

Parser und Tagger

7 Parser und Tagger

- TreeTagger

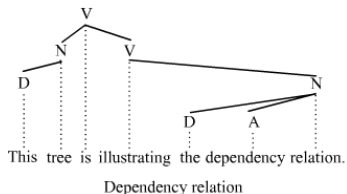
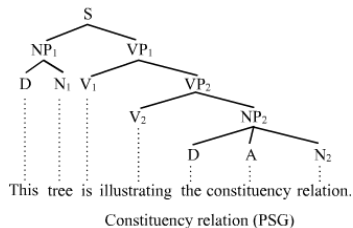
- MaltParser

- Stanford Parser

MaltParser

- Deterministischer Dependenz-Parser, implementiert in Java (seit Version 1.0)
- Parsing-Modell wird aus einer Treebank induziert
- Neue Daten werden mit dem induzierten Modell geparkt
- Webseite: <http://maltparser.org/download.html>
- Auf unseren Servern unter
/resources/processors/parser/maltparser-1.8/
- Modelle für Deutsch, Englisch, Französisch und Schwedisch
- Eingabe und Ausgabe im *CoNLL Format*

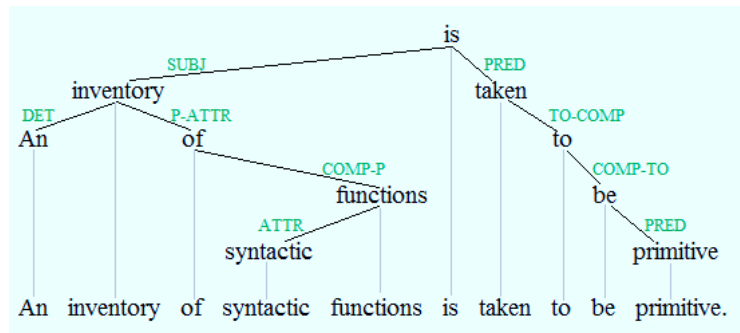
Dependenz- vs. Konstituenten-Parser



Quelle:

https://en.wikipedia.org/wiki/Phrase_structure_grammar

Dependenzbäume mit Kantenbeschriftungen



Quelle:

<https://commons.wikimedia.org/w/index.php?curid=21976793>

CoNLL Format

- Standardisiertes Format für die Eingaben/Ausgaben bei CoNLL Shared Task 2007
- Text-Dateien mit der Erweiterung `.conll`
- Eine Zeile besteht aus 10 durch `TAB` getrennten Feldern
- Ein Satz besteht aus einem oder mehreren Tokens, je eines auf einer Zeile
- Eine Datei kann mehrere Sätze enthalten (getrennt durch eine Leerzeile)
- Leerzeichen in den Feldern sind nicht erlaubt!

<http://universaldependencies.org/format.html>

CoNLL Format - Felder

- 1** ID - Der Zähler für Tokens (beginnt mit 1 in jedem Satz)
- 2** FORM - Wortform
- 3** LEMMA - Lemma des Wortes
- 4** CPOSTAG - POS-Tag (sprachabhängig)
- 5** POSTAG - Detailliertes POS-Tag (sprachabhängig)
- 6** FEATS - Syntaktische und/oder morphologische Merkmale (sprachabhängig)
- 7** HEAD - Head-Token für das aktuelle Token (ID oder 0)
- 8** DEPREL - Dependenzrelation mit Head-Token (von den Treebank-Annotationen abhängig)

CoNLL Format - Bemerkungen

- Die folgenden Felder *müssen* immer ausgefüllt werden:
 - Bei den Trainingsdaten:
ID, FORM, CPOSTAG, POSTAG, HEAD, DEPREL
 - Bei den Testdaten:
ID, FORM, CPOSTAG, POSTAG
- Andere Felder können mit dem Platzhalter „_“ besetzt werden
- Falls es keine detaillierten POS-Tags für eine bestimmte Sprache gibt, sind die Felder CPOSTAG und POSTAG gleich
- Mehrere Merkmale werden durch das Symbol „|“ getrennt
- Abhängigkeitsrelation für den Startknoten ist ROOT (mit dem HEAD-Wert 0)

CoNLL Format - Beispiel

Satz aus einem Trainingsset

1	Trouble	_	NN	NN	-	2	SBJ	-	-
2	is	-	VBZ	VBZ	-	0	ROOT	-	-
3	,	-	,	,	-	2	P	-	-
4	she	-	PRP	PRP	-	5	SBJ	-	-
5	has	-	VBZ	VBZ	-	2	PRD	-	-
6	lost	-	VC	VC	-	5	VC	-	-
7	it	-	PRP	PRP	-	6	OBJ	-	-
8	just	-	RB	RB	-	6	MNR	-	-
9	as	-	RB	RB	-	8	AMOD	-	-
10	quickly	-	RB	RB	-	8	AMOD	-	-
11	.	-	.	.	-	2	P	-	-

CoNLL Format - Beispiel

Satz aus einem Testset

1	Pierre	-	NNP	NNP	-
2	Vinken	-	NNP	NNP	-
3	,	-	,	,	-
4	61	-	CD	CD	-
5	years	-	NNS	NNS	-
6	old	-	JJ	JJ	-
7	,	-	,	,	-
8	will	-	MD	MD	-
9	join	-	VB	VB	-
10	the	-	DT	DT	-
11	board	-	NN	NN	-
12	as	-	IN	IN	-
13	a	-	DT	DT	-
14	nonexecutive	-	JJ	JJ	-
15	director	-	NN	NN	-
16	Nov.	-	NNP	NNP	-
17	29	-	CD	CD	-
18	.	-	.	.	-

CoNLL Format - Beispiel

Andere Sprache

1	Grundavdraget	-	N	NN	DD SS	2	SS	-	-
2	blir	-	V	BV	PS	0	ROOT	-	-
3	2500	-	R	RO	-	4	DT	-	-
4	kr	-	N	NN	-	2	SP	-	-
5	(-	I	IR	-	0	ROOT	-	-
6	=	-	I	IG	-	15	DT	-	-
7	4500	-	R	RO	-	8	DT	-	-
8	minskat	-	P	TP	PA	0	ROOT	-	-
9	med	-	PR	PR	-	14	AA	-	-
10	1/5	-	N	NN	-	9	PA	-	-
11	av	-	PR	PR	-	10	ET	-	-
12	10000	-	R	RO	-	13	DT	-	-
13	kr	-	N	NN	-	11	PA	-	-
14)	-	I	IR	-	0	ROOT	-	-
15	.	-	P	IP	-	2	IP	-	-

Erster Start des MaltParsers

- Starten des MaltParsers ohne Argumente:

```
:~$ java -jar /path/to/maltparser-1.8/maltparser-1.8.jar
```

- Gibt die Hauptoptionen für den MaltParser aus

- Starten der Hilfe für den MaltParser:

```
:~$ java -jar /path/to/maltparser-1.8/maltparser-1.8.jar -h
```

- Gibt alle Optionen für den MaltParser aus

Den MaltParser trainieren

```
:~$ java -jar /path/to/maltparser-1.8/maltparser-1.8.jar -c model \  
-i training-file.conll -m learn
```

- **model** : Name des Modells
Achtung: nur den Namen ohne Pfad angeben, die Erweiterung (.mco) wird *automatisch* angefügt!
- **training-file.conll** : Datei mit den Trainingsdaten im CoNLL-Format
- **-m learn** : der Parser wird im Lern-Modus gestartet

Das erstellte Modell wird (hier) in einer Datei namens **model.mco** im aktuellen Verzeichnis gespeichert.

Vortrainierte Modelle für den MaltParser

- Für diejenigen, die nicht mit unterschiedlichen Algorithmen und Modellen experimentieren möchten, sondern nur einen Dependenzparser brauchen
- Aktuell sind drei fertige Modelle vorhanden: `engmalt`, `fremalt`, `swemalt` und `espmalt` (alle jeweils `*.linear` und `*.poly`)
- Download: <http://maltparser.org/mco/mco.html>
- Basieren auf Penn Treebank, French Treebank und Swedisch Treebank

Parzen von Daten im MaltParser

```
:~$ java -jar /path/to/maltparser-1.8/maltparser-1.8.jar -c model\  
-i test-file.conll -o result-file.conll -m parse
```

- `model` : Name des Modells
Achtung: nur den Namen ohne die Erweiterung (`.mco`) angeben
- `test-file.conll` : Datei mit den Testdaten
- `result-file.conll` : Datei mit den Ergebnissen des Parsers
- `-m parse` : der Parser wird im Parsing-Modus gestartet

MaltParser – Weitere Informationen

- Userguide: <http://maltparser.org/userguide.html>
- (Technische) Informationen zur Optimierung:
[http:
//www.maltparser.org/guides/opt/quick-opt.pdf](http://www.maltparser.org/guides/opt/quick-opt.pdf)

Übung 9