

# Topics in Natural Language Processing

Shay Cohen

Institute for Language, Cognition and Computation

University of Edinburgh

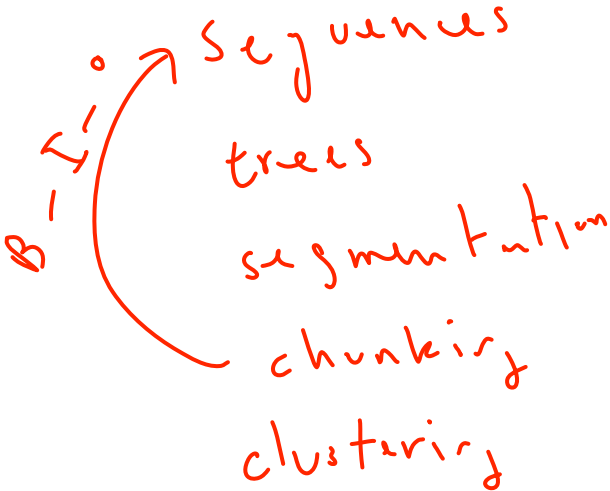
Lecture 6



# Last class

---

## Structure in NLP



# Today's class

---

- Grammars (CFGs and TAGs)
- Inference in NLP

# Context-free grammars

What is a CFG?

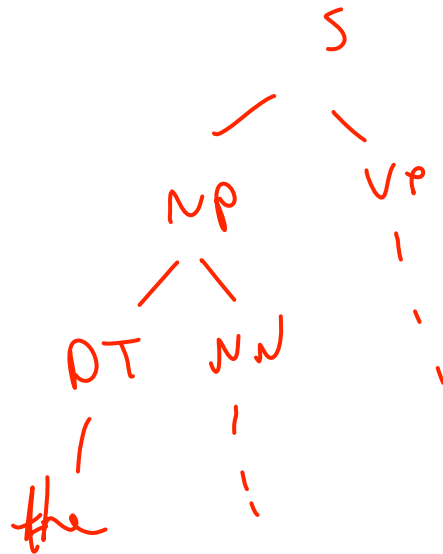
- A set of nonterminals  $N$
- A set of vocabulary terminals  $V$
- A start symbol  $S \in N$
- A set of production rules  $R$ ,  $a \rightarrow \alpha$  is such that  $a \in N$  and  $\alpha \in (V \cup N)^*$

$S \rightarrow NP VP$

$NP \rightarrow DT NN$

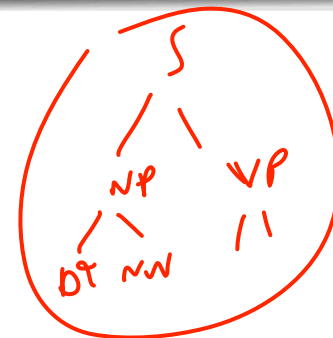
$DT \rightarrow the$

$S$



# Are CFGs sufficient for natural language?

Let  $G$  be a grammar.



$T(G)$  = set of parse trees allowable  
by  $G$

$L(G)$  = set of strings allowed by  $G$

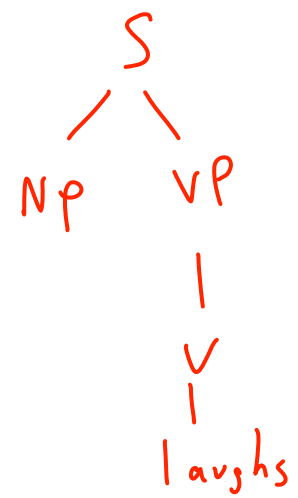
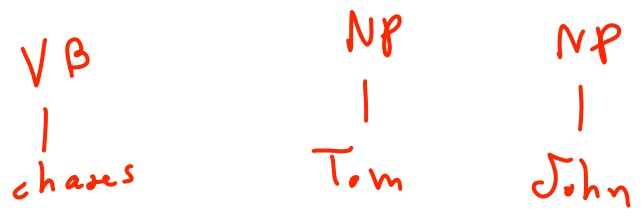
- Dutch - there are structures in Dutch which do not appear in any  $T(G)$  for any  $G$  context-free grammar
- Swiss-German - there are structures in Swiss-German which do not appear in any  $L(G)$  (and hence in any  $T(G)$ ) for any  $G$  CFG

The constructions are similar to demonstrate that. Swiss-German uses case markers.

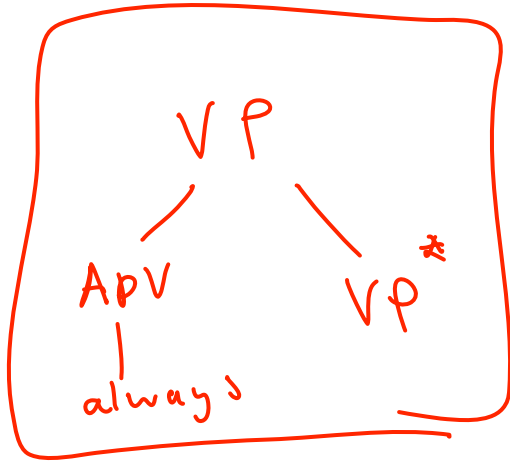
# Tree adjoining grammars

Josh

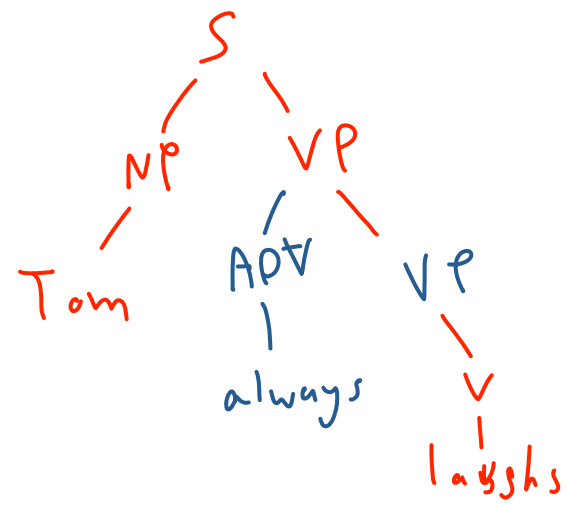
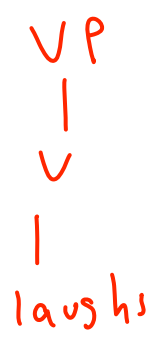
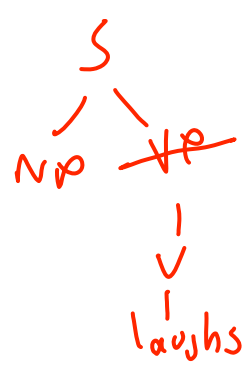
Initial trees:



Auxiliary trees:



Derivational process:



# Tree adjoining grammars

---

Quick question: is  $\{ww \mid w \in \Sigma^*\}$  a context-free language?

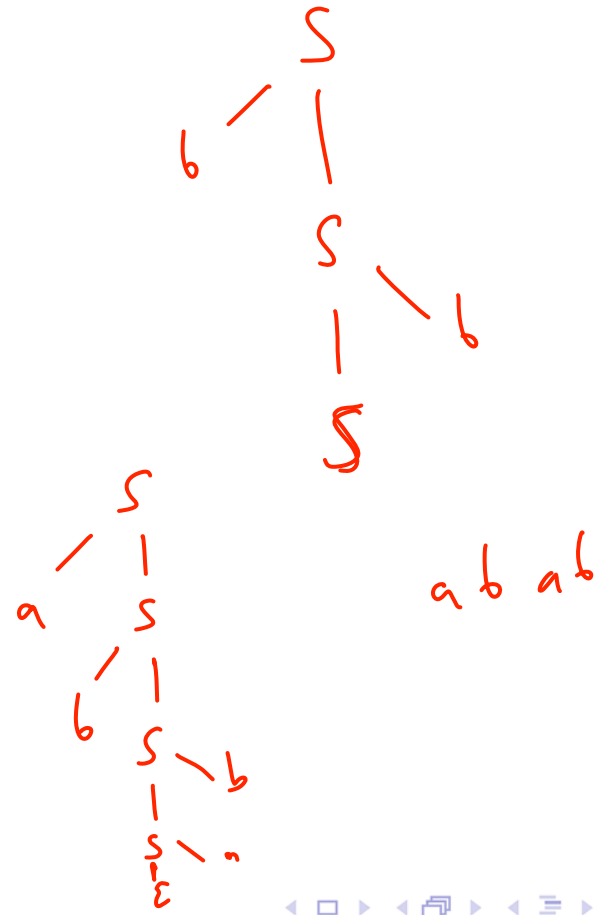
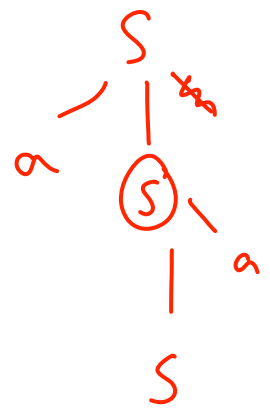
"copy"  
 $a^n b^n$

# Tree adjoining grammars

Quick question: is  $\{ww \mid w \in \Sigma^*\}$  a context-free language?

$\Sigma = \{a, b\}$

Is it a tree adjoining language?



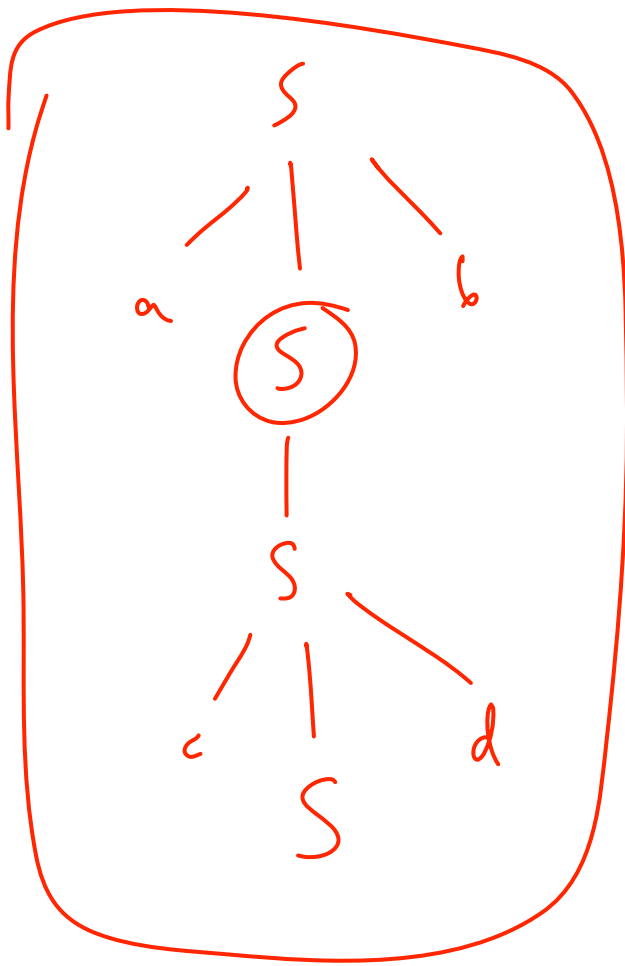


# Tree adjoining grammars

Another quick question: is  $\{a^n b^n c^n d^n \mid n \geq 1\}$  a context-free language?

Aux  
tree

Init. S  
tree |  
 $\epsilon$



# Tree adjoining grammars

---

Another quick question: is  $\{a^n b^n c^n d^n \mid n \geq 1\}$  a context-free language?

Is it a tree adjoining language?

# Tree adjoining grammars

---

They add the “minimum needed” in order to capture phenomena such as cross-serial dependencies

They are part of a family of grammar formalisms called “mildly context sensitive”

Other examples which are weakly equivalent: combinatory categorial grammars, head grammars, linear indexed grammars

# Canonical forms of grammars

---

Canonical form: (1) a specific form for writing a grammar; (2) every general CFG can be converted to an “equivalent” canonical form.

Important example: Chomsky normal form (or binarised form)

$$A \rightarrow B C$$

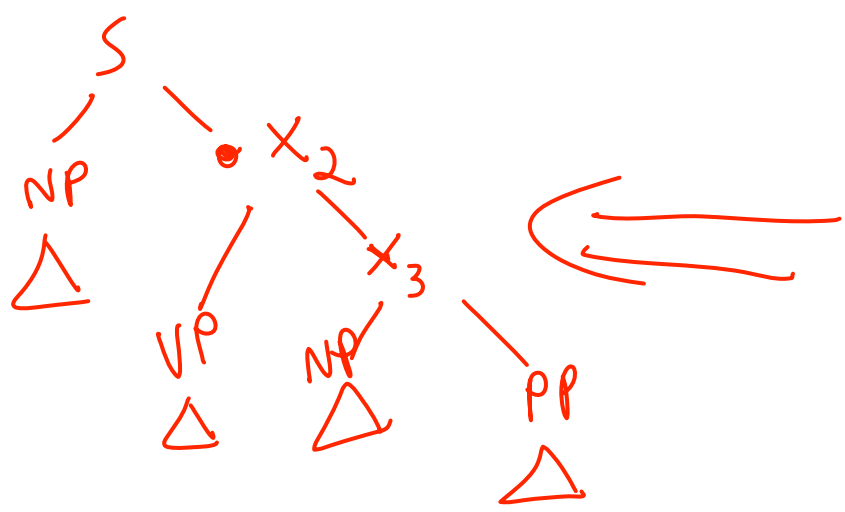
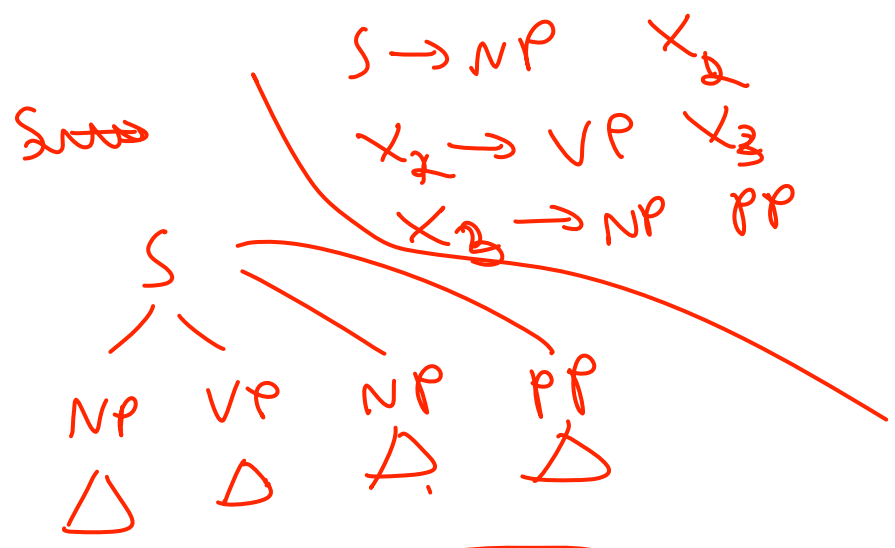
$$A \rightarrow w$$

$A, B, C$  nonterminals

$w$  word

# Why is CNF a normal form?

$S \rightarrow NP \ VP \ NP \ PP$



$G$  in arbitrary form

$\Downarrow$

$G'$  in binary form

can map any

$t \in T(G)$

to

$t' \in T(G')$

$S \rightarrow Y_1 \ PP$

$Y_1 \rightarrow Y_2 \ NP$

$Y_2 \rightarrow NP \ VP$

# Probabilistic grammars

$$p(A \rightarrow \beta) \geq 0$$

for every rule

$$\sum_{A \rightarrow \beta \in R(A)} p(A \rightarrow \beta) = 1 \quad \forall A \in N$$

↑  
rules of A

"Usually", in NLP

$$\sum_{\text{tree} \in T(G)} p(\text{tree}) = 1$$

$$p(\text{tree}) = \prod_{A \rightarrow \beta \in \text{tree}} p(A \rightarrow \beta | A)$$

$$S \rightarrow S S \quad 0.7$$

$$S \rightarrow a \quad 0.3$$

$$\sum p(\text{tree}) < 1$$

# Weighted grammars

Arbitrary positive weights to the rules

$$p(\text{tree}) \propto \prod_{\text{rule}} w(\text{rule})$$

$$p(\text{tree}) = \frac{\prod w(\text{rule})}{Z(G)}$$

$$Z(G) = \sum_{\text{tree} \in T(G)} \prod_{\text{rule} \in \text{tree}} w(\text{rule})$$

# Basic inference with grammars

---

The probability of a derivation:

$p(\text{tree})$  product of rule probabilities

We estimate a PCFG. How do we parse a sentence?



# Estimation

---

We learned how to do estimation:

- Maximum likelihood estimate
- Bayesian posterior summarisation
- ... There are many other ways

What's next?

# Inference

---

What's inference?



Given a statistical model, find probable structure, classification, etc. for the input

# Inference

---

Our  $\Omega$  was usually a cross-product of inputs and outputs

Now, given an input, we need to find the correct output

$$\begin{aligned} \text{output}^* &= \underset{\text{output}}{\text{arg max}} p(\text{output} | \text{input}) \\ &= \underset{\text{output}}{\text{arg max}} \frac{p(\text{output}, \text{input})}{p(\text{input})} = \underset{\text{output}}{\text{arg max}} p(\text{input}, \text{output}) \end{aligned}$$

doesn't depend on output

# Inference

---

Our  $\Omega$  was usually a cross-product of inputs and outputs

Now, given an input, we need to find the correct output

$$\arg \max_{\text{output}} p(\text{output}|\text{input})$$

# Linear Score Function

Consider a model which is a PCFG.

Probability of a tree:

$$p(t) = \prod_{i=1}^n p(r_i) = \prod_{r \in t} p(r) \text{freq}(r, t)$$

$p(r_i)$

“Best” tree  $y$  given sentence  $x$ :

$$t^*(x) = \underset{\text{yield}(t)=x}{\text{arg max}} p(t)$$

$\uparrow$   
sentence

# Linear Score Function

$$\log_2(a \cdot b) = \log_2(a) + \log_2(b)$$

apply  $\Rightarrow \log_2(a^y) = y \log_2 a$

"Best" tree given sentence  $x$ :

take  $\log$

$$y^* = \arg \max_{y: \text{yield}(y)=x} \prod_{r \in y} p(r)^{\text{freq}(y,r)}$$

$$\arg \max_{y: \text{yield}(y)=x} \sum (\log_2 p(r)) \times \text{freq}(r, y)$$

$$= \arg \max_{y, \text{yield}(y)=x} \sum_{r \in R} w(r) \times \text{freq}(y, r)$$

↑ parameters
↑ "structure"

$$\arg \max_y \theta^T f(y, x)$$