

Big attraction:

limited memory versions

BFGS:

$$x_{k+1} = x_k - \alpha_k H_k \nabla f_k$$

$$H_{k+1} = V_k H_k V_k + \rho_k s_k s_k^T$$

$$\rho_k = \frac{1}{y_k^T s_k}, \quad V_k = (I - \rho_k y_k s_k^T)$$

~~Now  $H_k$  is big so we don't want to store it~~

IDEA:

- store the last  $m$   $(s_i, y)$  pairs
- throw away oldest
- have an initial  $H_k^0$  for each step.

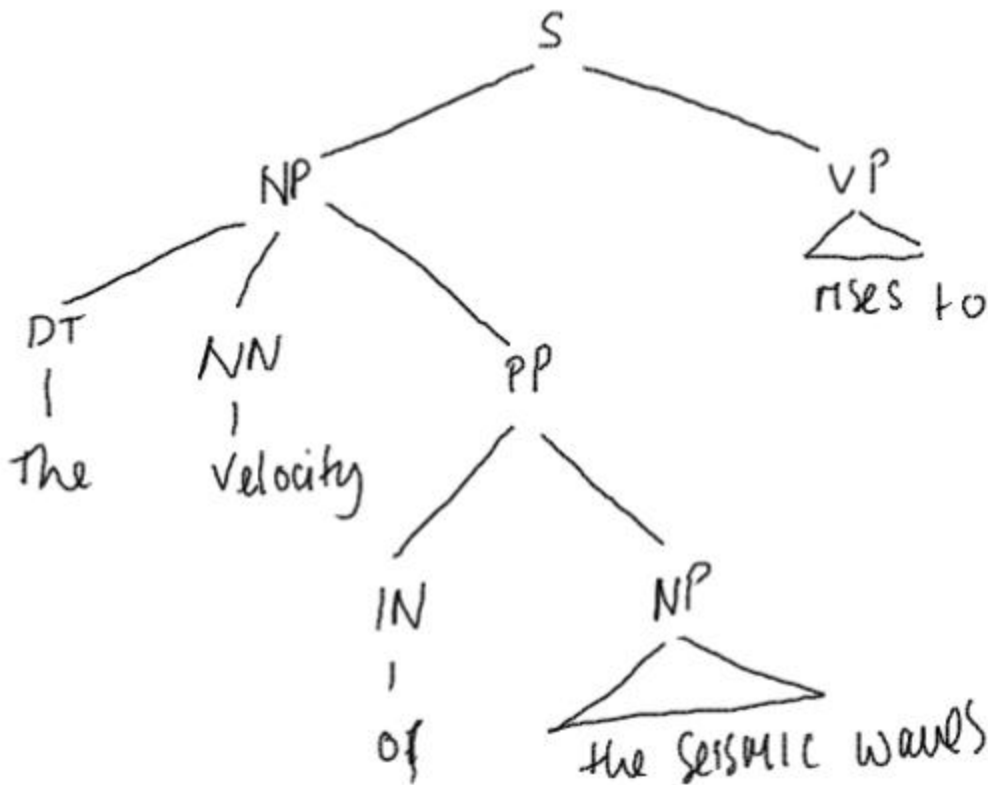
# Parsing.

- We think of strings of language as having structure at long scales:

"The velocity of the seismic waves rises to.

↑  
why is this singular?

- hard to capture with HMM



(2)

How can we extract this structure?

- complex qn with many partial answers
- One answer: Probabilistic Context Free Grammar

~~PCFG consists of~~

- a set of terminals  $\{w^k\}, k=1 \dots N$
- a set of nonterminals  $\{N^i\} \quad i=1 \dots n$
- a start symbol  $N^1$
- a set of rules  $N^i \rightarrow \zeta^j$   
 $\zeta^j$  is a sequence of terminals and non-terminals
- a set of probabilities on rules  
 $P(N^i \rightarrow \zeta^j | N^i) = P(N^i \rightarrow \zeta^j)$

note 
$$\sum_j P(N^i \rightarrow \zeta^j) = 1$$

Notation:

Sentence is  $w_1 \dots w_m$

span  $w_{ab}$  is  $w_a \dots w_b$

if by a set of rewrites we <sup>can</sup> go from  $N^j$  to  $w_{ab}$

then  $N^j$  dominates  $w_{ab}$  and  $w_{ab}$  is yield of  $N^j$



$N^j_{ab}$  means  $N^j$  dominates  $a \dots b$

$$P(w_{im}) = \sum_t P(w_{im}, t)$$

↑  
all trees that yield  $w_{im}$

④

Conditions:

need

Place - invariance:

$$\forall_k \quad P(N_{k(k+c)}^j \rightarrow \zeta) \text{ is the same}$$

Context free:

$$P(N_{ke}^j \rightarrow \zeta \mid \text{anything outside } k-e) \\ = P(N_{ke}^j \rightarrow \zeta)$$

Ancestor free:

$$P(N_{ke}^j \rightarrow \zeta \mid \text{any ancestors outside } N_{ke}^j) \\ = P(N_{ke}^j \rightarrow \zeta)$$

⑤

Under these conditions, computing  $P(w, t)$  is  
straightforward

Properties

- Predictive power tends to be greater for language than HMM's w/ same # of parameters
- PCFG's favour smaller trees over larger trees.
- Probability can be wasted on infinite trees

(6)

Example:

$$S \rightarrow \text{rhubarb} \quad P = 1/3$$

$$S \rightarrow SS \quad P = 2/3$$

$$P(\omega) = \sum_t P(\omega, t)$$

 $\omega$   
 rhubarb

 $1/3$ 

$$\text{rhubarb rhubarb} \quad 2/3 \times 1/3 \times 1/3 = 2/27$$

$$\text{rhubarb rhubarb rhubarb} \quad (2/3)^2 \left(\frac{1}{3}\right)^3 \times 2 = 8/243$$

$\hookrightarrow \wedge$   
 $\hookrightarrow \wedge$

$$\text{So } \sum_{\omega} P(\omega) = \sum_{\omega, t} P(\omega, t)$$

$$= 1/3 + 2/27 + 8/243 + \dots$$

$$= 1/2 \quad !$$

(Half of the probability has gotten stuck in infinite trees)

issues:

eval:

$P(w) ?$

inference:

$\text{argmax}_t P(w, t)$

learning:

rule probs.

- consider only Chomsky Normal Form grammars

rules are:

$$N^i \rightarrow N^j N^k$$

$$N^i \rightarrow w^j$$

} only forms available

parameters are:

$$P(N^i \rightarrow N^j N^k)$$

$$P(N^i \rightarrow w^j)$$

$n^3$  table for  $n$  nonterminals

$nV$  table for  $n$  nonterminals  $V$  terminal

and

$$\sum_{j,k} P(N^i \rightarrow N^j N^k) + \sum_j P(N^i \rightarrow w^j) = 1$$

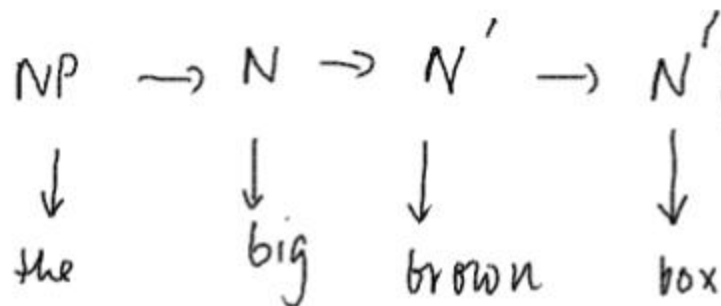


Example:

Probabilistic regular grammar  
(which is rather like an HMM)

$$N^i \rightarrow w^j N^k$$

$$N^i \rightarrow w^j$$



if we add this,  
we have something  
that can be imple-  
mented with an  
HMM

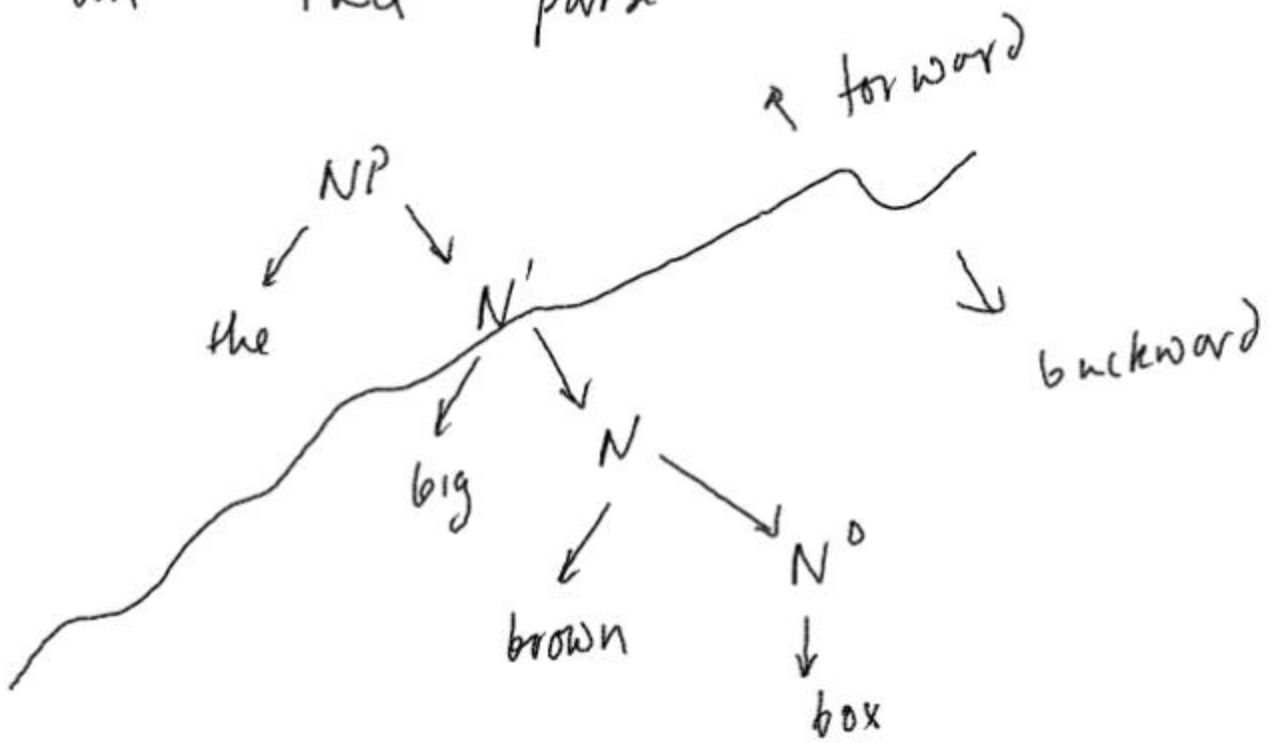
Now in an HMM, we worked with

$$P(w_{1:t-1}, X_t = c)$$

$$P(w_{t:T} | X_t = i)$$

9

built PRG parse



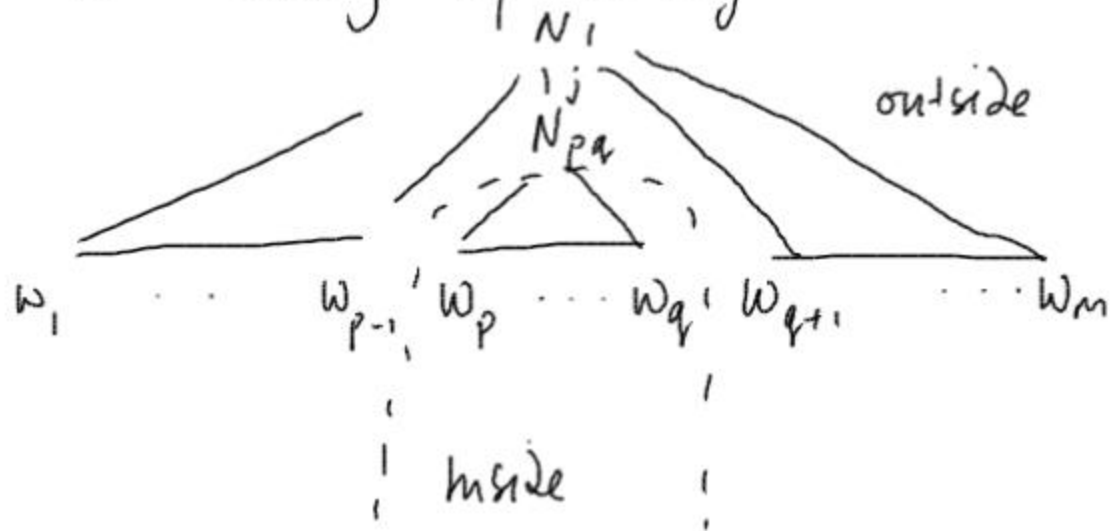
All this suggests an approach to PCFG's

Outside prob  $\alpha_j(p, q)$  ← forward  
 $= P(w_{1:p-1} | N_{pq}^j, w_{p+1:q})$

Inside prob  $\beta_j(p, q) = P(w_{pq} | N_{pq}^j)$

(10)

this is cutting up string



Notice

$$\alpha_j(p, q) \cdot \beta_j(p, q)$$

$$= P(w_1 \dots w_m, \text{tree with } N_{pq}^j \text{ yielding})$$

$$\sum_j \alpha_j(p, q) \beta_j(p, q)$$

$$= P(w_1 \dots w_m, \text{tree with some } N \text{ yielding } w_p q)$$

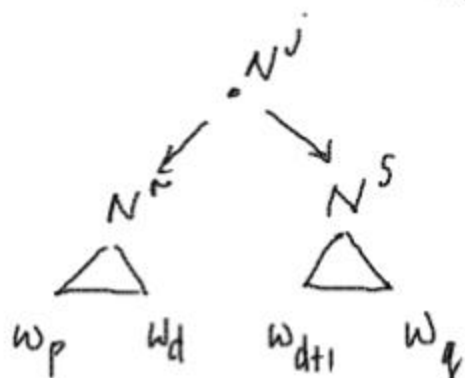
(11)

How to calculate probabilities:

Inside probs:

$$\beta_j(k, k) = P(N^j \rightarrow \omega_k)$$

now consider  $\beta_j(p, q)$



This is a picture of all possible cases.

so

$$\beta_j(p, q) = \sum_{r, s} \left[ \sum_{d=p+1}^{q-1} \beta_r(p, d) \cdot \beta_s(d+1, q) \right] P(N^j \rightarrow N^r N^s)$$

Outside probs:

$$\alpha_1(1, m) = 1$$

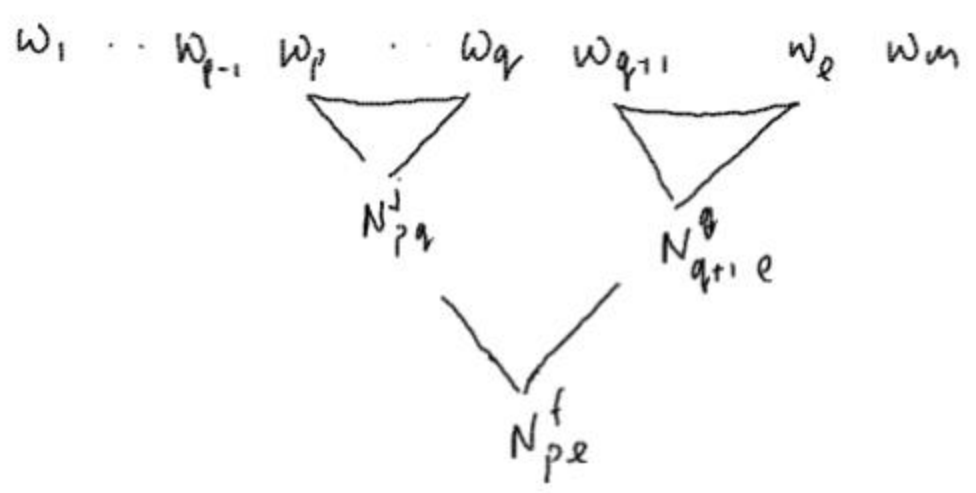
$$= P(\text{tree yielding } w_1 \dots w_m, \text{ rooted at } N^1)$$

$$\alpha_j(1, m) = 0$$

Interior strings

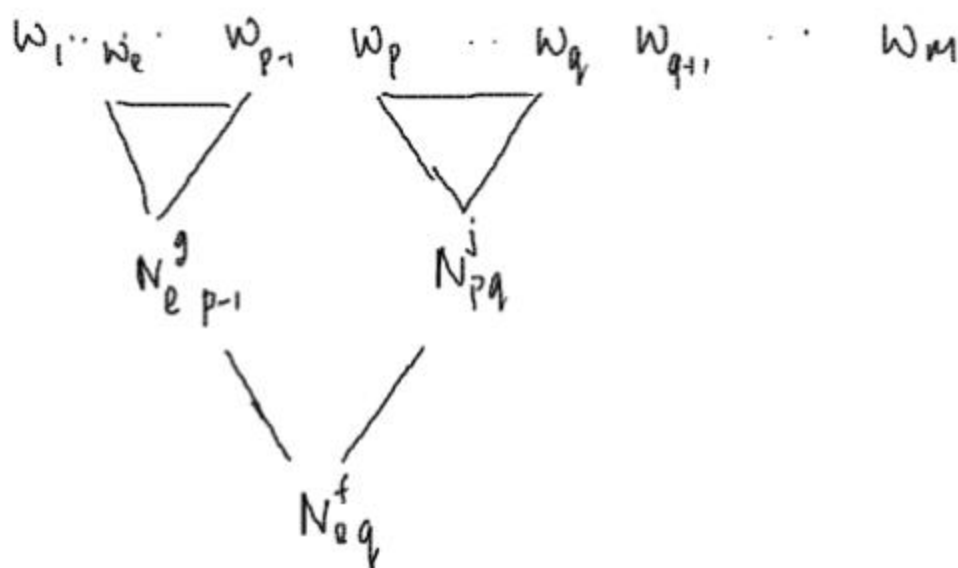
$w_1 \dots w_{p-1} w_p \quad w_q w_{q+1} \dots w_m$

we could have



(13)

So  
or we could have



be careful of  $N^j \rightarrow N^j N^j$

$$\alpha_j(p, q) = \left[ \sum_{f, g \neq j} \sum_{e=q+1}^m \left( \alpha_f(p, e) P(N^f \rightarrow N^j N^g) \right) \right. \\ \left. \times \beta_g(q+1, e) \right] \\ + \left[ \sum_{fg} \sum_{e=1}^{p-1} \alpha_f(e, q) P(N^f \rightarrow N^j N^g) \right. \\ \left. \times \beta_g(e, p-1) \right]$$

So we have an algorithm for computing  $P(\omega)$  (19)

- compute  $\beta$ 's bottom up
- (then compute  $\alpha$ 's top down)

$$\begin{aligned} P(\omega) &= \alpha_1(1 \dots m) \cdot \beta_1(1 \dots m) \\ &= \sum_j \alpha_j(k, k) \cdot \beta_j(k, k) \end{aligned}$$

Inference:

- recall HMM
- key is to keep track of the highest probability path from state  $j$  at time  $t$ , forward
- call this an accumulator

(15)

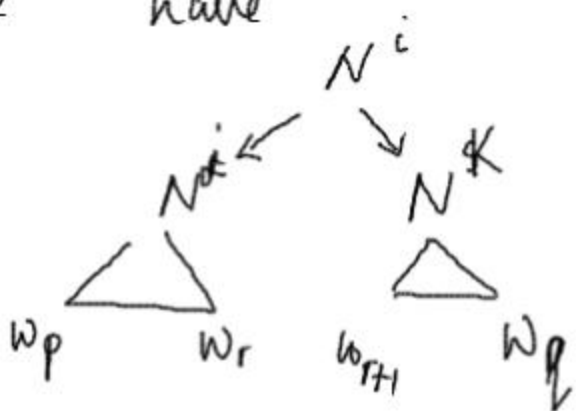
In this case, we want best parse tree with a non-terminal spanning a substring rooted

write

$S_i(p, q) = \text{prob. of highest inside prob. parse tree spanning } p, q, \text{ using } N^i \text{ at top (i.e. } N_{pq}^i)$

1)  $S_i(p, p) = P(N^i \rightarrow w_p)$

2) we have





(16)

so we must have

$$\delta_i(p, q) = \max_{\substack{j, k, \\ p \leq r < q}} \left\{ P(N^i \rightarrow N^j N^k) \times \begin{cases} \delta_j(p, r) \\ \delta_k(r+1, q) \end{cases} \right.$$

how good the best tree is

also, keep

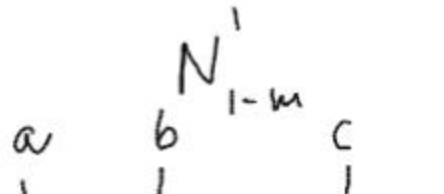
$$\psi_i(p, q) = \text{arg max}$$

which is ~~the tree~~  $(j, k, r)$

3)  $\delta_1(1, M)$  is prob of most prob parse

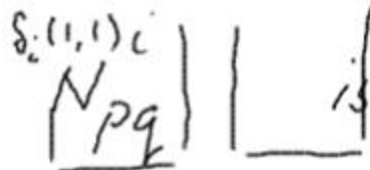
Example

We know root is



now

assume we know



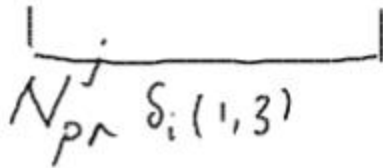
in  $\gamma_i(p,q) = (j,k,l)$

then

$S_i(1,2)$  or  $S_i(2,3)$

left, we must

have



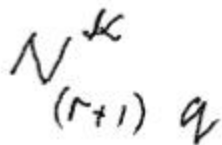
which is either

$S_i(1,1) S_j(2,3)$  or

and

on the right

$S_i(1,2) S_j(3,3)$







## Training:

- assume we know terminals, non-terminals, start, all rules
- must now determine rule probs from data
- Do this with EM

## Single sentence

$$\hat{P}(N^j \rightarrow \zeta) = \frac{\#(N^j \rightarrow \zeta)}{\sum_{\gamma} \#(N^j \rightarrow \gamma)}$$

Now

$$\begin{aligned} \alpha_j(p, q) \beta_j(p, q) &= P(N^i \Rightarrow w_{1m}, \\ &\quad N^j \Rightarrow w_{pq} \mid \text{Grammar}) \\ &= P(N^i \Rightarrow w_{1m} \mid G) \times \\ &\quad P(N^j \Rightarrow w_{pq} \mid N^i \Rightarrow w_{1m}, G) \end{aligned}$$

$$P(N^i \Rightarrow w_{1m} \mid G) = \alpha_1(1, m) \beta_1(1, m) = \pi$$

So:

number ~~fraction~~ of times  $N^j$  is used

$$\begin{aligned} &= \sum_{q=1}^m \sum_{p=1}^m P(N^j \Rightarrow w_{pq} \mid N^i \Rightarrow w_{1m}, G) \\ &= \sum_{q=1}^m \sum_{p=1}^m \frac{\alpha_j(p, q) \beta_j(p, q)}{\pi} \end{aligned}$$

(21)

Two kinds of rule

|              |   |              |    |
|--------------|---|--------------|----|
| nonterminals | → | nonterminals | I  |
|              |   | term         | II |

I: need:

$$\begin{aligned}
 & P(N^j \rightarrow N^r N^s \Rightarrow w_{pq} \mid N^i \Rightarrow w_{im}, G) \\
 &= \sum_{d=p}^{q-1} \alpha_j(p, q) P(N^j \rightarrow N^r N^s) \beta_r(p, d) \beta_s(d+1, q)
 \end{aligned}$$


---


$$\pi$$

then

number of times  $N^j \rightarrow N^r N^s$  given  $N^i$

$$= \sum_{p=1}^{m-1} \sum_{q=p+1}^m \left[ P(N^j \rightarrow N^r N^s \Rightarrow w_{pq} \mid N^i \Rightarrow w_{im}, G) \right]$$

(22)

which yields a re-estimation formula

$$\begin{aligned}
 \text{II} \quad \text{count}(N^j \rightarrow \omega^k \mid N^i \Rightarrow \omega_{1:m}, G) \\
 &= \sum_{h=1}^m \alpha_j(h, h) P(N^j \rightarrow \omega_h, \omega_h = \omega^k) \\
 &= \sum_{h=1}^m \alpha_j(h, h) \delta(\omega_h = \omega^k) \cdot \beta_j(h, h)
 \end{aligned}$$

which gives the re-est formula.



(23)

More than one sentence is only slightly more complex — one must count over all sent.

### Important facts:

- slow:

- each iteration for 1 sentence

- costs  $O(m^3 n^3)$

- sent length  $\uparrow$   $\uparrow$  number of ~~non-term~~

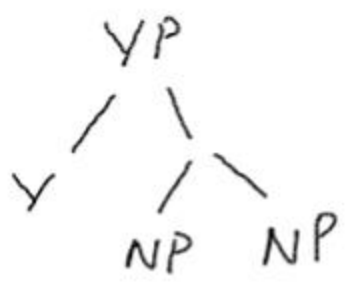
- ~~total maxima are a major~~ problem; eg 300 trials give 300 different local maxima (Charniak)

- Some evidence that method ~~prova~~ needs more non-terminals than are strictly needed

Some options :

- lexicalization :

- Current parser does not know what verb is involved in, say



We should have rules that know what the word is

25

- known as lexicalization
- but how much should be known?
- some tension here between lexicalization (one or more rules per word?) and estimation (many instances of a ~~word~~<sup>rule</sup>)
- try to break words into classes?

Partially unsupervised learning:

- Hard for us to build an activity treebank (Don't know rules)
- But we could segment
-

- Some evidence grammar learning works better on a segmented corpus (Pereira + Schabes; Schabes et al)

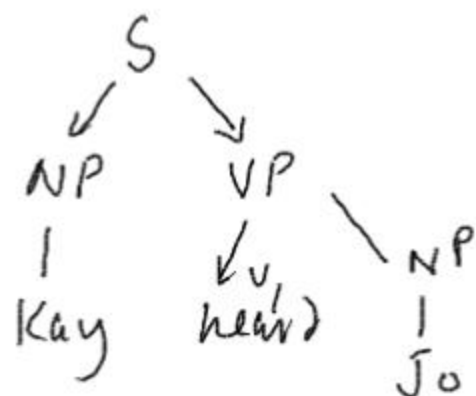
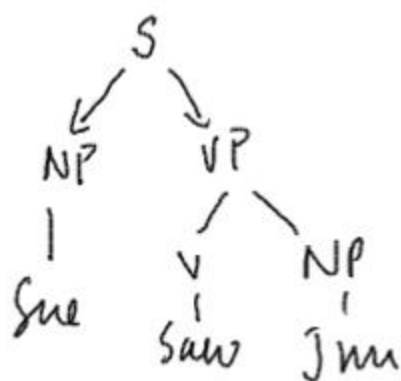
### Learning Structure:

- Some evidence MDL methods apply

### Data oriented parsing:

- Parse by cut + paste of parse trees

We have



we can (a) make new sentns by  
cut + paste

- This is the essence of parsing  
expose what syntactically acceptable  
cut + pastes are

(b) parse sentns (with a very  
nasty search)