

ACS Introduction to NLP

Lecture 1: Automatic Linguistic Annotation



UNIVERSITY OF
CAMBRIDGE

Stephen Clark

Natural Language and Information Processing (NLIP) Group

`sc609@cam.ac.uk`

Overall Goal: Automatic Annotation of Linguistic Structure 2

England's fencers won gold on day 4 in Delhi with a medal-winning performance. This is Prof. Briscoe's second gold of the Games.

England's fencers won gold on day 4 in Delhi with a medal-winning performance.

This is Prof. Briscoe's second gold of the Games.

England 's fencers won gold on day 4 in Delhi with a medal -winning performance .

This is Prof. Briscoe 's second gold of the Games .

England|NNP 's|POS fencers|NNS won|VBD gold|NN on|IN
day|NN 4|CD in|IN Delhi|NNP with|IN a|DT medal|JJ
-winning|JJ performance|NN .|.

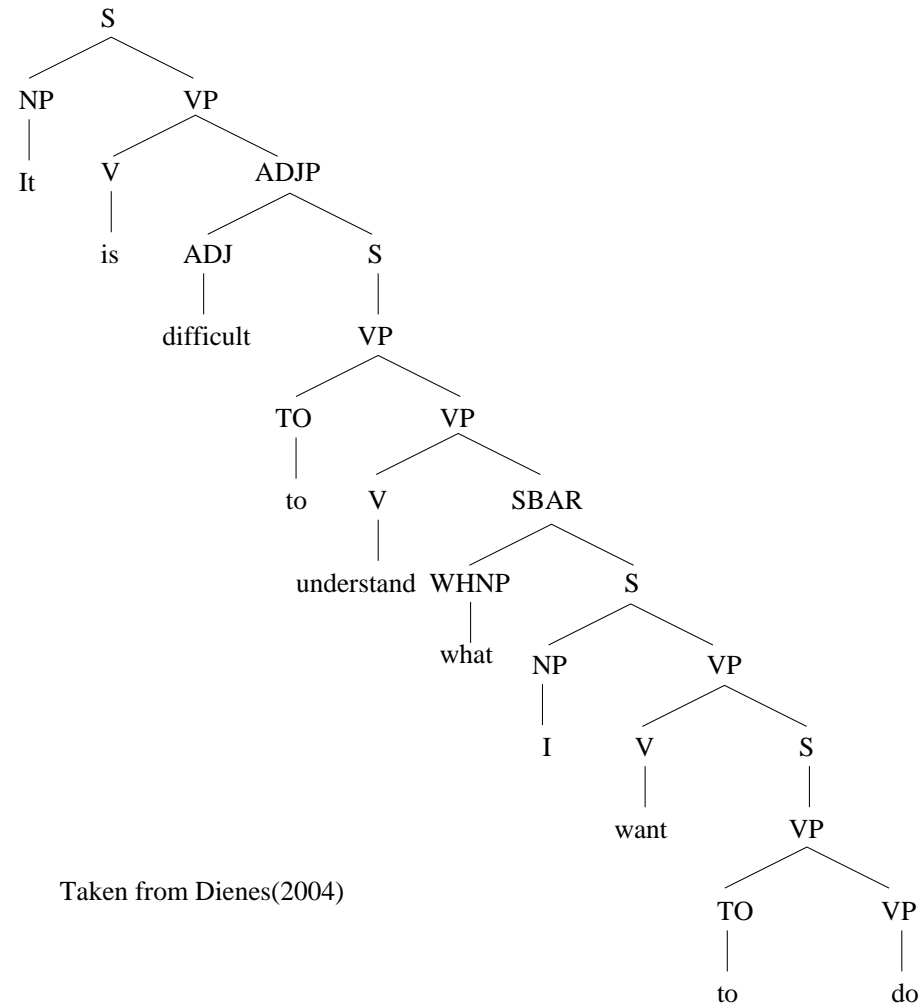
This|DT is|VBZ Prof.|NNP Briscoe|NNP 's|POS second|JJ
gold|NN of|IN the|DT Games|NNP .|.

England|I-LOC 's|O fencers|O won|O gold|O on|O
day|I-TIME 4|I-TIME in|O Delhi|I-LOC with|O a|O medal|O
-winning|O performance|O .|O

This|O is|O Prof.|I-PER Briscoe|I-PER 's|O second|O
gold|O of|O the|O Games|O .|O

England|I-NP 's|I-NP fencers|I-NP won|I-VP gold|I-NP
on|I-PP day|I-NP 4|I-NP in|I-PP Delhi|I-NP with|I-PP
a|I-NP medal|I-NP -winning|I-NP performance|I-NP .|. .

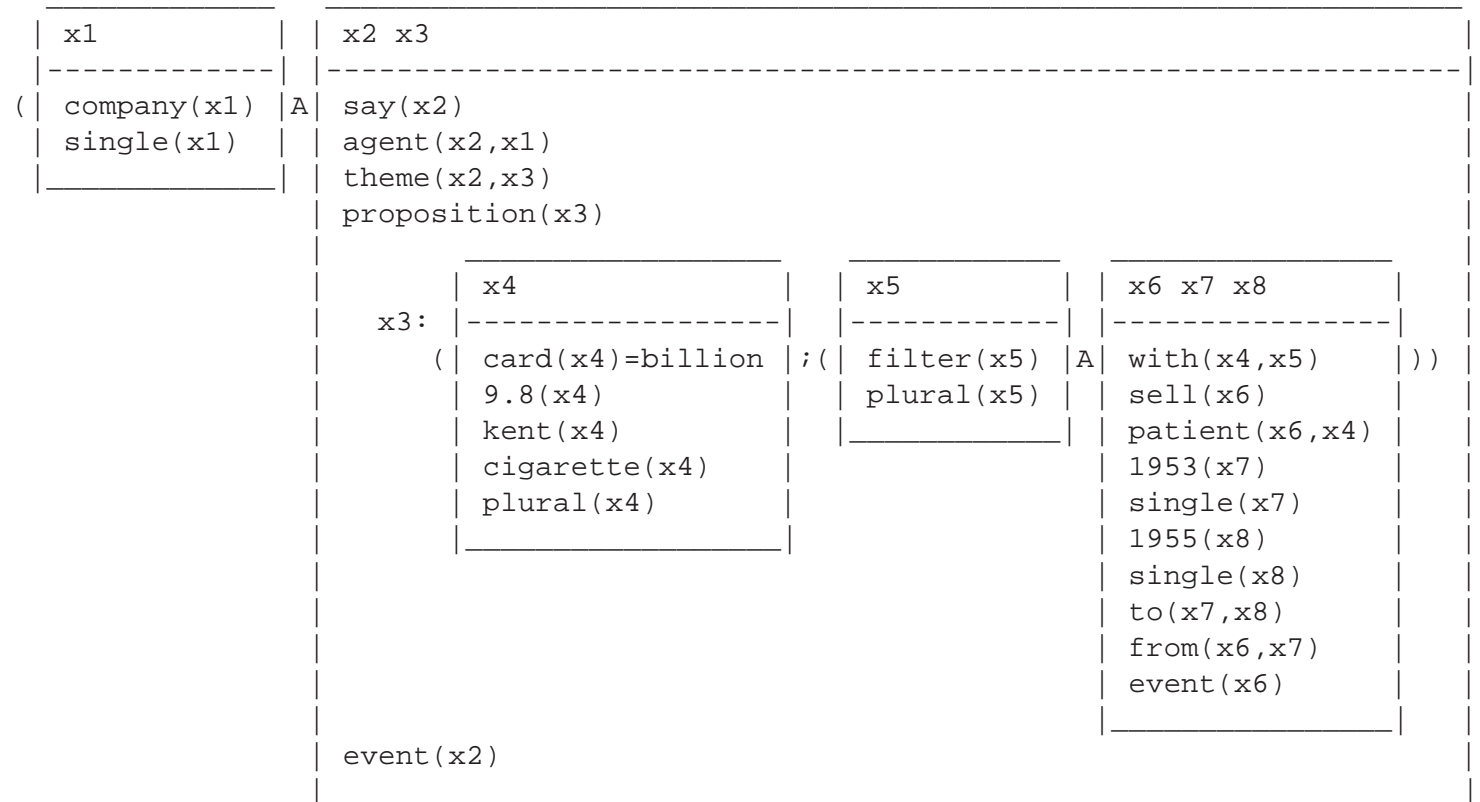
This|I-NP is|I-VP Prof.|I-NP Briscoe|I-NP 's|I-NP
second|I-NP gold|I-NP of|I-PP the|I-NP Games|I-NP .|. .



Taken from Dienes(2004)

Semantic Parsing - Logical Form

From 1953 to 1955 , 9.8 billion Kent cigarettes with the filters were sold , the company said .



-
- Allows the computer access to (elements of) the **meaning** of the sentence (or document)
 - Allows the computer a (rudimentary) “understanding” of the sentence
 - Allows the computer to reason (to some extent) about the sentence

[DEMO: <http://svn.ask.it.usyd.edu.au/trac/candc/wiki/Demo>]

-
- Task: given a set of POS tags and a sentence, assign a POS tag to each word
 - or a set of tags to each word with a probability distribution
 - What are the tags?
 - How does the computer decide which tag to assign to each word?
 - what knowledge is required and where does it come from?
 - What's the algorithm for assigning the tags?

-
- What are the POS tags used for?
 - to provide basic grammatical information, e.g. noun or verb
 - to provide input to more complex annotation, e.g. parsing
 - Example tag sets
 - Penn Treebank set is the most common
 - Others exist, e.g. CLAWS (<http://ucrel.lancs.ac.uk/claws6tags.html>)
 - Choice of tag set may depend to some extent on the algorithm being used to assign the tags

[LOOK AT THE PTB TAG SET]

- AMBIGUITY
- e.g. *can* can be a noun or a (modal) verb

[DEMO]

$$y^* = \arg \max_{y \in Y} \text{score}(y, x)$$

where x is a sentence and Y is the set of possible tag sequences for x

- In machine learning this is known as a *sequence labelling* problem
- There are many possible solutions (HMM, CRF, perceptron, ...)

$$y^* = \arg \max_{y \in Y} P(y|x)$$

where x is a sentence and Y is the set of possible tag sequences for x

- More on the motivation for the probabilistic (statistical) approach in the next lecture
- But for now: we use probabilities because the computer is having to make a guess at the correct tag for a word on the basis of incomplete information
- Probability theory is perhaps the best theory we have for *reasoning under uncertainty*

$$y^* = \arg \max_{y \in Y} P(y|x)$$

where $x = (x_1, \dots, x_n)$ is a sentence and $y = (y_1, \dots, y_n) \in Y$ is a possible tag sequence for x

- Two problems:
 - where do the probabilities come from? (age-old question in statistical approaches to AI)
 - how do we find the arg max?
- Problem 1 is the problem of *model estimation*
- Problem 2 is the *search problem*

-
- Penn Treebank POS tag manual (<http://www.cis.upenn.edu/treebank/>)
 - Jurafsky and Martin, *Speech and Language Processing*, Chapter on Word Classes and Part of Speech Tagging