

Automatic annotations : the syntactic information

Stéphane Rauzy

Laboratoire Parole et Langage
CNRS & Université de Provence

`stephane.rauzy@lpl-aix.fr`

In input, a transcription of spontaneous speech in interaction (e.g. CID copora)

In output, the associated syntactic information :

- Part-Of-Speech tagging, e.g. **Det**, **Noun**, ...
- Groups and syntactic tree structure, e.g. **NP**, **VP**, ...
- Functional relations, e.g. **SUB**, **OBJ**, ...

Strategy :

- Apply the existing tools developed for written textual input
- Adapt them for the treatment of speech transcription
- Investigate and manually correct the annotations in output
- Propose a new model for treating spontaneous speech

The LPL resources and tools for the syntactic treatment of french textual entries (see for example Rauzy & Blache 2009) :

- Lexicon and tokenizer
- Stochastic tagger
- Stochastic chunker and deep parser

Based on Patterns model approach (stochastic model), the grammar is learned/extracted from an annotated corpus. The quality of the tools depends on :

- The size and the coverage of the training corpus
- The quality of the Gold Standard annotations

LPL tagger, chunker, and parser for written french

Stochastic models (e.g. Patterns model) allows to associate a probability to any sequence of categories (e.g. POS tags), and thus describes the regularities found in their distribution.

Training stage :

- Learn the parameters of the model on an annotated corpus, e.g. 853 occurrences of **Noun Verb Det Adj**, followed by :

pattern context	category	occurrences	proba
Noun Verb Det Adj	Pct	12	0.015
	Coord	34	0.045
	Noun	807	0.94

Automatic annotation stage :

- Apply the model on raw data, the solution is the one maximizing the overall probability of the sequence.

Part-Of-Speech tagging

A lexicon allows to associate to each form (Part-of-Speech) of the sentence its corresponding lexical tags distribution, e.g.

form	lemma	sampa	tag	frequency
est	être	E	Aux	1671
est	être	E	Verb	21395
est	est	Est	Noun	422

The tagger operates the desambiguation process by choosing the most probable solution, e.g.

Sentence : La valise est dans le coffre .
Propositions : Det Noun Noun Noun Det Noun Pct
 Noun Verb Prep Pro Verb
 Pro Aux

Maximal probability : Det Noun Verb Prep Det Noun Pct

Part-Of-Speech tagging

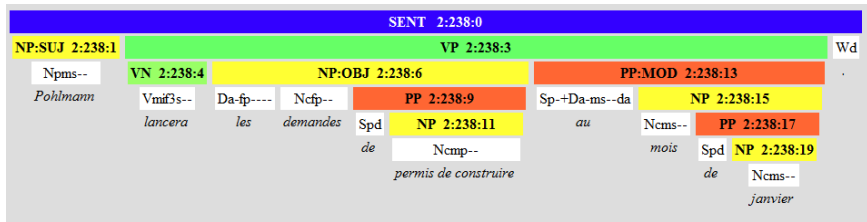
Form	Solution	Score	A-score	Propositions
<i>La</i>	Da-fs--d-	-6.6362705	A	Pp3fsj- Da-fs--d- Ncm---
<i>définition</i>	Ncfs--	-6.4358025	A	Ncfs--
<i>connaît</i>	Vmip3s--	0.74682426	C	Vmip3s--
<i>des</i>	Spd+Da-mp--id	-4.5408936	A	Da-mp--i- Sp-+Da-fp--dd Spd+Da-fp--id Spd+Da-mp--dd Spd+Da-mp--id
<i>nuances</i>	Ncfp--	-5.697094	A	Ncfp-- Vmip2s-- Vmsp2s--
<i>importantes</i>	Afpfp-	-1.0718307	A	Afpfp-
<i>selon</i>	Sp-	-3.0093784	A	Sp-
<i>ces</i>	Dd-mp----	-0.7771969	A	Dd-fp---- Dd-mp----
<i>différents</i>	Afpmp-	-0.046440125	A	Afpmp- Ai-mp- Di-mp----
<i>domaines</i>	Ncmp--	-5.917076	A	Ncmp--
.	Wd	-1.1673737	A	Wd

- LPL french lexicon : 595.000 entries with frequency computed from a 140 Megawords corpus (mainly newspapers).
- LPL french tagger : Trained on the 700.000 manually tag corrected LPL-Grace corpus. Tags set of 51 categories, score of 0.975 (F-Measure) on written texts.

GN	NV	GN	GA	GP	Wd
Da-fs--d- Ncfs-- La définition	Vmip3s-- connaît	Spd+Da-mp--id Ncfp-- des nuances	Afpfp- importantes	Sp- Dd-mp---- Afpmp- Nemp-- selon ces différents domaines	.

Objective : Insert frontiers and label chunks to form syntactic constituents with flat structure.

- LPL french chunker : Trained on the 100.000 Easy gold standard. Easy grammar of seven constituents (**GN**, **GP**, **NV**, ...). Score of 0.93 (F-Measure) on written texts.



Objective : Form syntactic constituents and tree structure, and indicate functional relations between the constituents.

- LPL french parser : Trained on the 100.000 words LPL-FT corpus, a grammar of 15 constituents (NP, VP, VN, ...) and 9 relations (SUB, OBJ, COORD, ...). Evaluation in progress.

Remove from the transcription the phenomena which are not found in written french :

- hesitation, pause and filled pause (e.g. *heu*, #, +, ...)
- laughs (e.g. @)
- truncations (e.g. *remp-*)

TOE : # *alors moi j'y étais allée déjà je comprends rien à ce qu'ils font*
+ *mais euh j'y étais allée pour remp- pour remplir la salle #*

Filtered TOE : *alors moi j'y étais allée déjà je comprends rien à ce qu'ils font*
mais j'y étais allée pour pour remplir la salle

Speech transcription treatment - Punctuation marks

Punctuation marks are not found in the transcription. Allow the tagger to insert these marks based on written french model :

- **Wd** : Strong punctuations (e.g. ., !) delimiters of sentences
- **Wm** : Weak punctuations (e.g. ,) delimiters of smaller syntactic units

Example of pattern with punctuation marks insertion :

pattern context	insertion	category	proba
Verb Det Noun	-	SubPro	0.01
	Wd	SubPro	0.34
	Wm	SubPro	0.75

Punctuation marks are inserted if they increase the value of the overall sequence probability.

Pseudo-sentence units

je	sai	pas	pourquoi	ils	nous	ont	pris	mais	nous	s'	est	dit	mais	qu'
je	sai	pas	pourquoi	ils	nous	ont	pris	mais	nous	s'	est	dit	mais	qu'
/d										Wm			Wm	
Pp	V	Rgd	Cs	P	Pp1-pj-	Va	Vmps-sm-	Cc	Pp1-po-	Px3fp	V	Vmps-s	Cc	Cs
l-s	mi			p		ip				--	ei	m-		
pr	ver	adver	conjunction	pr	pronoun	au	verb	conjunction	pronoun	prono	a	verb	conj	conju
on	b	b		o		xil				un	u		uncti	nction
NV	GR				NV		NV		GN	NV		NV		
A	A	D	A	A	A	A	A	C	E	A	A	A	B	B

The insertion process allows to identify pseudo-sentence units (and smaller units corresponding to weak punctuations), based on the syntactic content of the transcription.

Pseudo-sentence and prosodic annotation

ip				ip				ip									
ap				ap				ap	ap			ap					
pronoun	verb	pronoun	verb	pronoun	pronoun	auxiliary	verb	adverb	pronoun	verb	pronoun	verb	preposition	determiner	noun	noun	conj
<i>j'</i>	<i>irai</i>	<i>le</i>	<i>voir</i>	<i>je</i>	<i>l'</i>	<i>ai</i>	<i>noté</i>	<i>enfait</i>	<i>je</i>	<i>devais</i>	<i>y</i>	<i>aller</i>	<i>avec</i>	<i>ma</i>	<i>copine</i>	<i>Sabine</i>	<i>par</i>
Id																	
Im				Im				Im									
NV		NV		NV		NV	GR	NV		NV		GP			GN		cc
pronoun	verb	pronoun	verb	pronoun	pronoun	auxiliary	verb	adverb	pronoun	verb	pronoun	verb	preposition	determiner	noun	noun	cc
<i>j'</i>	<i>irai</i>	<i>le</i>	<i>voir</i>	<i>je</i>	<i>l'</i>	<i>ai</i>	<i>noté</i>	<i>enfait</i>	<i>je</i>	<i>devais</i>	<i>y</i>	<i>aller</i>	<i>avec</i>	<i>ma</i>	<i>copine</i>	<i>Sabine</i>	<i>et</i>

Comparison of the automatic syntactic annotation and manual annotation of prosodic intonational phrase (ip) and accentual phrase (ap). Pseudo-sentence frontiers match ip frontiers at a level of 60% (Nesterenko et al. 2010).

Speech transcription treatment - Lexicon adaptation

Some forms do not play their usual syntactic role. These forms are few but frequent (10 % of the corpus). Their lexical tags distribution has to be modified :

- Add or transform lexical distribution frequencies (e.g. *quoi*, *bon*, *bien*, *putain*, *parce que*, ...)

Some other use are more problematic, e.g. *tu vois*, *je crois*, ...

[Light blue bar]														
[Yellow bar]														
[Light yellow bar]			[Light yellow bar]				[Light yellow bar]			[Light yellow bar]				
NV			NV		GP		GP		NV		GP		conjunction	GR
pronoun	pronoun	verb	pronoun	verb	preposition	noun	preposition	noun	pronoun	verb	preposition	noun	<i>comme</i>	adverb
<i>il</i>	<i>nous</i>	<i>restait</i>	<i>je</i>	<i>crois</i>	<i>des</i>	<i>feux</i>	<i>d'</i>	<i>artifice</i>	<i>tu</i>	<i>vois</i>	<i>des</i>	<i>trucs</i>		<i>ça</i>

Syntactic annotations for the transcriptions of the CID corpus :

- Tagging and chunking steps are efficient.
- Identification of larger units (pseudo-sentence) corresponding to punctuation marks found in written texts.
- Evaluation of the performance and manual correction in progress.

The deep parser does not work properly (about 15 % of tree formation) :

- Disfluencies do not allow tree structure formation ?
- Not appropriate for the description of speech in interaction ?