

Syntax: tools

(POS tagging, chunking)

Ing. Roberto Tedesco, PhD

roberto.tedesco@polimi.it



arcslab
adaptable, relational and cognitive software environments

NLP – AA 17-18
Prof. L. Sbattella

POS tagging

- Morphological analysis considers one word at a time
- Often, this is not enough to disambiguate part-of-speech
 - E.g. *Talk*: verb or name?
- The context is *needed*
- POS taggers do just that

Methodologies for POS tagging

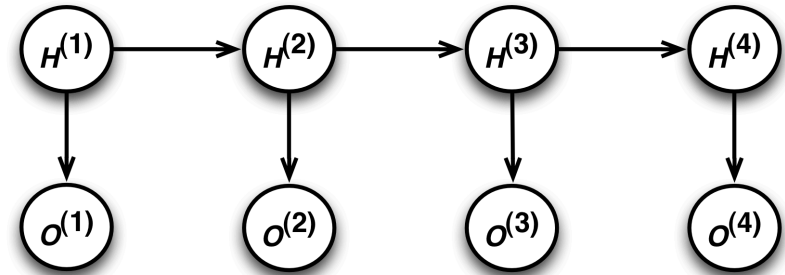
- *A language model* is needed
- Several approaches
 - Stochastic and non stochastic
- Two popular stochastic models:
 - HMM-based
 - Entropy Maximization-based

HMM POS tagging

- Works a sequence at a time
- Hidden variable: tag set
 - E.g, stochastic var. $H=\{t_a, t_b, \dots\}$: tag set (hidden state)
- Observable variable: the set of the word forms (lexicon)
 - E.g. stochastic var. $O=\{w_a, w_b, w_c, \dots\}$: word set (output)
- Given a sequence of words $\langle w^{(i)} \rangle$, Viterbi calculates the most probable sequence of tags $\langle t^{(i)} \rangle$

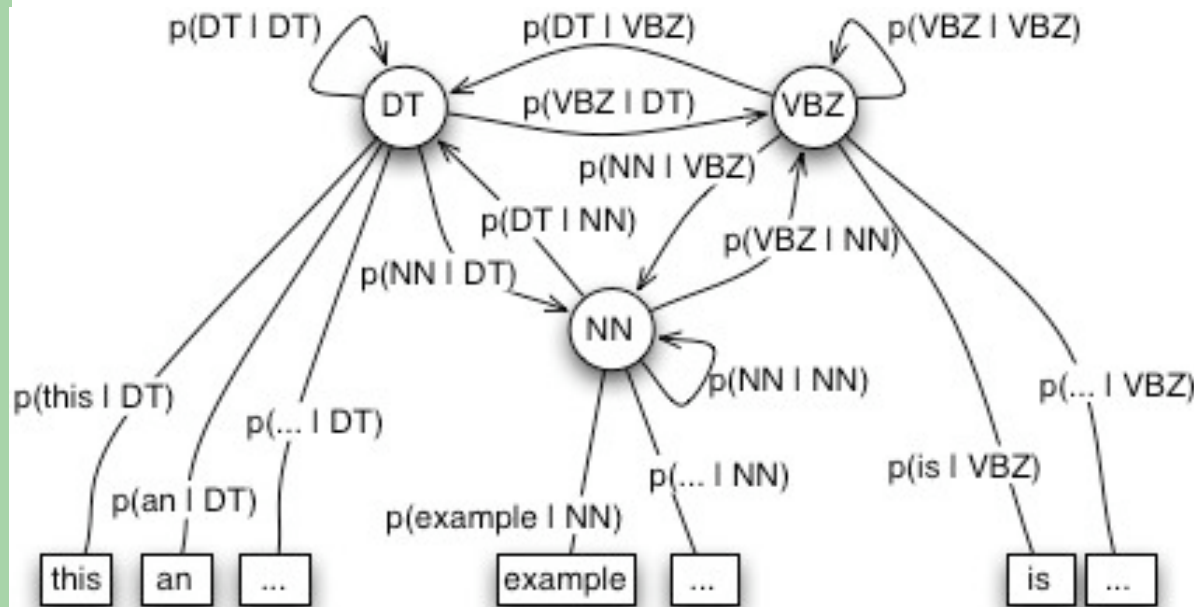
Hidden Markov Models (HMM)

$\langle H^{(1)}=DT, H^{(2)}=VBZ, H^{(3)}=DT, H^{(4)}=NN \rangle$



$\langle O^{(1)}=this, O^{(2)}=is, O^{(3)}=an, O^{(4)}=example \rangle$

Unrolled view



Graph view

- $H^{(t)} = \{DT, VBZ, NN, \dots\}$
- $O^{(t)} = \{\dots, an, \dots, example, \dots, is, \dots, this, \dots\}$

FreeLing

- As a morphologic analyzer

Write your sentences

This is a, quite simple, example

Analysis options

- Multiword detection
- Number recognition
- Date/Time recognition
- Quantities, ratios, and percentages
- Named Entity detection
- Named Entity classification
- No sense annotation
- WN sense annotation: Frequency sorted (MFS disambiguation)
- WN sense annotation: PageRank sorted (UKB disambiguation)

Select language

English

Select output

Morphological Analysis

Submit

Analysis Results

Sentence #1

<u>This</u>	<u>is</u>	<u>a</u>	<u>,</u>	<u>quite</u>	<u>simple</u>	<u>,</u>	<u>example</u>
<i>this</i> DT 0.999824	<i>be</i> VBZ 1	<i>1</i> Z 0.999969	<i>,</i> Fc 1	<i>quite</i> RB 0.935714	<i>simple</i> JJ 0.864583	<i>,</i> Fc 1	<i>example</i> NN 1
<i>this</i> PRP 0.0001755		<i>a</i> DT 1.01887e-05		<i>quite</i> PDT 0.0642857	<i>simple</i> NN 0.135417		
		<i>a</i> NN 1.01887e-05					
		<i>a</i> NNS 1.01887e-05					

FreeLing

- POS tagging
- HMM

Write your sentences

This is a, quite simple, example

Analysis options

- Multiword detection
- Number recognition
- Date/Time recognition
- Quantities, ratios, and percentages
- Named Entity detection
- Named Entity classification
- No sense annotation
- WN sense annotation: Frequency sorted (MFS disambiguation)
- WN sense annotation: PageRank sorted (UKB disambiguation)

Select language English **Select output** PoS Tagging

Analysis Results

Sentence #1

This	is	a	,	quite	simple	,	example
<i>this</i>	<i>be</i>	<i>1</i>	<i>,</i>	<i>quite</i>	<i>simple</i>	<i>,</i>	<i>example</i>
DT	VBZ	Z	Fc	RB	JJ	Fc	NN

Stanford POS tagger

- Stanford POS Tagger
- Entropy Maximization
 - Uses a CMM, basically a MEMM
- Java based

Chunking (aka shallow parsing)

- Identifying and classifying the flat, non-overlapping segments of a sentence
 - This set typically includes noun phrases, verb phrases, adjective phrases, and prepositional phrases
 - [_{NP} The morning flight] [_{PP} from] [_{NP} Denver] [_{VP} has arrived.]
- Leverages POS tagging
- Two approaches:
 - Finite-state rules able to catch phrase segments (FST)
 - Machine learning. We present this approach

Tags (Penn treebank corpus)

TAG	DESCRIPTION	WORDS	EXAMPLE	%
NP	noun phrase	DT+RB+JJ+NN + PR	<i>the strange bird</i>	51
PP	prepositional phrase	TO+IN	<i>in between</i>	19
VP	verb phrase	RB+MD+VB	<i>was looking</i>	9
ADVP	adverb phrase	RB	<i>also</i>	6
ADJP	adjective phrase	CC+RB+JJ	<i>warm and cosy</i>	3
SBAR	subordinating conjunction	IN	<i><u>whether</u> or not</i>	3
PRT	particle	RP	<i><u>up</u> the stairs</i>	1
INTJ	interjection	UH	<i>hello</i>	0

CoNLL corpus

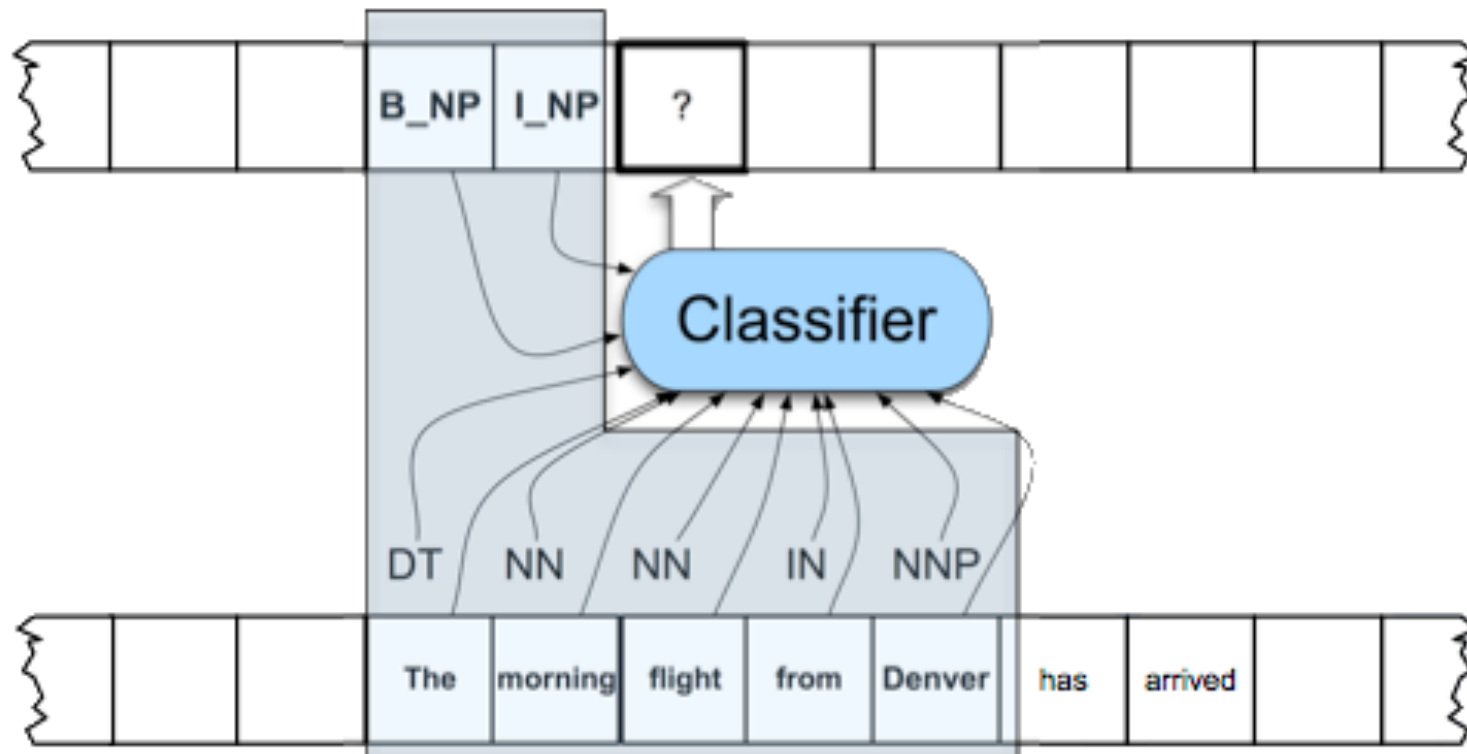
- To train stochastic chunkers
- token, POS, and chunk type
- IBO tagging, for chunk types:

B_ begin of a chunk
I_ inside the chunk
O not part of a chunk

He	PRP	B_NP
reckons	VBZ	B_VP
the	DT	B_NP
current	JJ	I_NP
account	NN	I_NP
deficit	NN	I_NP
will	MD	B_VP
narrow	VB	I_VP
to	TO	B_PP
only	RB	B_NP
#	#	I_NP
1.8	CD	I_NP
billion	CD	I_NP
in	IN	B_PP
September	NNP	B_NP

Machine learning based chunking

- The chunker slides a context window over the sentence classifying words as it proceeds
- At this point the classifier is attempting to label *flights*



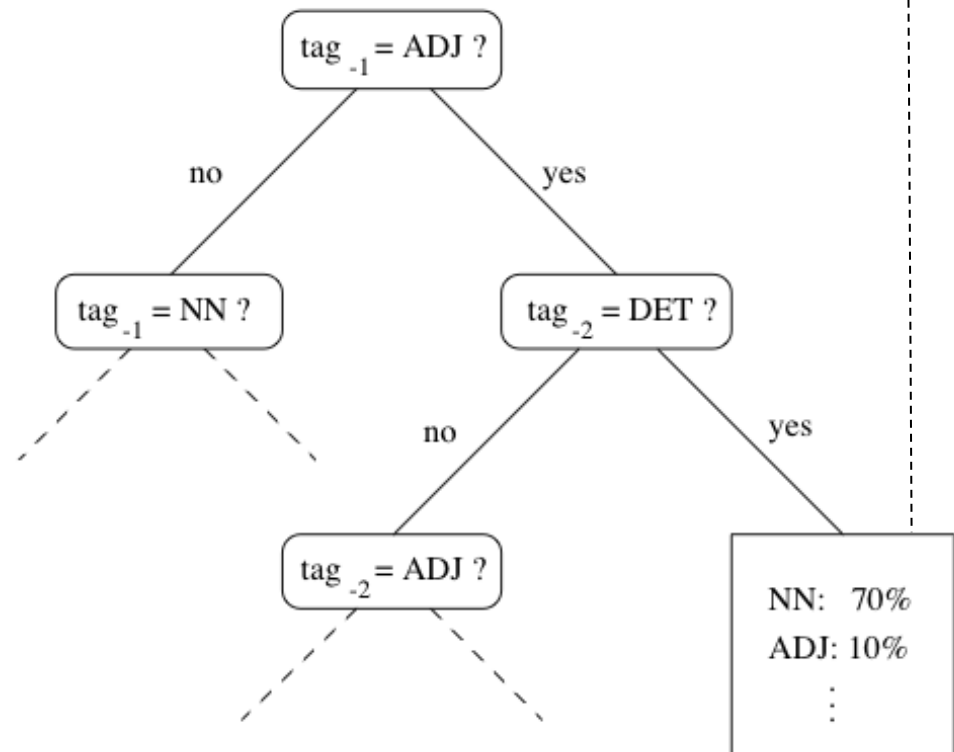
TreeTagger

$$p(T^{(t)}=NN \mid T^{(t-1)}=ADJ, T^{(t-2)}=DET)=0.7$$
$$p(T^{(t)}=ADJ \mid T^{(t-1)}=ADJ, T^{(t-2)}=DET)=0.1$$

...

- POS tagging
- Uses a 2nd order HMM; estimates transition probability by means of a model (not an n-gram)
- Uses a binary decision tree
 - Built from a training corpus of trigrams with POS's

$$p(T^{(t)} \mid T^{(t-1)}, T^{(t-2)})$$



Illinois Chunker: POS + chunking

- HMM-based chunking
- Extends the HMM model: transition probability depends on observation

$$P(T^{(t)} | T^{(t-1)}, W^{(t)})$$



REFERENCES

POS Tagging

- TreeTagger:
 - <http://www.ims.uni-stuttgart.de/projekte/corplex/TreeTagger/DecisionTreeTagger.html>
- Stanford
 - <http://nlp.stanford.edu/software/index.shtml>
- FreeLing (and parser, morpho analyzer, ...)
 - <http://nlp.lsi.upc.edu/freeling/>

Shallow parsers

- CHAOS
 - <http://ai-nlp.info.uniroma2.it/external/chaosproject/>
- TreeTagger:
 - <http://www.ims.uni-stuttgart.de/projekte/corplex/TreeTagger/DecisionTreeTagger.html>
- Illinois
 - http://cogcomp.cs.illinois.edu/page/software_view/13