

- create parsing table & parse sentence
- convert between dependency tree \Leftrightarrow phrase structure tree
- fill chart parsing table, for given PCFG
- mišić tan hrt i juha "spječite M!" nebo "spječite H"...

1. The Basics: Probability

Suppose your training data contains (only) the sentence (<http://www.petpainter.com>, slightly adapted)

We specialize in petbirds including but not limited to dogs, cats, and horses.

You are building a bigram language model $p(w_i|w_{i-1})$ using maximum likelihood estimation. You do **not** add any special symbols/words to the beginning and end of the use separate unigram distribution for the first position in any given data).

- What is the (raw) probability of $p(\text{cats} | ,)$?
- What is the definition of the training data probability (using, let's say, the fourgram approximation)?

a) $1/2$

$$b) P(y) \approx \prod_i (y_i | y_{i-3}, y_{i-2}, y_{i-1})$$

\rightarrow would be $P(y_0) \cdot P(y_1 | y_0) \cdot P(y_2 | y_0, y_1) \cdot P(y_3 | y_0, y_1, y_2) \dots$
 therefore we generalize into

2. Tagging

- Describe three different approaches to tagging. Include formulas if/where appropriate.

a) HMM tagging

- represented as HMM, where states are actually hidden and can represent some semantic properties of the language. Transitions simulate relations between the language properties. Each state can emit output (tag) as well, closing the gap between relevant (hidden) properties of the language and the desired output.
- is supported by a hypothesis that there are some hidden properties that might very well generate our output.

b) Rule-based tagging

- implemented usually as a set of Reg. Exp.
- Viterbi-like searches (without probs.) are made
- rules \sim FSA ?
must be created by hand
- if still used, then only to very quickly check some basic rules of the language.
- there is problem when it returns multiple possible parses of the input, as there is no additional probability given

c)

Maximum-entropy rule-based tagging

- rules are generated very simply to create big pool
- then we always try to add one new rule
 - we want to have the highest entropy when picking new feature set. \sim high recall
- then we take such combination that gives the lowest entropy. \sim high accuracy

3. Shift-Reduce Parsing

Let's assume the following Context-Free Grammar G (here with seven numbered rules):

- #1 $S \rightarrow K W$
- #2 $K \rightarrow K B$
- #3 $B \rightarrow b$
- #4 $K \rightarrow k$
- #5 $W \rightarrow B W C$
- #6 $C \rightarrow c$
- #7 $W \rightarrow w$

- Decide which of the following strings belong to the language defined by the grammar G :

kkkkbw Yes/No: X
kbbbwcc Yes/No: ✓
kw Yes/No: ✓

- Create the shift-reduce parse tables for this grammar.

0 $S \rightarrow \cdot U W$ $U1$

$U \rightarrow \cdot h B$

$U \rightarrow \cdot h$ $U2$

1 $S \rightarrow U \cdot W$ $W3$

$W \rightarrow \cdot B W C$ $B4$

$U \rightarrow U \cdot B$

$W \rightarrow \cdot w$ $w5$

$B \rightarrow \cdot b$ $b6$

2 $U \rightarrow U \cdot$ $\#4$

3 $S \rightarrow U W \cdot$ $\#1$

4 $W \rightarrow B \cdot W C$ $W7$

$W \rightarrow \cdot B W C$ $B4$

$W \rightarrow \cdot w$ $w5$

$B \rightarrow \cdot b$ $b6$

$U \rightarrow U B \cdot$ $\#2$

5 $W \rightarrow w \cdot$ $\#7$

6 $B \rightarrow b \cdot$ $\#3$

7 $W \rightarrow B W \cdot C$ $C8$

$C \rightarrow \cdot c$ $c9$

8 $W \rightarrow B W C \cdot$ $\#5$

9 $C \rightarrow c \cdot$ $\#6$

4. Statistical Parsing

- What is a Probabilistic CFG? Use five sentences at most.

It is way to introduce language parsing, for example for tagging, with probability included. Problem with typical CFG is that even the most obscure sentence within the training data will have the same weight as other branches of parsing, even though it does not mirror reality. Therefore we introduce:

$$P(T) = \prod_i P(r_i) \rightarrow \text{prob. of parse tree, applying rules } r_i$$

$$P(w) = \sum_j P(T_j) \rightarrow \text{prob. of parse of sentence } w.$$

Now suppose the following "treebank" data, where the pair of symbols (...)XX denotes a (sub)tree rooted in a nonterminal XX.

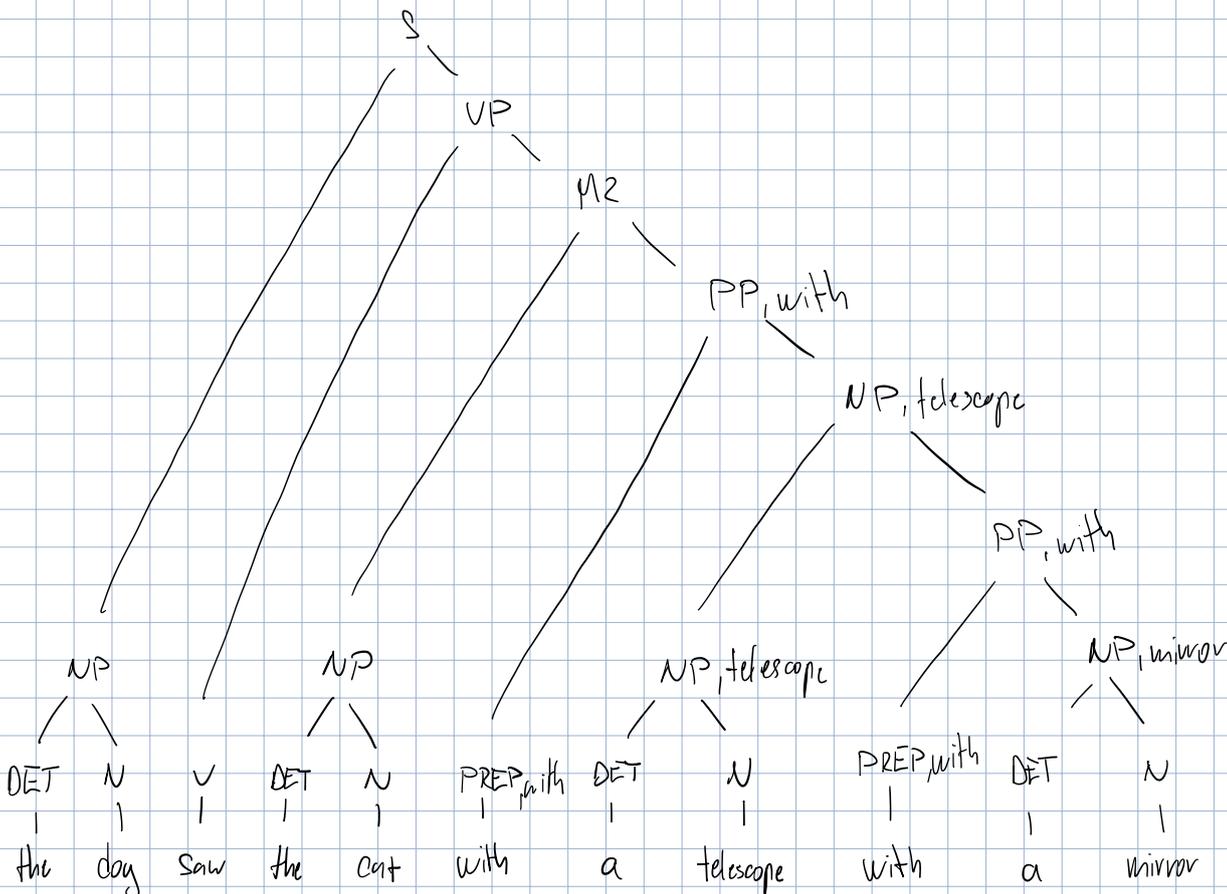
Sentence 1:

```
(( (the) DET
 (dog) N) NP
 (saw) V
 (( (the) DET
 (cat) N) NP
 (with) PREP
 ((a) DET
 (telescope) N) NP
 (with) PREP
 ((a) DET
 (mirror) N) NP) PP) NP) M2) VP) S
```

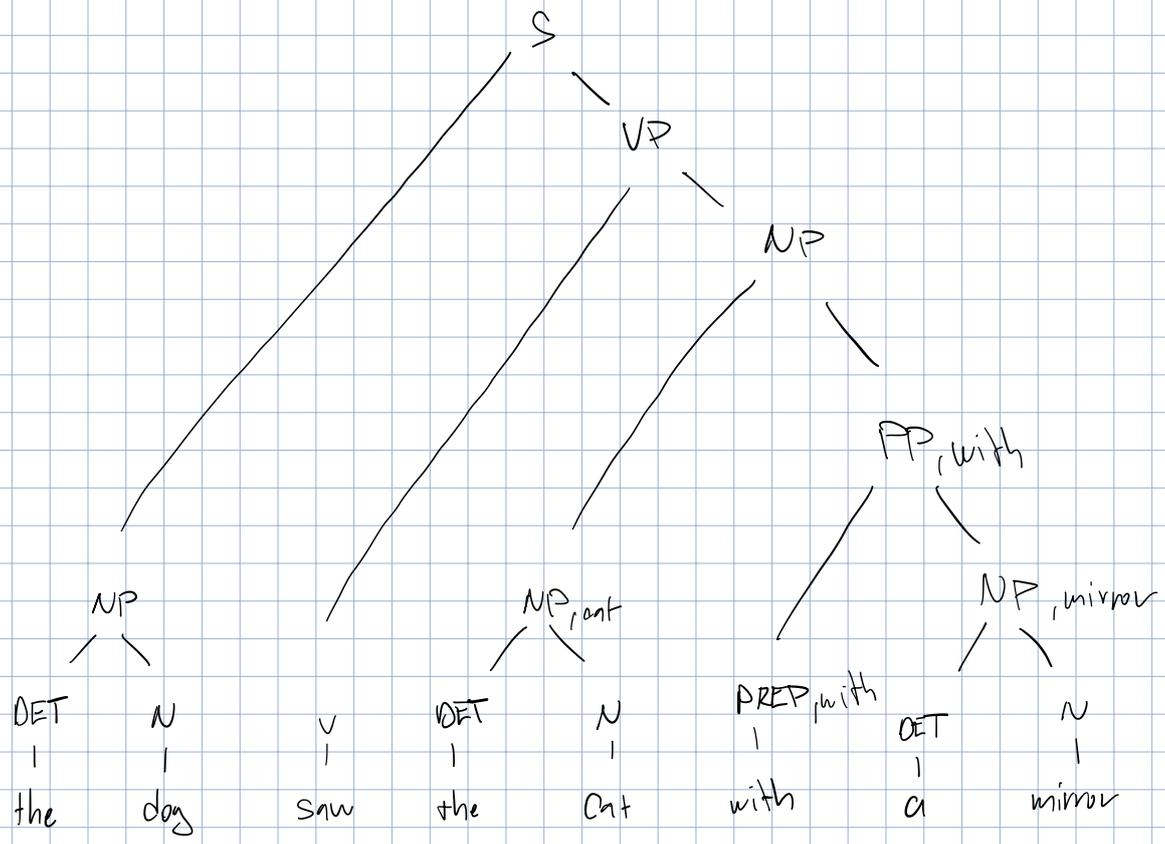
Sentence 2:

```
(( (the) DET
 (dog) N) NP
 (saw) V
 (( (the) DET
 (cat) N) NP
 (with) PREP
 (a) DET
 (mirror) N) NP) PP) NP) VP) S
```

- Draw the usual parse trees, with terminal symbols (words) as leaves, and other nodes labeled by nonterminals (use the next page):



((the)DET
 (dog)N)NP
 ((saw)V
 ((the)DET
 (cat)N)NP
 ((with)PREP
 (a)DET
 (mirror)N)NP)PP)NP)VP)S



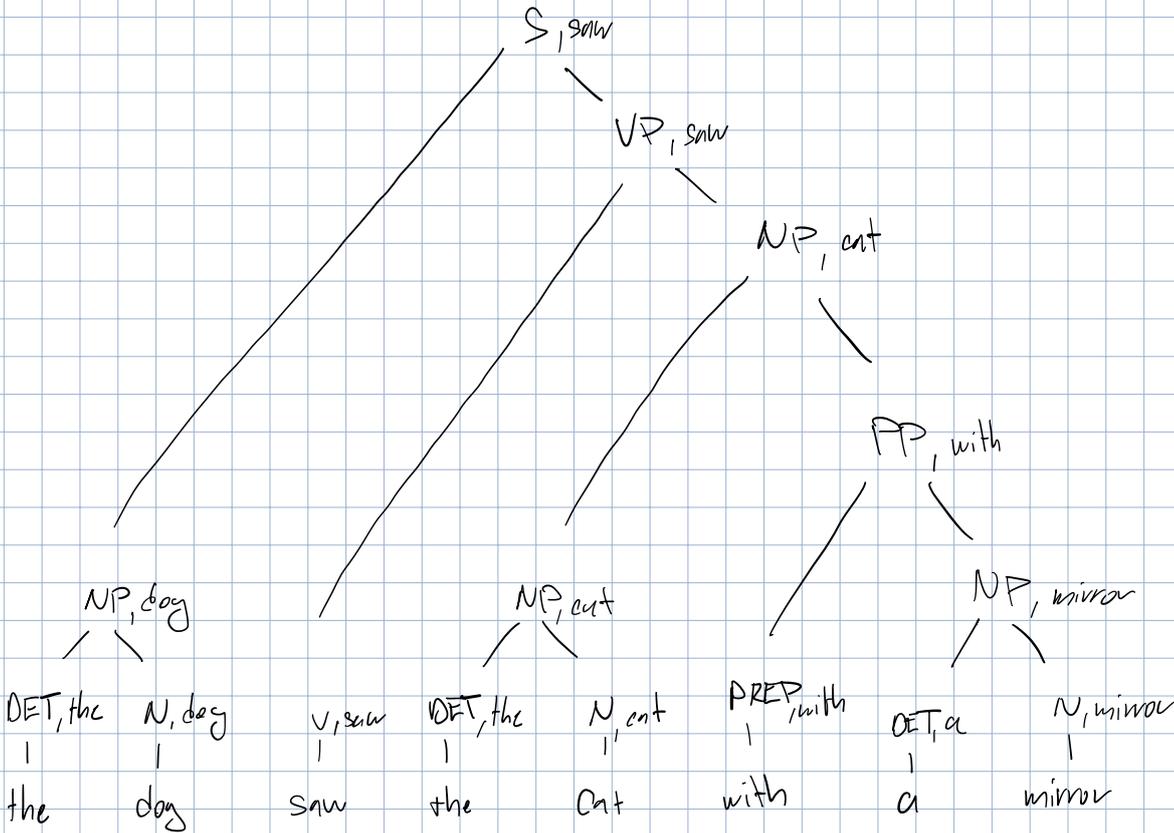
• Write down the (general, non-lexicalized) CFG grammar G as extracted from this training data.

#01	$S \rightarrow NP VP$	2
#02	$VP \rightarrow V M2$	1
#03	$VP \rightarrow V NP$	1
#04	$NP \rightarrow DET N$	7
#05	$NP \rightarrow NP PP$	2
#06	$V \rightarrow Saw$	2
#07	$M2 \rightarrow NP PP$	1
#08	$PP \rightarrow PREP NP$	3
#09	$PREP \rightarrow with$	3
#10	$DET \rightarrow the$	4
#11	$DET \rightarrow a$	2
#12	$N \rightarrow dog$	2
#13	$N \rightarrow cat$	2
#14	$N \rightarrow mirror$	2
#15	$N \rightarrow telescope$	1

OK. Now you have everything to answer these "numerical" questions:

- Estimate $p(NP \rightarrow NP PP) = \frac{2}{9}$
- Estimate $p(NP \rightarrow PREP N) = 0$
- Estimate $p(VP \rightarrow V NP) = \frac{1}{2}$
- Estimate $p(VP \rightarrow V M2) = \frac{1}{2}$

- Now draw the **lexicalized** parse tree for the **second** sentence, assuming that whenever there is a choice of heads in any given subtree, chose (in the following precedence order): VP, V, PREP, N, NP.



Suppose we stick to the usual three independence assumptions and define a probability of a parse tree \underline{g} (of a sentence W) as a simple product of $p(\langle \text{rule} \rangle)$ over all rules used in the course of generation of the sentence W .

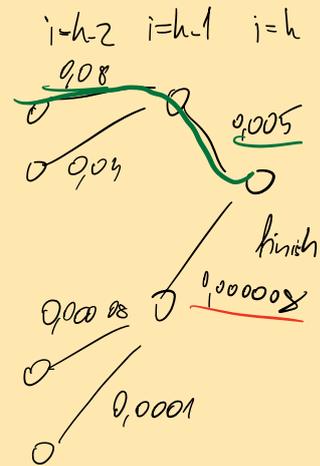
Describe the main idea of the algorithm to compute:

- (a) the best parse given a sentence and a model. Use five sentences at most.

We can modify Viterbi to always remember for each state the most likely path to it. If we get to the final word, we then backtrack and go over the most likely paths.

- (b) the probability of a string given the model. Use five sentences at most.

Now instead of always taking just one argument, if we come to a state from two paths, we sum the paths, making the state more likely to happen as more paths lead to it. We get the probability as just the prob. at the final state.



- What is then the probability of the following parse tree, using the estimated nonlexicalized PCFG:

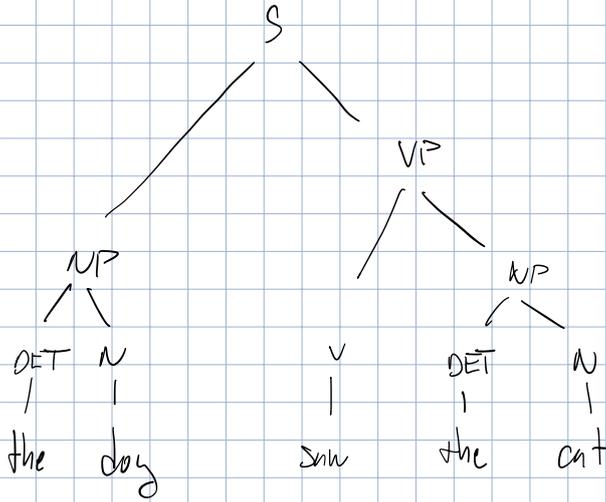
$p(((\text{the})\text{DET} (\text{dog})\text{N})\text{NP} ((\text{saw})\text{V} ((\text{the})\text{DET} (\text{cat})\text{N})\text{NP})\text{VP})\text{S}) = \underline{\hspace{2cm}} \quad 0,008062$

- What is the probability of the string "the dog saw the cat with a mirror"?

$p(\text{sentence}_2) = \underline{\hspace{2cm}} \quad 0,0009385$

- Now a lexicalized grammar question. What is the (raw) probability estimate of the lexicalized rule $p(\text{PP}(\text{with}) \rightarrow \text{PREP}(\text{with}) \text{NP}(\text{mirror}))$? [Imagine you generate the lexicalized rules, too, under the same independence assumptions; however, to answer this question, there is no need to generate them all, of course. However, do not forget to look at **both** training sentences!]

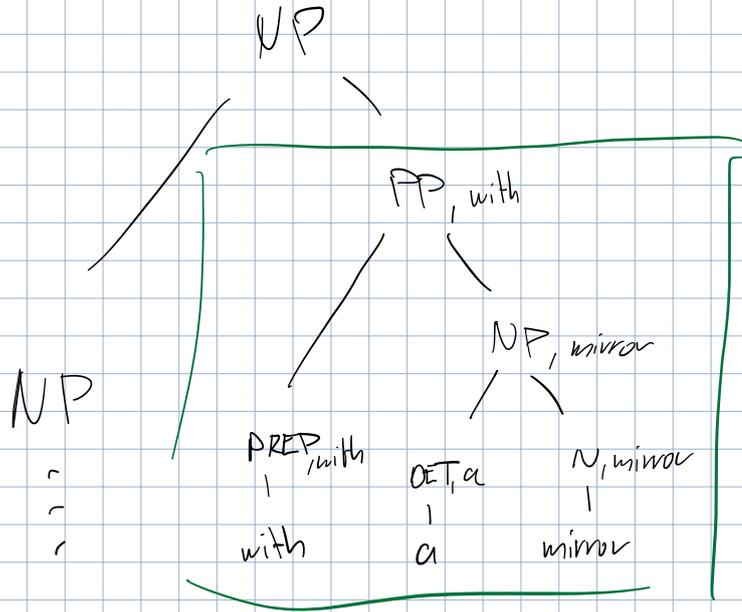
$p(\text{PP}(\text{with}) \rightarrow \text{PREP}(\text{with}) \text{NP}(\text{mirror})) = \underline{\hspace{2cm}} \quad 2/3$



- 0) $S \rightarrow \text{NP VP}$ 1
- 1) $\text{NP} \rightarrow \text{DET N}$ 7/9
- 2) $\text{DET} \rightarrow \text{the}$ 4/2
- 3) $N \rightarrow \text{dog}$ 2/7
- 4) $\text{VP} \rightarrow \text{V NP}$ 1/2
- 5) $V \rightarrow \text{saw}$ 1
- 6) $\text{NP} \rightarrow \text{DET N}$ 7/9
- 7) $\text{DET} \rightarrow \text{the}$ 4/7
- 8) $N \rightarrow \text{cat}$ 2/7

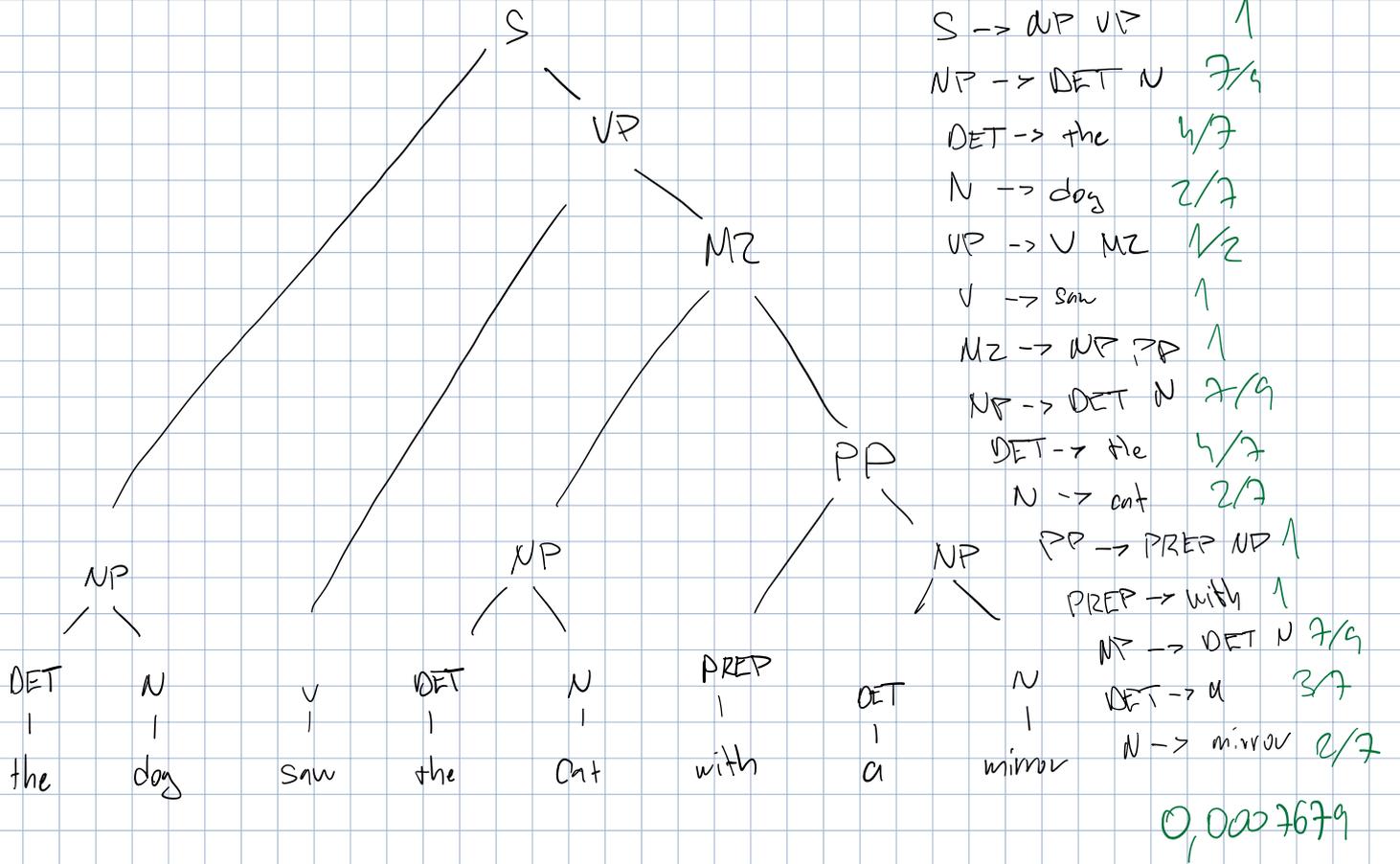
||

$0,008062 \dots$



- $\text{NP} \rightarrow \text{NP PP}$ 2/9
- $\text{PP} \rightarrow \text{PREP NP}$ 1
- $\text{NP} \rightarrow \text{DET N}$ 7/9
- $\text{DET} \rightarrow a$ 3/7
- $N \rightarrow \text{mirror}$ 2/7

$0,0001706 =$



Another shift-reduce example:

0	$S \rightarrow \cdot A R$	A1
	$A \rightarrow \cdot A B$	
	$A \rightarrow \cdot a$	a2
1	$S \rightarrow A \cdot R$	R3
	$A \rightarrow A \cdot B$	B4
	$R \rightarrow \cdot B R C$	
	$R \rightarrow \cdot r$	r5
	$B \rightarrow \cdot b$	b6
2	$A \rightarrow a \cdot$	#3
3	$S \rightarrow A R \cdot$	#1
4	$A \rightarrow A B \cdot$	#2
	$R \rightarrow B \cdot R C$	R7
	$R \rightarrow \cdot B R C$	B4
	$R \rightarrow \cdot r$	r5

aaaaabr	x	
abbbrc	✓	
ar	✓	
stack	states	Comments
abbbrc	0	
bbbr	20	
A bbbrc	0	#3
bbbr	10	
bbrc	610	
B bbr	10	#4
bbrc	410	
brc	6410	
B brc	410	#4
brc	4410	
rc	64410	
Brc	4410	#4
rc	44410	

	B	→ .S	b6
5	R	→ r.	#6
6	B	→ b.	#4
7	R	→ BR.C	C8
	C	→ .c	c9
8	R	→ BRc.	#5
9	C	→ c.	#7

c	54410	#6
Rc	4410	#6
c	74410	#7
-	974410	#7
C	74410	#7
-	874410	#7
R	4410	#5
-	74410	#7
AR	10	#2
ARc	410	#2
Arc	410	#2
ABrc	10	#2
Abrc	10	#2
ABbre	0	#2
Bbre	10	
brc	410	
rc	6410	
Brc	410	#4
rc	4410	
c	54410	
Rc	4410	#6
c	74410	
-	974410	
C	74410	#7
-	874410	
R	410	#5
-	7410	

Entropy: $H(P) = - \sum_x P(x) \cdot \log P(x)$ NLL

Perplexity = $2^{H(P)}$

$$H(P_{\text{bigram}}) = - \sum_x P(x,y) \cdot \log P(x|y)$$

Cross Entropy:

$$H(P, Q) = - \sum_x P(x) \cdot \log Q(x)$$

$$H(P, X) \geq H(P) \quad \forall X$$

$$H(P, P) = H(P)$$

we: $1/16$

specialize: $1/16$

in: $1/16$

petbirds: $1/16$

including: $1/16$

but: $1/16$

not: $1/16$

limited: $1/16$

to: $1/16$

dogs: $1/16$

1: $2/16$

cats: $1/16$

and: $1/16$

horses: $1/16$