Language Processing with Perl and Prolog Chapter 11: Syntactic Formalisms

Pierre Nugues

Lund University Pierre.Nugues@cs.lth.se http://cs.lth.se/pierre_nugues/



Pierre Nugues

Syntax has been the core of linguistics in the US and elsewhere for many years

Noam Chomsky, professor at the MIT, has had an overwhelming influence, sometimes misleading

Syntactic structures (1957) has been a cult book for the past generation of linguists

Syntax can be divided into two parts:

- Formalism How to represent syntax
- Parsing How to get the representation of a sentence

Syntactic Formalisms

The two most accepted formalisms use a tree representation:

- One is based on the idea of constituents
- Another is based on dependencies between words. Trees have originally been called stemmas

They are generally associated respectively to Chomsky and Tesnière. Later, constituent grammars evolved into unification grammars

Constituency

Constituency can be expressed by context-free grammars. They are defined by

- A set of designated start symbols, Σ, covering the sentences to parse. This set can be reduced to a single symbol, such as sentence, or divided into more symbols: declarative_sentence, interrogative_sentence.
- A set of nonterminal symbols enabling the representation of the syntactic categories. This set includes the sentence and phrase categories.
- A set of terminal symbols representing the vocabulary: words of the lexicon, possibly morphemes.
- A set of rules, *F*, where the left-hand-side symbol of the rule is rewritten in the sequence of symbols of the right-hand side.

These grammars can be mapped to DCG rules as for The boy hit the ball

```
sentence --> np, vp.
np --> t, n.
vp -- verb, np.
t --> [the].
n --> [man] ; [ball] ; etc.
verb --> [hit] ; [took] ; etc.
```

Generation of sentences is one of the purposes of grammar according to Chomsky

Perl and Pro

Chomsky Normal Form

In some parsing algorithms, it is necessary to have rules in the Chomsky normal form (CNF) with two right-hand-side symbols Non-CNF rules:

lhs --> rhs1, rhs2, rhs3.

can be converted into a CNF equivalent:

lhs --> rhs1, lhs_aux. lhs_aux --> rhs2, rhs3.

> Language Processing with Peri and Prolog Sector

Transformations

Rearrangement of sentences according to some syntactic relations: active/passive, declarative/interrogative, etc. Transformations use rules – transformational rules or T rules – *The boy will hit the ball/the ball will be (en) hit by the boy*

```
T1: np1, aux, v, np2 --->
    np2, aux, [be], [en], v, [by], np1
```



Transformations



Syntactic Categories (Penn Treebank)

	Categories	Description
1.	ADJP	Adjective phrase
2.	ADVP	Adverb phrase
3.	NP	Noun phrase
4.	PP	Prepositional phrase
5.	S	Simple declarative clause
6.	SBAR	Clause introduced by subordinating conjunction or 0
7.	SBARQ	Direct question introduced by wh-word or phrase
8.	SINV	Declarative sentence with subject-aux inversion
9.	SQ	Subconstituent of SBARQ excluding wh-word or phrase
10.	VP	Verb phrase
11.	WHADVP	wh-adverb phrase
12.	WHNP	wh-noun phrase
13.	WHPP	wh-prepositional phrase Processing with Perland Prology
14.	Х	Constituent of unknown or uncertain category

A Hand-Parsed Sentence using the Penn Treebank Annotation

Battle-tested industrial managers here always buck up nervous newcomers with the tale of the first of their countrymen to visit Mexico, a boatload of samurai warriors blown ashore 375 years ago.

```
( (S
   (NP Battle-tested industrial managers
    here)
   always
   (VP buck
    up
        (NP nervous newcomers)
        (PP with
            (NP the tale
            (PP of
```

A Hand-Parsed Sentence using the Penn Treebank Annotation

```
(NP
    (NP the
        (ADJP first
               (PP of
                   (NP their countrymen)))
        (S (NP *)
           to
           (VP visit
                (NP Mexico))))
    ,
    (NP (NP a boatload
            (PP of
                 (NP (NP samurai warriors)
                     (VP-1 blown
                         ashore
                     (ADVP (NP 375 years)
                           ago)))))
        (VP-1 *nseudo-attach*り))的
```

Language Processing with Perl and Prolog

rocessina wi

Unification-based Grammars

Grammatical features such as case modify the word morphology

Cases	Noun groups
Nominative	der kleine Ober
Genitive	des kleinen Obers
Dative	dem kleinen Ober
Accusative	den kleinen Ober

```
The rule
```

```
np --> det, adj, n.
outputs ungrammatical phrases as:
?-np(L, []).
[der, kleinen, Ober]; %wrong
[der, kleinen, Obers]; %wrong
[dem, kleine, Obers] %wrong
```



. . .

Representing Features

```
A possible solution is to use arguments: np(case:C) where the C value is a member of list [nom, gen, dat, acc]
```

```
np(gend:G, num:N, case:C, pers:P, det:D)
np(gend:G, num:N, case:C, pers:P, det:D) -->
det(gend:G, num:N, case:C, pers:P, det:D),
adj(gend:G, num:N, case:C, pers:P, det:D),
n(gend:G, num:N, case:C, pers:P).
```



A Small Fragment of German

det(gend:masc, num:sg, case:nom, pers:3, det:def) --> [der]. det(gend:masc, num:sg, case:gen, pers:3, det:def) --> [des]. det(gend:masc, num:sg, case:dat, pers:3, det:def) --> [dem]. det(gend:masc, num:sg, case:acc, pers:3, det:def) --> [den]. adj(gend:masc, num:sg, case:nom, pers:3, det:def) --> [kleine] adj(gend:masc, num:sg, case:gen, pers:3, det:def) --> [kleinen]. adj(gend:masc, num:sg, case:dat, pers:3, det:def) --> [kleinen]. adj(gend:masc, num:sg, case:acc, pers:3, det:def) --> [kleinen]. n(gend:masc, num:sg, case:nom, pers:3) --> ['Ober']. n(gend:masc, num:sg, case:gen, pers:3) --> ['Obers']. n(gend:masc, num:sg, case:dat, pers:3) --> ['Ober']. Processing with Perl and Prolog n(gend:masc, num:sg, case:acc, pers:3) --> ['Ober']. ma a

A Unification-based Formalism

Unification-based grammars use a notation close to that of DCGs





Some Rules



Feature Structures are Graphs

Structures can be embedded



Feature Structures are Graphs



Unification-based Formalism

The feature notation is based on the name, not on the position

are equivalent

Unification is a generalization of Prolog unification See the course book for the implementation



Dependency Grammars

Dependency grammars (DG) describe the structure in term of links



Each word has a head or "régissant" except the root of the sentence. A head has one or more modifiers or dependents: *Cat* is the head of *big* and *the*; *big* is the head of *very*. DG can be more versatile with a flexible word order language like German Russian, or Latin.

A Sentence Tree – Stemma





Language Technology

Properties of Dependency Graphs



separated by direct or indirect dependents of Dep or Head



Language Technology

Chapter 11: Syntactic Formalisms

Nonprojective Graphs (McDonald and Pereira)







Tesnière makes a distinction between essential and circumstantial complements Essential – or core – complements are for instance subject and objects. Circumstantial – or noncore – complements are the adjuncts Valence corresponds to the verb saturation of its essential complements



Valence Examples

Val.	Examples	Frames
0	it's raining	raining []
1	he's sleeping	sleeping [subject : he]
2	she read this book	read [subject : she object : book]
3	Elke gave a book to Wolfgang	gave subject : Elke object : book iobject : Wolfgang
4	<i>I moved the car from here to the street</i>	moved subject : I object : car source : here destination : streat
D'auro		Yetan Prod Pertan Prod · · · · · · · · · · · · · · · · · · ·

Subcategorization Frames

Valence is a model of verb construction. It can be extended to more specific patterns as in the *Oxford Advanced Learner's Dictionary* (OALD).

Verb	Complement structure	Example
slept	None (Intransitive)	l slept
bring	NP	The waiter brought the meal
bring	NP + to + NP	The waiter brought the meal to the patron
depend	on + NP	It depends on the waiter
wait	for $+ NP + to + VP$	<i>I am waiting for the waiter to bring the meal</i>
keep	VP(ing)	He kept working
know	that + S	The waiter knows that the Datron loves fish

Subcategorization Frames in German

Verb	Complement structure	Example
schlafen	None (Intransitive)	Ich habe geschlafen
bringen	NP(Accusative)	Der Ober hat eine Speise ge- bracht
bringen	NP(Dative) + NP(Accusative)	Der Ober hat dem Kunde eine Speise gebracht
abhängen	von + NP(Dative)	Es hängt vom Ober ab
warten	auf + S	Er wartete auf dem Ober, die Speise zu bringen
fortsetzen	NP	Er hat die Arbeit fortgesetzt
wissen	NP(Final verb)	Der Ober weiß, das der Kunde
		Fisch liebt
		Language Processing with Peri and Prolog

Dependencies and Grammatical Functions

The dependency structure generally reflects the traditional syntactic representation

The links can be annotated with grammatical function labels.

In a simple sentence, it corresponds to the subject and the object



Probably a more natural description to tie syntax to semantics

Dependencies and Functions (II)

Adjuncts form another class of functions that modify the verb They include prepositional phrases whose head is set arbitrarily to the front preposition

Adjuncts include adverbs that modify a verb



Dependency Parse Tree



Word	Word	Direction	Head	Head	Function
pos.				position	
1	Bring	*		Root	Main verb
2	the	>	meal	3	Determiner
3	meal	<	Bring	1	Object
4	to	<	Bring	1	Location
5	the	>	table	6	Determiner
6	table	<	to	4	Prepositional com

Representing Dependencies

 $D = \{ < \mathsf{Head}(1), \mathsf{Rel}(1) >, < \mathsf{Head}(2), \mathsf{Rel}(2) >, ..., < \mathsf{Head}(n), \mathsf{Rel}(n) > \},$

The representation of *Bring the meal to the table*:



Annotation: MALT XML

```
<sentence id="24">
<word id="1" form="Dessutom" postag="ab" head="2"</pre>
  deprel="ADV"/>
<word id="2" form="höjs" postag="vb.prs.sfo" head="0"</pre>
  deprel=""/>
<word id="3" form="åldergränsen" postag="nn.utr.sin.def.nom"</pre>
 head="2" deprel="SUB"/>
<word id="4" form="till" postag="pp" head="2" deprel="ADV"/>
<word id="5" form="18" postag="rg.nom" head="6" deprel="DET"/>
<word id="6" form="år" postag="nn.neu.plu.ind.nom" head="4"</pre>
  deprel="PR"/>
<word id="7" form="." postag="mad" head="2" deprel="IP"/>
</sentence>
```

TMALT XML is an extended annotation

Annotation: CoNLL

The CoNLL shared tasks organize evaluations of machine-learning systems for natural language processing.

They define formats to share data between participants.

1	Dessutom	_	AB	AB	_	2	+A	_	_
2	höjs	_	VV	VV	_	0	ROOT	_	_
3	åldergränsen	_	NN	NN	_	2	SS	_	_
4	till		PR	PR		2	OA		
5	18		RO	RO		6	DT		
6	år		NN	NN		4	PA		
7		_	IP	IP	_	2	IP		_



Annotation: CoNLL

#	Name	Description
1	ID	Token index, starting at 1 for each sentence.
2	FORM	Word form or punctuation.
3	LEMMA	Lemma or stem.
4	CPOSTAG	Part-of-speech tag.
5	POSTAG	Fine-grained part-of-speech tag.
6	FEATS	Unordered set of morphological features separated by a vertical
		bar ().
7	HEAD	Head of the current token, which is either a value of ID or zero
		(0) if this is the root.
8	DEPREL	Dependency relation to the HEAD.
9	PHEAD	Projective head of current token, which is either a value of ID or
		zero (0). The dependency structure resulting from the PHEAD
		column is guaranteed to be projective, when available in the
		Corpus.
10	PDEPREL	Dependency relation to the PHEAD.

Visualizing Dependencies

Using What's Wrong With My NLP (https://code.google.com/p/whatswrong/):



Language Technology

Chapter 11: Syntactic Formalisms

Function Annotation Tagset (Järvinen and Tapanainen 1997)

Name	Description	Example				
		Main	n function	S		
main	Main element		He doesn	n't know who	ether to se	end a gift
qtag	Question tag		Let's pla	y another ga	ame, <u>shal</u>	I we?
		Intran	uclear lin	ks		
v-ch	Verb chain		lt may h	nave been b	oeing exa	mined
pcomp	Prepositional	comple-	They	played	the	game
	ment		in a diffe	erent way		
phr	Verb particle		He asked	l me who wo	uld look a	after the
			baby			Parenti Lisagan
						Language Processing with Perl and Prolog

Function Annotation Tagset (Järvinen and Tapanainen 1997)

Verb complementation				
subj	Subject			
obj	Object	I gave him my address		
comp	Subject complement.	It has become marginal		
dat	Indirect object	Pauline gave it <u>to Tom</u>		
ос	Object complement	His friends call him Ted		
copred	Copredicative	We took a swim naked		
voc	Vocative	Play it again, Sam		
	Determi	native functions		
qn	Quantifier	I want more money		
det	Determiner	Other members will join		
neg	Negator	It is not coffee that I like, but te		
		Sound State		

Language Technology

Chapter 11: Syntactic Formalisms

Function Annotation Tagset (Järvinen and Tapanainen 1997)

Modifiers				
attr	Attributive nominal	Knowing no French, I couldn't express		
		my thanks		
mod	Other postmodifiers	The baby, <u>Frances Bean,</u> was		
		The people <u>on the bus</u> were singing		
ad	Attributive adverbial	She is more popular		
Junctives				
сс	Coordination	<u>Two or</u> more cars		



Dependency vs. Constituency

- Constituency (most textbooks) is a declining formalism
- It cannot properly handle many languages: Swedish, Russian, Czech, Arabic, etc.
- Dependency parsing can handle all these languages as well as English, German, French, etc.
- Dependency parsing has improved considerably over the last 4 years: see CoNLL 2006 and 2007.
- CoNLL 2008 and 2009 extend it to semantic parsing
- However, constituency and dependency are (weakly) compatible provided that we restrict us to projective dependency graphs



From Constituency to Dependency

It is possible to convert constituent trees into dependency graphs We need to identify a headword in all the PS rules, here with a star:

```
s --> np, vp*.
vp --> verb*, np.
np --> det, noun*.
```

Parsers by Magerman and Collins used this to convert the Penn Treebank constituent annotation for their dependency parsers When projective, dependency structures are loosely compatible with constituent grammars.



From Constituency to Dependency (II)

A constituent tree with head-marked rules:



The resulting dependency graph:

