction, and hence there is no direct way to perform an *edge-flip* by parameter modification. However, setting extreme values for the parameters controlling the direction of a certain edge type creates a strong preference towards one of the directions, and effectively determines the edge direction. This procedure is henceforth termed *parameter-flip*.

We show that by performing a few *parameter-flips*, a substantial improvement in the attachment score can be obtained. Results are reported for three

algorithms.

3.1 Parameter Changes

All the works experimented with in this work are not lexical and use sequences of POS tags as their input. In addition, they all use the DMV parameter set (P_{STOP} and P_{ATTACH}) for parsing. We will henceforth refer to this set, conditioned on POS tags, as the model parameter set.

We show how an edge in the dependency graph is encoded using the DMV parameters. Say the model prefers setting “to” (POS tag: TO) as a dependent of the infinitive verb (POS tag: VB) to its right (e.g., “to eat”). This is reflected by a high value of $P_{ATTACH}(TO|VB, L)$, a low value of $P_{ATTACH}(VB|TO, R)$, since “to” tends to be a left dependent of the verb and not the other way around, and a low value of $P_{STOP}(VB, L, N)$, as the verb usually has at least one left argument (i.e., “to”).

A *parameter-flip* of $w_1 \rightarrow w_2$ is hence performed by setting $P_{ATTACH}(w_2|w_1, R)$ to a very low value and $P_{ATTACH}(w_1|w_2, L)$ to a very high value. When the modifications to P_{ATTACH} are insufficient to modify the edge direction, $P_{STOP}(w_2, L, N)$ is set to a very low value and $P_{STOP}(w_1, R, N)$ to a very high value¹.

Table 3.1 describes the changes made for the three algorithms. The ‘+’ signs in the table correspond to edges in which the algorithm disagreed with the gold standard, and were thus modified. Similarly, the ‘-’ signs in the table correspond to edges in which the algorithm agreed with the gold standard, and were thus not modified. The number of modified parameters does not exceed 18 (out of a few thousands).

The *Freq.* column in the table shows the percentage of the tokens in sections 2-21 of PTB WSJ that participate in each structure. Equivalently, the

¹Note that this yields unnormalized models. Again, this is justified since the resulting model is only used as a basis for discussion and is not a fully fledged algorithm.

percentage of edges in the corpus which are of either of the types appearing in the *Orig. Edge* column. As the table shows, the modified structures cover a significant portion of the tokens. Indeed, 42.9% of the tokens in the corpus participate in at least one of them².

Structure	Freq.	Orig. Edge	<i>km04</i>	<i>cs09</i>	<i>saj10a</i>
Coordination ("John & Mary")	2.9%	$CC \rightarrow NNP$	–	+	–
Prepositional Phrase ("in the house")	32.7%	$DT \rightarrow NN$	+	+	+
		$DT \rightarrow NNP$	–	+	+
		$DT \rightarrow NNS$	–	–	+
		$IN \rightarrow DT$	+	+	–
		$IN \leftarrow NN$	+	+	–
		$IN \leftarrow NNP$	+	–	–
		$IN \leftarrow NNS$	–	+	–
		$PRP\$ \rightarrow NN$	–	–	+
Modal Verb ("can eat")	2.4%	$MD \leftarrow VB$	–	+	–
Infinitive Verb ("to eat")	4.5%	$TO \rightarrow VB$	–	+	+
Proper Name Sequence ("John Doe")	18.5%	$NNP \rightarrow NNP$	+	–	–

Table 3.1: Parameter changes for the three algorithms. The Freq. column shows what percentage of the tokens in sections 2-21 of PTB WSJ participate in each structure. The Orig. column indicates the original edge. The modified edge is of the opposite direction. The other columns show the different algorithms: *km04*: basic DMV model (replication of (Klein and Manning, 2004)); *cs09*: (Cohen and Smith, 2009); *saj10a*: (Spitkovsky et al., 2010a).

²Some tokens participate in more than one structure.

3.2 Experimenting with Edge Flipping

We experimented with three DMV-based algorithms: a replication of (Klein and Manning, 2004), as appears in (Cohen et al., 2008) (henceforth, *km04*), Cohen and Smith (2009) (henceforth, *cs09*), and Spitkovsky et al. (2010a) (henceforth, *saj10a*). Decoding is done using the Viterbi algorithm³. For each of these algorithms we present the performance gain when compared to the original parameters.

The training set is sections 2-21 of the Wall Street Journal Penn Tree-Bank (Marcus et al., 1993). Testing is done on section 23. The constituency annotation was converted to dependencies using the rules of (Yamada and Matsumoto, 2003)⁴.

Following standard practice, we present the attachment score (i.e., percentage of words that have a correct head) of each algorithm, with both the original parameters and the modified ones. We present results both on all sentences and on sentences of length ≤ 10 , excluding punctuation.

Table 3.2 shows results for all algorithms⁵. The performance difference between the original and the modified parameter set is considerable for all data sets, where differences exceed 9.3%, and go up to 15.1%. These are enormous differences from the perspective of current algorithm evaluation results.

³<http://www.cs.cmu.edu/~scohen/parser.html>.

⁴<http://www.jaist.ac.jp/~h-yamada/>

⁵Results are slightly worse than the ones published in the original papers due to the different decoding algorithms (*cs09* use MBR while we used Viterbi) and a different conversion procedure (*saj10a* used (Collins, 1999) and not (Yamada and Matsumoto, 2003)); see Chapter 5.

Algo.	≤ 10			$\leq \infty$		
	<i>Orig.</i>	<i>Mod.</i>	Δ	<i>Orig.</i>	<i>Mod.</i>	Δ
<i>km04</i>	45.8	59.8	14	34.6	43.9	9.3
<i>cs09</i>	60.9	72.9	12	39.9	54.6	14.7
<i>saj10a</i>	54.7	69.8	15.1	41.6	54.3	12.7

Table 3.2: Results of the original (*Orig.* columns), the modified (*Mod.* columns) parameter sets and their difference (Δ columns) for the three algorithms.

Chapter 4

Linguistically Problematic

Annotations

In this chapter, we discuss the controversial nature of the annotation in the modified structures (Kübler et al., 2009). We remind the reader that structures for which no linguistic consensus exists as to their correct annotation are referred to as (linguistically) problematic.

We begin by showing that all the structures modified are indeed linguistically problematic. We then note that these controversies are reflected in the evaluation of this task, resulting in three, significantly different, gold standards currently in use.

Coordination Structures are composed of two proper nouns, separated by a conjunctive (e.g., “John and Mary”). It is not clear which token should be the head of this structure, if any (Nilsson et al., 2006).

Prepositional Phrases (e.g., “in the house” or “in Rome”), where every word is a reasonable candidate to head this structure. For example, in the annotation scheme used by (Collins, 1999) the preposition is the head, in the scheme used by (Johansson and Nugues, 2007) the noun is the head, while TUT annotation, presented in (Bosco and Lombardo, 2004), takes the determiner to be the noun’s head.

Verb Groups are composed of a verb and an auxiliary or a modal verb (e.g., “can eat”). Some schemes choose the modal as the head (Collins, 1999), others choose the verb (Rambow et al., 2002).

Infinitive Verbs (e.g., “to eat”) are also in controversy, as in (Yamada and Matsumoto, 2003) the verb is the head while in (Collins, 1999; Bosco and Lombardo, 2004) the “to” token is the head.

Sequences of Proper Nouns (e.g., “John Doe”) are also subject to debate, as PTB’s scheme takes the last proper noun as the head, and BIO’s scheme defines a more complex scheme (Dredze et al., 2007).

4.1 Evaluation Inconsistency Across Papers

A fact that may not be recognized by some readers is that comparing the results of unsupervised dependency parsers across different papers is not directly possible, since different papers use different gold standard annotations *even when they are all derived from the Penn Treebank constituency annotation*. This happens because they use different rules for converting constituency annotation to dependency annotation. A probable explanation for this fact is that people have tried to correct linguistically problematic annotations in different ways, which is why we note this issue here¹.

There are three different annotation schemes in current use: (1) Collins head rules (Collins, 1999), used in e.g., (Berg-Kirkpatrick et al., 2010; Spitkovsky et al., 2010a); (2) Conversion rules of (Yamada and Matsumoto, 2003), used in e.g., (Cohen and Smith, 2009; Gillenwater et al., 2010); (3) Conversion rules of (Johansson and Nugues, 2007) used, e.g., in the CoNLL shared task 2007 (Nivre et al., 2007) and in (Blunsom and Cohn, 2010).

The differences between the schemes are substantial. For instance, 14.4%

¹Indeed, half a dozen flags in the LTH Constituent-to-Dependency Conversion Tool (Johansson and Nugues, 2007) are used to control the conversion in problematic cases.

of section 23 is tagged differently by (1) and (2)².

²In our experiments we used the scheme of (Yamada and Matsumoto, 2003), see Chapter 3. The significant effects of edge flipping were observed with the other two schemes as well.

The Neutral Edge Direction (NED) Measure

As shown in the previous chapters, the annotation of problematic edges can substantially affect performance. This was briefly discussed in (Klein and Manning, 2004), which used undirected evaluation as a measure which is less sensitive to alternative annotations. Undirected accuracy was commonly used since to assess the performance of unsupervised parsers (e.g., (Smith and Eisner, 2006; Headden III et al., 2008; Spitkovsky et al., 2010a)) but also of supervised ones (Wang et al., 2005; Wang et al., 2006). In this chapter we discuss why this measure is in fact not indifferent to *edge-flips* and propose a new measure, Neutral Edge Direction (NED).

5.1 Undirected Evaluation

The measure is defined as follows: traverse over the tokens and mark a correct attachment if the token's induced parent is either (1) its gold parent or (2) its gold child. The score is the ratio of correct attachments and the number of tokens.

We show that this measure does not ignore *edge-flips*. Consider Figure 5.1 that shows a dependency structure on the words w_1, w_2, w_3 before (Fig-

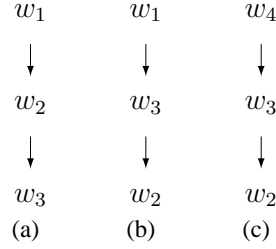


Figure 5.1: A dependency structure on the words w_1, w_2, w_3 before (Figure 5.1(a)) and after (Figure 5.1(b)) an *edge-flip* of $w_2 \rightarrow w_3$, and when the direction of the edge between w_2 and w_3 is switched and the new parent of w_3 is set to be some other word, w_4 (Figure 5.1(c)).

ure 5.1(a)) and after (Figure 5.1(b)) an *edge-flip* of $w_2 \rightarrow w_3$. Assume that 5.1(a) is the gold standard and that 5.1(b) is the induced parse. Consider w_2 . Its induced parent (w_3) is its gold child, and thus undirected evaluation does not consider it an error. On the other hand, w_3 is assigned w_2 's gold parent, w_1 . This is considered an error, since w_1 is neither w_3 's gold parent (as it is w_2), nor its gold child¹. Therefore, one of the two tokens involved in the *edge-flip* is penalized by the measure.

Recall the example “I want to eat” and the *edge-flip* of the edge “to” \leftarrow “eat” (Figure 3.2). As “to”'s parent in the induced graph (“want”) is neither its gold parent nor its gold child, the undirected evaluation measure marks it as an error. This is an example where an *edge-flip* in a problematic edge, which should not be considered an error, was in fact considered an error by undirected evaluation.

5.2 Neutral Edge Direction (NED)

The NED measure is a simple extension of the undirected evaluation measure². Unlike undirected evaluation, NED ignores all errors directly resulting

¹Otherwise, the gold parse would have contained a $w_1 \rightarrow w_2 \rightarrow w_3 \rightarrow w_1$ cycle.

²An implementation of NED is available at <http://www.cs.huji.ac.il/~roys02/software/ned.html>

from an *edge-flip*.

NED is defined as follows: traverse over the tokens and mark a correct attachment if the token’s induced parent is either (1) its gold parent (2) its gold child or (3) its gold grandparent. The score is the ratio of correct attachments and the number of tokens.

NED, by its definition, ignores *edge-flips*. Consider again Figure 5.1, where we assume that 5.1(a) is the gold standard and that 5.1(b) is the induced parse. Much like undirected evaluation, NED will mark the attachment of w_2 as correct, since its induced parent is its gold child. However, unlike undirected evaluation, w_3 ’s induced attachment will also be marked as correct, as its induced parent is its gold grandparent.

Now consider another induced parse in which the direction of the edge between w_2 and w_3 is switched and the w_3 ’s parent is set to be some other word, w_4 (Figure 5.1(c)). This should be marked as an error, even if the direction of the edge between w_2 and w_3 is controversial, since the structure $[w_2, w_3]$ is no longer a dependent of w_1 . It is indeed a NED error. Note that undirected evaluation gives the parses in Figure 5.1(b) and Figure 5.1(c) the same score, while if the structure $[w_2, w_3]$ is problematic, there is a major difference in their correctness.

5.3 Discussion

Problematic structures are ubiquitous, with more than 40% of the tokens in PTB WSJ appearing in at least one of them (see Chapter 3). Therefore, even a substantial difference in the attachment between two parsers is not necessarily indicative of a true quality difference. However, an attachment score difference that persists under NED is an indication of a true quality difference, since generally problematic structures are local (i.e., obtained by an *edge-flip*) and NED ignores such errors.

Reporting NED alone is insufficient, as obviously the edge direction does

matter in some cases. For example, in adjective–noun structures (e.g., “big house”), the correct edge direction is widely agreed upon (“big” \leftarrow “house”) (Kübler et al., 2009), and thus choosing the wrong direction should be considered an error. Therefore, we suggest evaluating using both NED and attachment score in order to get a full picture of the parser’s performance.

A possible criticism on NED is that it is only indifferent to alternative annotations in structures of size 2 (e.g., “to eat”) and does not necessarily handle larger problematic structures, such as coordinations (see Chapter 4). For example, Figure 5.2(a) and Figure 5.2(b) present two alternative annotations of the sentence “John and Mary”. Assume the parse in Figure 5.2(a) is the gold parse and that in Figure 5.2(b) is the induced parse. The word “Mary” is a NED error, since its induced parent (“and”) is neither its gold child nor its gold grandparent. Thus, NED does not accept all possible annotations of structures of size 3. On the other hand, using a method which accepts all possible annotations of structures of size 3 seems too permissive. A better solution may be to modify the gold standard annotation, so to explicitly annotate problematic structures as such. We defer this line of research to future work.

NED is therefore an evaluation measure which is indifferent to *edge-flips*, and is consequently less sensitive to alternative annotations. We now show that NED is indifferent to the differences between the structures originally learned by the algorithms mentioned in Chapter 3 and the gold standard annotation in all the problematic cases we consider.

Most of the modifications made are *edge-flips*, and are therefore ignored by NED. The exceptions are coordinations and prepositional phrases which are structures of size 3. In the former, the alternative annotations differ only in a single *edge-flip* (i.e., $CC \rightarrow NNP$), and are thus not NED errors. Regarding prepositional phrases, Figure 5.2(c) presents the gold standard of “in the house”, Figure 5.2(d) the parse induced by *km04* and *saj10a* and

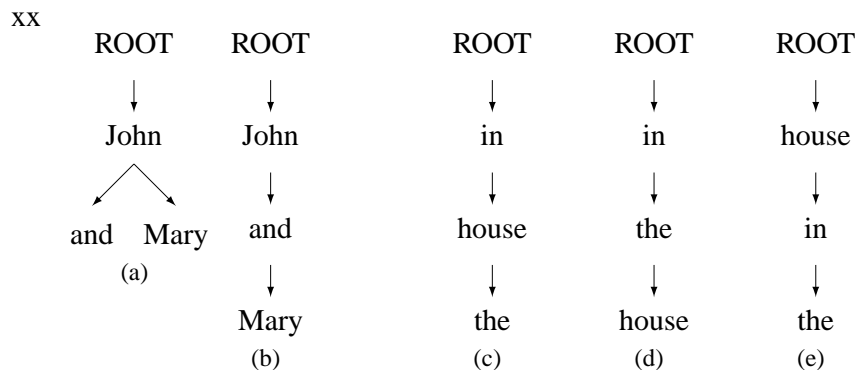


Figure 5.2: Alternative parses of “John and Mary” and “in the house”. Figure 5.2(a) follows (Collins, 1999), Figure 5.2(b) follows (Johansson and Nugues, 2007). Figure 5.2(c) follows (Collins, 1999; Yamada and Matsumoto, 2003). Figure 5.2(d) and Figure 5.2(e) show induced parses made by (*km04,saj10a*) and *cs09*, respectively.

Figure 5.2(e) the parse induced by *cs09*. As the reader can verify, both induced parses receive a perfect NED score.

In order to further demonstrate NED’s insensitivity to alternative annotations, we took two of the three common gold standard annotations (see Chapter 4) and evaluated them one against the other. We considered section 23 of WSJ following the scheme of (Yamada and Matsumoto, 2003) as the gold standard and of (Collins, 1999) as the evaluated set. Results show that the attachment score is only 85.6%, the undirected accuracy is improved to 90.3%, while the NED score is 95.3%. This shows that NED is significantly less sensitive to the differences between the different annotation schemes, compared to the other evaluation measures.

Chapter 6

Experimenting with NED

In this chapter we show that NED indeed reduces the performance difference between the original and the modified parameter sets, thus providing empirical evidence for its validity. For brevity, we present results only for the entire WSJ corpus. Results on WSJ10 are similar. The datasets and decoding algorithms are the same as those used in Chapter 3.

Table 6.1 shows the score differences between the parameter sets using attachment score, undirected evaluation and NED. A substantial difference persists under undirected evaluation: a gap of 7.7% in *cs09*, of 3.5% in *saj10a* and of 1.3% in *km04*.

The differences are further reduced using NED. This is consistent with our discussion in Chapter 5, and shows that undirected evaluation only ignores some of the errors inflicted by *edge-flips*.

For *cs09*, the difference is substantially reduced, but a 4.2% performance gap remains. For *km04* and *saj10a*, the original parameters outperform the new ones by 3.6% and 1% respectively.

We can see that even when ignoring *edge-flips*, some difference remains, albeit not necessarily in the favor of the modified models. This is because we did not directly perform *edge-flips*, but rather *parameter-flips*. The difference is thus a result of second-order effects stemming from the *parameter-*

flips. In the next chapter, we explain why the remaining difference is positive for some algorithms (*cs09*) and negative for others (*km04*, *saj10a*).

Algo.	<i>Mod. – Orig.</i>		
	Attach.	Undir.	NED
<i>km04</i>	9.3 (43.9–34.6)	1.3 (54.2–52.9)	–3.6 (63–66.6)
<i>cs09</i>	14.7 (54.6–39.9)	7.7 (56.9–49.2)	4.2 (66.8–62.6)
<i>saj10a</i>	12.7 (54.3–41.6)	3.5 (59.4–55.9)	–1 (66.8–67.8)

Table 6.1: Differences between the modified and original parameter sets when evaluated using attachment score (*Attach.*), undirected evaluation (*Undir.*), and NED.

For completeness, Table 6.2 shows a comparison of some of the current state-of-the-art algorithms, using attachment score, undirected evaluation and NED. The training and test sets are those used in Chapter 3. The table shows that the relative orderings of the algorithms under NED is different than under the other measures. This is an indication that NED provides a different perspective on algorithm quality¹.

¹Results may be different than the ones published in the original papers due to the different conversion procedures used in each work. See Chapter 4 for discussion.

Algo.	Att_{10}	Att_{∞}	Un_{10}	Un_{∞}	NED ₁₀	NED _{∞}
<i>bdk10</i>	66.1	49.6	70.1	56.0	75.5	61.8
<i>bc10</i>	67.2	53.6	73	61.7	81.6	70.2
<i>cs09</i>	61.5	42	66.9	50.4	81.5	62.9
<i>ggtp10</i>	57.1	45	62.5	53.2	80.4	65.1
<i>km04</i>	45.8	34.6	60.3	52.9	78.4	66.6
<i>saj10a</i>	54.7	41.6	66.5	55.9	78.9	67.8
<i>saj10c</i>	63.8	46.1	72.6	58.8	84.2	70.8
<i>saj10b*</i>	67.9	48.2	74.0	57.7	86.0	70.7

Table 6.2: A comparison of recent works, using Att (*attachment score*) Un (*undirected evaluation*) and NED, on sentences of length ≤ 10 (excluding punctuation) and on all sentences. The gold standard is obtained using the rules of (Yamada and Matsumoto, 2003). *bdk10*: (Berg-Kirkpatrick et al., 2010), *bc10*: (Blunsom and Cohn, 2010), *cs09*: (Cohen and Smith, 2009), *ggtp10*: (Gillenwater et al., 2010), *km04*: A replication of (Klein and Manning, 2004), *saj10a*: (Spitkovsky et al., 2010a), *saj10c*: (Spitkovsky et al., 2010c), *saj10b**: A lightly-supervised algorithm (Spitkovsky et al., 2010b).

Discussion

In this work we explored two ways of dealing with cases in which there is no clear theoretical justification to prefer one dependency structure over another. Our experiments suggest that it is crucial to deal with such structures if we would like to have a proper evaluation of unsupervised parsing algorithms against a gold standard.

The first way was to modify the parameters of the parsing algorithms so that in cases where such problematic decisions are to be made they follow the gold standard annotation. Indeed, this modification leads to a substantial improvement in the attachment score of the algorithms.

The second way was to change the evaluation. The NED measure we proposed does not punish for differences between gold and induced structures in the problematic cases. Indeed, in Chapter 6 (Table 6.1) we show that the differences between the original and modified models are much smaller when evaluating with NED compared to when evaluating with the traditional attachment score.

As Table 6.1 reveals, however, even when evaluating with NED, there is still some difference between the original and the modified model, for each of the algorithms we consider. Moreover, for two of the algorithms (*km04* and *saj10a*) NED prefers the original model while for one (*cs09*) it prefers

the modified version. In this chapter we explain these patterns and show that they are both consistent and predictable.

Our hypothesis, for which we provide empirical justification, is that in cases where there is no theoretically preferred annotation, NED prefers the structures that are more learnable by DMV. That is, NED gives higher scores to the annotations that better fit the assumptions and modeling decisions of DMV, the model that lies in the basis of the parsing algorithms.

To support our hypothesis we perform an experiment requiring two preparatory steps for each algorithm. First, we construct a supervised version of the algorithm. This supervised version consists of the same statistical model as the original unsupervised algorithm, but the parameters are estimated to maximize the likelihood of a *syntactically annotated* training corpus, rather than of a plain text corpus.

Second, we construct two corpora for the algorithm, both consist of the same text and differ only in their syntactic annotation. The first is annotated with the gold standard annotation. The second is similarly annotated except in the linguistically problematic structures. We replace these structures with the ones that would have been created with the unsupervised version of the algorithm (see Table 3.1 for the relevant structures for each algorithm)¹. Each corpus is divided into a training and a test set.

We then train the supervised version of the algorithms on each of the training sets. We parse the test data twice, once with each of the resulting models. We evaluate both parsed corpora against the corpus annotation from which they originated.

¹In cases the structures are comprised of a single edge, the second corpus is obtained from the gold standard by an *edge-flip*. The only exceptions are the cases of the prepositional phrases. Their gold standard and the learned structures for each of the algorithms are shown in Figure 5.2. In this case, the second corpus is obtained from the gold standard by replacing each prepositional phrase in the gold standard with the corresponding learned structure.

The training set of each corpus consists of sections 2–21 of WSJ20 (i.e., WSJ sentences of length ≤ 20 , excluding punctuation)² and the test set is section 23 of WSJ $^\infty$. Evaluation is performed using both NED and attachment score. The patterns we observed are very similar for both. For brevity, we report only attachment score results.

	<i>km04</i>		<i>cs09</i>		<i>saj10a</i>	
	<i>Orig.</i>	<i>Gold</i>	<i>Orig.</i>	<i>Gold</i>	<i>Orig.</i>	<i>Gold</i>
NED, Unsup.	66.6	63	62.6	66.8	67.8	66.8
Sup.	71.3	69.9	63.3	69.9	71.8	69.9

Table 7.1: The first line shows the NED results from Chapter 6, when using the original parameters (*Orig.* columns) and the modified parameters (*Gold* columns). The second line shows the results of the supervised versions of the algorithms using the corpus which agrees with the unsupervised model in the problematic cases (*Orig.*) and the gold standard (*Gold*).

The results of our experiment are presented in Table 7.1 along with a comparison to the NED scores from Chapter 6. The table clearly demonstrates that a set of parameters (original or modified) is preferred by NED in the unsupervised experiments reported in Chapter 6 (top line) if and only if the structures produced by this set are better learned by the supervised version of the algorithm (bottom line).

This observation supports our hypothesis that in cases where there is no theoretical preference for one structure over the other, NED (unlike the other measures) prefers the structures that are more consistent with the modeling assumptions lying in the basis of the algorithm. We consider this to be

²In using WSJ20, we follow (Spitkovsky et al., 2010a), which showed that training the DMV on sentences of bounded length yields a higher score than using the entire corpus. We use it as we aim to use an optimal setting.

a desired property of a measure since a more consistent model should be preferred where no theoretical preference exists.

Conclusion

In this work we showed that the standard evaluation of unsupervised dependency parsers is highly sensitive to problematic annotations. We modified a small set of parameters that controls the annotation in such problematic cases in three leading parsers. This resulted in a major performance boost, which is unindicative of a true difference in quality.

We also presented another problem with the evaluation of this task – the fact that three substantially different gold standard are currently used in evaluating it.

We presented *Neutral Edge Direction* (NED), a measure that is less sensitive to the annotation of local structures. As the problematic structures are generally local, NED is less sensitive to their alternative annotations. In the future, we suggest reporting NED along with the current measures.

References

- Taylor Berg-Kirkpatrick and Dan Klein. 2010. Phylogenetic grammar induction. In *Proc. of ACL*.
- Taylor Berg-Kirkpatrick, Alexandre Bouchard-Côté, John DeNero, and Dan Klein. 2010. Painless unsupervised learning with features. In *Proc. of NAACL*.
- Phil Blunsom and Trevor Cohn. 2010. Unsupervised induction of tree substitution grammars for dependency parsing. In *Proc. of EMNLP*.
- Cristina Bosco and Vincenzo Lombardo. 2004. Dependency and relational structure in treebank annotation. In *Proc. of the Workshop on Recent Advances in Dependency Grammar at COLING'04*.
- Alexander Clark. 2001. *Unsupervised language acquisition: theory and practice*. Ph.D. thesis, University of Sussex.
- Shay B. Cohen and Noah A. Smith. 2009. Shared logistic normal distributions for soft parameter tying. In *Proc. of HLT-NAACL*.
- Shay B. Cohen, Kevin Gimpel, and Noah A. Smith. 2008. Logistic normal priors for unsupervised probabilistic grammar induction. In *Proc. of NIPS*.
- Michael J. Collins. 1999. *Head-driven statistical models for natural language parsing*. Ph.D. thesis, University of Pennsylvania, Philadelphia.
- Hal Daumé III. 2009. Unsupervised search-based structured prediction. In *Proc. of ICML*.
- Mark Dredze, John Blitzer, Partha Pratim Talukdar, Kuzman Ganchev, João V. Graça, and Fernando Pereira. 2007. Frustratingly hard domain adaptation for dependency parsing. In *Proc. of the CoNLL 2007 Shared Task. EMNLP-CoNLL*.
- Gregory Druck, Gideon Mann, and Andrew McCallum. 2009. Semi-supervised learning of dependency parsers using generalized expectation criteria. In *Proc. of ACL*.
- Jennifer Gillenwater, Kuzman Ganchev, João V. Graça, Ben Taskar, and Fernando Pereira. 2010. Sparsity in dependency grammar induction. In *Proc. of ACL*.

- William P. Headden III, David McClosky, and Eugene Charniak. 2008. Evaluating unsupervised part-of-speech tagging for grammar induction. In *Proc. of COLING*.
- William P. Headden III, Mark Johnson, and David McClosky. 2009. Improving unsupervised dependency parsing with richer contexts and smoothing. In *Proc. of HLT-NAACL*.
- Richard Johansson and Pierre Nugues. 2007. Extended constituent-to-dependency conversion for english. In *Proc. of NODALIDA*.
- Dan Klein and Christopher Manning. 2004. Corpus-based induction of syntactic structure: Models of dependency and constituency. In *Proc. of ACL*.
- Dan Klein. 2005. *The unsupervised learning of natural language structure*. Ph.D. thesis, Stanford University.
- Sandra Kübler, R. McDonald, and Joakim Nivre. 2009. *Dependency Parsing*. Morgan And Claypool Publishers.
- Sandra Kübler. 2005. How do treebank annotation schemes influence parsing results? or how not to compare apples and oranges. In *Proc. of RANLP*.
- Mitchell Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a large annotated corpus of english: The penn treebank. *Computational Linguistics*, 19:313–330.
- Igor Mel’čuk. 1988. *Dependency Syntax: Theory and Practice*. State University of New York Press.
- Tahira Naseem, Harr Chen, Regina Barzilay, and Mark Johnson. 2010. Using universal linguistic knowledge to guide grammar induction. In *Proc. of EMNLP*.
- Jens Nilsson, Joakim Nivre, and Johan Hall. 2006. Graph transformations in data-driven dependency parsing. In *Proc. of ACL*.
- Joakim Nivre, Johan Hall, and Jens Nilsson. 2006. Maltparser: A data-driven parser-generator for dependency parsing. In *Proc. of LREC-2006*.
- Joakim Nivre, Johan Hall, Sandra Kübler, Ryan McDonald, Jens Nilsson, Sebastian Riedel, and Deniz Yuret. 2007. The conll 2007 shared task on dependency parsing. In *Proc. of the CoNLL Shared Task, EMNLP-CoNLL, 2007*.

- Joakim Nivre. 2006. *Inductive Dependency Parsing*. Springer.
- Owen Rambow, Cassandre Creswell, Rachel Szekely, Harriet Tauber, and Marilyn Walker. 2002. A dependency treebank for english. In *Proc. of LREC*.
- Roy Schwartz, Omri Abend, Roi Reichart, and Ari Rappoport. 2011. Neutralizing linguistically problematic annotations unsupervised dependency parsing evaluation. In *Proc. of ACL-HLT*.
- Noah A. Smith and Jason Eisner. 2005. Guiding unsupervised grammar induction using contrastive estimation. In *Proc. of IJCAI*.
- Noah A. Smith and Jason Eisner. 2006. Annealing structural bias in multilingual weighted grammar induction. In *Proc. of ACL*.
- Valentin I. Spitzkovsky, Hiyen Alshawi, and Daniel Jurafsky. 2010a. From baby steps to leapfrog: How “less is more” in unsupervised dependency parsing. In *Proc. of NAACL-HLT*.
- Valentin I. Spitzkovsky, Hiyen Alshawi, and Daniel Jurafsky. 2010b. Profiting from mark-up: Hyper-text annotations for guided parsing. In *Proc. of ACL*.
- Valentin I. Spitzkovsky, Hiyen Alshawi, Daniel Jurafsky, and Christopher D. Manning. 2010c. Viterbi training improves unsupervised dependency parsing. In *Proc. of CoNLL*.
- Qin Iris Wang, Dale Schuurmans, and Dekang Lin. 2005. Strictly lexical dependency parsing. In *Proc. of IWPT*.
- Qin Iris Wang, Colin Cherry, Dan Lizotte, and Dale Schuurmans. 2006. Improved large margin dependency parsing via local constraints and laplacian regularization. In *Proc. of CoNLL*.
- Hiroyasu Yamada and Yuji Matsumoto. 2003. Statistical dependency analysis with support vector machines. In *Proc. of the International Workshop on Parsing Technologies*.