

Language Technology

<http://cs.lth.se/edan20/>
Chapter 10: Techniques for Sequence Prediction

Pierre Nugues

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

September 22, 2022



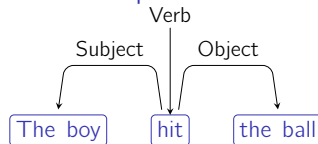
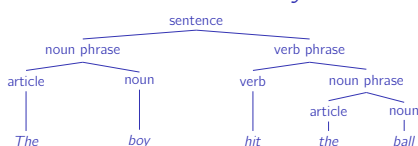
Sequence Prediction

- In this section, we will see how we can apply sequence prediction for part-of-speech tagging to partial parsing.
- We will extend this to other another task we already saw: tokenization



Partial Parsing

- Parsing is the analysis of the relation between the words of a sentence using constituents or dependencies.
- For the sentence *The boy hit the ball*, this corresponds to:



- Parsing might be difficult and useless
- The analysis of parts of a sentence may be enough for many tasks



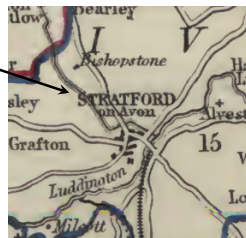
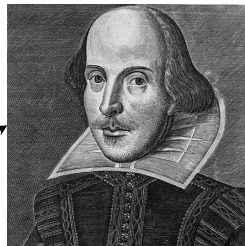
Multiwords

Type	English	French
Prepositions	<i>to the left hand side</i>	<i>À gauche de</i>
Adverbs	<i>because of</i>	<i>à cause de</i>
Conjunctions		
Names	<i>British gas plc.</i>	<i>Compagnie générale d'électricité SA</i>
Titles	<i>Mr. Smith</i> <i>The President of the United States</i>	<i>M. Dupont</i> <i>Le président de la République</i>
Verbs	<i>give up</i> <i>go off</i>	<i>faire part</i> <i>rendre visite</i>



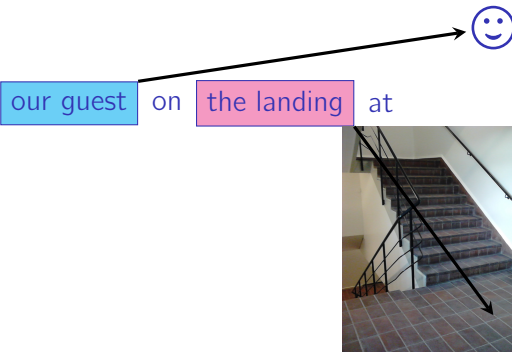
Named Entities: Proper Nouns

William Shakespeare was born and brought
up in Stratford-upon-Avon



Others Entities: Common Nouns

Meeting with our guest on the landing at
lunchtime



Multiword Annotation

The Message Understanding Conferences (MUC), a benchmarking competition organized by the US military, defined an annotation scheme. The MUC annotation restricts the annotation to information useful to the funding source: names (named entities), time expressions, and money quantities.

The annotation scheme defines an XML element for three classes: `<ENAMEX>`, `<TIMEX>`, and `<NUMEX>` with which it brackets the relevant phrases in a text.

The phrases can be real multiwords, consisting of two or more words, or restricted to a single word.



<ENAMEX>

The <ENAMEX> element identifies proper nouns and uses a TYPE attribute with three values to categorize them: ORGANIZATION, PERSON, and LOCATION as in

- The <ENAMEX TYPE="PERSON">Clinton</ENAMEX> government
- <ENAMEX TYPE="ORGANIZATION">Bridgestone Sports Co.</ENAMEX>
- <ENAMEX TYPE="ORGANIZATION">European Community</ENAMEX>
- <ENAMEX TYPE="ORGANIZATION">University of California</ENAMEX> in <ENAMEX TYPE="LOCATION">Los Angeles</ENAMEX>



Named Entities

The detection of named entities and multiwords with regular expressions is an extension of word spotting.

Just as for word spotting, we store them in a Python dictionary.

To get a list of names, we can use geographical or name dictionaries, called **gazetteers**.

We can also model patterns, for example for:

```
'<ENAMEX> M. Dupont </ENAMEX>'
```

and

```
'<NUMEX> 200 euros </NUMEX>'
```



Noun Groups

English	French	German
The waiter <i>is bringing</i> the very big dish <i>on</i> the table	Le serveur <i>apporte</i> le très grand plat <i>sur</i> la table	Der Ober <i>bringt</i> die sehr große Speise <i>an</i> den Tisch
Charlotte <i>has eaten</i> the meal <i>of</i> the day	Charlotte <i>a mangé</i> le plat <i>du</i> jour	Charlotte <i>hat</i> die Tagesspeise <i>gegessen</i>



Verb Groups

English	French	German
<i>The waiter is bringing the very big dish on the table</i>	<i>Le serveur apporte le très grand plat sur la table</i>	<i>Der Ober bringt die sehr große Speise an den Tisch</i>
<i>Charlotte has eaten the meal of the day</i>	<i>Charlotte a mangé le plat du jour</i>	<i>Charlotte hat die Tagesspeise gegessen</i>



Segment Recognition

Group detection – chunking –:

Brackets: [_{NG} The government _{NG}] has [_{NG} other agencies and instruments _{NG}] for pursuing [_{NG} these other objectives _{NG}] .

Tags: *The/I government/I has/O other/I agencies/I and/I instruments/I for/O pursuing/O these/I other/I objectives/I ./O*

Brackets: Even [_{NG} Mao Tse-tung _{NG}] [_{NG} 's China _{NG}] began in [_{NG} 1949 _{NG}] with [_{NG} a partnership _{NG}] between [_{NG} the communists _{NG}] and [_{NG} a number _{NG}] of [_{NG} smaller, non-communists parties _{NG}] .

Tags: *Even/O Mao/I Tse-tung/I 's/B China/I began/O in/O 1949/I with/O a/I partnership/I between/O the/I communists/I and/O a/I number/I of/O smaller/I ,/I non-communists/I parties/I ./O*



Other Chunking Schemes

Tjong and Venstra (1999) created 3 other schemes: IOB1, IOB2, IOE1, and IOB2. A 5th tagset, BIOES, is gaining popularity:

IOB1 : Inside, Outside, Between

IOB2 : Begin, Inside, Outside, possibly the most popular

IOE1 : Inside, Outside, End (between two chunks)

IOE2 : Inside, Outside, End

BIOES : Begin, Inside, Outside, End, and Singleton, the most efficient one.



Other Chunking Schemes

IOB1	Even/O	Mao/I	Tse-tung/I	's/B	China/I	began/O	in/O	1949/I	with/O	a/I	partnership/I	between/O	the/I	communists/I	and/O	a/I	number/I	of/O	smaller/I,	non-communists/I	parties/I											
IOB2	Even/O		Mao/B	Tse-tung/I		's/B	China/I		began/O	in/O	1949/B		with/O		a/B	partnership/I		between/O		the/B	communists/I		and/O		a/B	number/I		of/O		smaller/B,	non-communists/I	parties/I
IOE1	Even/O	Mao/I	Tse-tung/E	's/I	China/I	began/O	in/O	1	949/I	with/O	a/I	partnership/I	between/O	the/I	communists/I	and/O	a/I	number/I	of/O	smaller/I,	non-communists/I	parties/I										
BIOES	Even/O		Mao/B	Tse-tung/E		's/B	China/E		began/O	in/O	1949/S		with/O		a/B	partnership/E		between/O		the/B	communists/E		and/O		a/B	number/E		of/O		smaller/B,	non-communists/I	parties/E



IOB Annotation for Named Entities

CoNLL 2002		CoNLL 2003			
Words	Named entities	Words	POS	Groups	Named entities
Wolff	B-PER	U.N.	NNP	I-NP	I-ORG
,	O	official	NN	I-NP	O
currently	O	Ekeus	NNP	I-NP	I-PER
a	O	heads	VBZ	I-VP	O
journalist	O	for	IN	I-PP	O
in	O	Baghdad	NNP	I-NP	I-LOC
Argentina	B-LOC	.	.	O	O
,	O				
played	O				
with	O				
Del	B-PER				
Bosque	I-PER				
in	O				
the	O				
final	O				
years	O				
of	O				
the	O				
seventies	O				
in	O				
Real	B-ORG				
Madrid	I-ORG				
.	O				



Segment Categorization

Tags extendible to any type of chunks: nominal, verbal, etc.
 For the IOB scheme, this means tags such as I.Type, O.Type, and B.Type, Types being NG, VG, PG, etc.
 In CoNLL 2000, ten types of chunks

Word	POS	Group	Word	POS	Group
<i>He</i>	PRP	B-NP	<i>to</i>	TO	B-PP
<i>reckons</i>	VBZ	B-VP	<i>only</i>	RB	B-NP
<i>the</i>	DT	B-NP	<i>£</i>	#	I-NP
<i>current</i>	JJ	I-NP	<i>1.8</i>	CD	I-NP
<i>account</i>	NN	I-NP	<i>billion</i>	CD	I-NP
<i>deficit</i>	NN	I-NP	<i>in</i>	IN	B-PP
<i>will</i>	MD	B-VP	<i>September</i>	NNP	B-NP
<i>narrow</i>	VB	I-VP	<i>.</i>	.	O

Noun groups (NP) are in red and verb groups (VP) are in blue.



Evaluation

There are different kinds of measures to evaluate the performance of machine learning techniques, for instance:

- Precision and recall in information retrieval and natural language processing;
- The *receiver operating characteristic* (ROC) in medicine.

	Positive examples: P	Negative examples: N
Classified as P	True positives: A	False positives: B
Classified as N	False negatives: C	True negatives: D

More on the receiver operating characteristic here: http://en.wikipedia.org/wiki/Receiver_operating_characteristic



Recall, Precision, and the F-Measure

The **accuracy** is $\frac{|AUD|}{|PUN|}$.

Recall measures how much relevant examples the system has classified correctly, for P :

$$\text{Recall} = \frac{|A|}{|A \cup C|}.$$

Precision is the accuracy of what has been returned, for P :

$$\text{Precision} = \frac{|A|}{|A \cup B|}.$$

Recall and precision are combined into the **F-measure**, which is defined as the harmonic mean of both numbers:

$$F = \frac{2 \cdot \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}.$$



Evaluation: Accuracy, precision, and recall

For noun groups with the predicted output:

Word	POS	Group	Predicted		Word	POS	Group	Predicted
He	PRP	B-NP	B-NP		to	TO	B-PP	B-PP
reckons	VBZ	B-VP	B-VP		only	RB	B-NP	B-NP
the	DT	B-NP	B-NP	X	£	#	I-NP	I-NP
current	JJ	I-NP	B-NP	X	1.8	CD	I-NP	B-NP
account	NN	I-NP	I-NP	X	billion	CD	I-NP	I-NP
deficit	NN	I-NP	I-NP	X	in	IN	B-PP	B-PP
will	MD	B-VP	B-VP		September	NNP	B-NP	B-NP
narrow	VB	I-VP	I-VP		.	.	O	O

There are 16 chunk tags, 14 are correct: $\text{Accuracy} = \frac{14}{16} = 0.875$

There are 4 noun groups, the system retrieved 2 of them: $\text{Recall} = \frac{2}{4} = 0.5$

The system identified 6 noun groups, two are correct: $\text{Precision} = \frac{2}{6} = 0.33$

Harmonic mean = $2 \times \frac{0.33 \times 0.5}{0.33 + 0.5} = 0.4$



Message Understanding Conferences

The Message Understanding Conferences (MUCs) measure the performance of information extraction systems.

They are competitions organized by an agency of the US department of defense, the DARPA

The competitions have been held regularly until MUC-7 in 1997.

The performances improved dramatically in the beginning and stabilized then.

MUCs are divided into a set of tasks that have been changing over time.

The most basic task is to extract people and company names.

The most challenging one is referred to as information extraction.



Information Extraction

Information extraction consists of:

- The analysis of pieces of text ranging from one to two pages,
- The identification of entities or events of a specified type,
- The filling of a pre-defined template with relevant information from the text.

Information extraction then transforms free texts into tabulated information.



An Example

San Salvador, 19 Apr 89 (ACAN-EFE) – [TEXT] Salvadoran President-elect Alfredo Cristiani condemned the terrorist killing of Attorney General Roberto Garcia Alvarado and accused the Farabundo Marti National Liberation Front (FMLN) of the crime...

Garcia Alvarado, 56, was killed when a bomb placed by urban guerrillas on his vehicle exploded as it came to a halt at an intersection in downtown San Salvador...

Vice President-elect Francisco Merino said that when the attorney general's car stopped at a light on a street in downtown San Salvador, an individual placed a bomb on the roof of the armored vehicle...

According to the police and Garcia Alvarado's driver, who escaped unscathed, the attorney general was traveling with two bodyguards. One of them was injured.



The Template

Template slots	Information extracted from the text
Incident: Date	19 Apr 89
Incident: Location	El Salvador: San Salvador (city)
Incident: Type	Bombing
Perpetrator: Individual ID	<i>urban guerrillas</i>
Perpetrator: Organization ID	<i>FMLN</i>
Perpetrator: Organization confidence	Suspected or accused by authorities: <i>FMLN</i>
Physical target: Description	<i>vehicle</i>
Physical target: Effect	Some damage: <i>vehicle</i>
Human target: Name	<i>Roberto Garcia Alvarado</i>
Human target: Description	Attorney general: <i>Roberto Garcia Alvarado</i>
	<i>driver</i>
	<i>bodyguards</i>
Human target: Effect	Death: <i>Roberto Garcia Alvarado</i>
	No injury: <i>driver</i>
	Injury: <i>bodyguards</i>



FASTUS

The FASTUS system has been designed at the Stanford Research Institute to extract information from free-running text

FASTUS uses partial parsers that are organized as a cascade of finite-state automata.

It includes a tokenizer, a multiword detector, and a group detector as first layers.

Verb groups are tagged with active, passive, gerund, and infinitive features.

Then FASTUS combines some groups into more complex phrases and uses extraction patterns to fill the template slots.



FASTUS' Architecture

Sentence

Tokenizer

Multiwords

Part-of-speech
tagging

Group detection
(or chunking)



Tokenization Revisited

- Some Asian languages do not include tokenization marks as in:
然而，這樣的處理也衍生了一些問題。
'However, this treatment also created some problems.'
From Universaldependencies.org
- Tokenized as: 然而||，||這樣||的||處理||也||衍生||了||一些||問題||。
- Shao proposed the tokenization with the tagset: B, I, E, and S, where
 - B is the beginning of a word, I is inside, and E is the end.
 - S is for a single-character word.

然 而 ， 這 樣 的 處 理 也 衍 生 了 一 些 問 題 。
B E S B E S B E S B E S B E S B E S



Adaptation to Other Languages

In other languages, we have tokenization markers, mostly spaces.
We mark them with the X tag.

An example in French:

Chars: On considère qu'environ 50 000 Allemands du Wartheland ont péri pendant la période.
Tags: BEXBIIIIIIIEXBIEBIIIIIEXBIIIIIEXBIIIIIIIEXBEXBIIIIIIIEXBIEBIIIEXBIIIIIEXBEXBIIIIIES

Finally, we can use a final tag T to mark the end of a sentence.
This will enable us to carry out jointly tokenization and the sentence segmentation of a text.



Training the Model

The sentence # sent_id = test-s1

text = 然而，這樣的處理也衍生了一些問題。

The tokenized version from universal dependencies:

ID	FORM	LEMMA	UPOS	XPOS	FEATS	HEAD	DEPREL	DEPS	MISC
1	然而	然而	ADV	RB	—	7	mark	—	SpaceAfter=No
2	,	,	PUNCT	,	—	7	punct	—	SpaceAfter=No
3	這樣	這樣	PRON	PRD	—	5	det	—	SpaceAfter=No
4	的	的	PART	DEC	Case=Gen	3	case	—	SpaceAfter=No
5	處理	處理	NOUN	NN	—	7	nsubj	—	SpaceAfter=No
6	也	也	ADV	RB	—	7	mark	—	SpaceAfter=No
7	衍生	衍生	VERB	VV	—	0	root	—	SpaceAfter=No
8	了	了	AUX	AS	Aspect=Perf	7	aux	—	SpaceAfter=No
9	一些	一些	ADJ	JJ	—	10	amod	—	SpaceAfter=No
10	問題	問題	NOUN	NN	—	7	obj	—	SpaceAfter=No
11	。	。	PUNCT	.	—	7	punct	—	SpaceAfter=No

