

POS Tagging

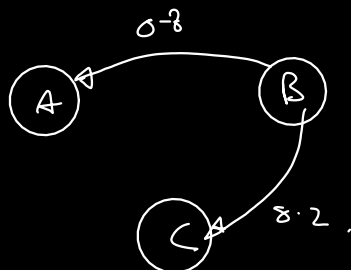
— Most frequent class baseline.

Markov Assumption : $P(q_i = a \mid q_1 \dots q_{i-1}) = P(q_i = a \mid q_{i-1})$

$$P(A|B) = 0.8$$

$$P(C|B) = 0.2$$

\Rightarrow



A	a_1	a_2	...	a_n
a_1	a_{11}	a_{12}		a_{1n}
a_2	a_{21}	a_{22}		
\vdots				
\vdots				
\vdots				
a_n	a_{n1}			a_{nn}

$$Q: q_1 q_2 \dots q_n$$

$$A = a_{11} a_{12} \dots a_{nn}$$

$$\pi = \pi_1 \pi_2 \dots \pi_n$$

HMM Model

Tags: "hidden". Words observed.

$$t: n$$

$$w: o$$

word tags

$$Q: q_1 q_2 \dots q_n$$

$$A: a_{11} a_{12} \dots a_{nn}$$

$$O: o_1 o_2 \dots o_T$$

$$B: b_i(o_t)$$

$$\pi = \pi_1 \pi_2 \dots \pi_n$$

A	a_1	a_2	...	a_n
a_1	a_{11}	a_{12}		a_{1n}
a_2	a_{21}	a_{22}		
\vdots				
\vdots				
\vdots				
a_n	a_{n1}			a_{nn}

sequence of observations (words)

$$b_i(o_t): P(o = o_t \mid q = q_i)$$

each o_t : N 'b's

T 'o's

$$P(q_i \mid a_1 \dots a_n) = P(q_i \mid a_{i-1})$$

$$P(o_i \mid q_1 \dots q_T, o_1 \dots o_T) = P(o_i \mid q_i)$$

Basic HMM problems

1. Likelihood of a sequence

given observation sequence $O = \{o_1, o_2, \dots, o_T\}$

efficiently estimate $P(O|\lambda)$

$$O = \{o_1, o_2, \dots, o_T\}$$

$$Q = \{q_1, q_2, \dots, q_T\}$$

Complexity:

$$O(2TN^3)$$

Now,

$$P(O|Q, \lambda) = \prod_{i=1}^T P(o_i | q_i, \lambda)$$

$$= b_{q_1}(o_1) \dots b_{q_T}(o_T)$$

$$P(Q|\lambda) = \prod_q a_{q_1 q_2} \dots a_{q_{T-1} q_T} a_{q_T}$$

$$P(O|\lambda) = \sum_{\prod_q a_{q_1 q_2} \dots a_{q_{T-1} q_T} b_{q_1}(o_1) \dots b_{q_T}(o_T)} [P(A, B) = P(A|B) \cdot P(B)]$$

Forward Procedure :

DEF:

$$\alpha_t(i) = P(o_1, \dots, o_t, q_t = S_i | \lambda)$$

given partial obs. seq. \rightarrow Prob. that state at position t is S_i

STEP:

1. Init

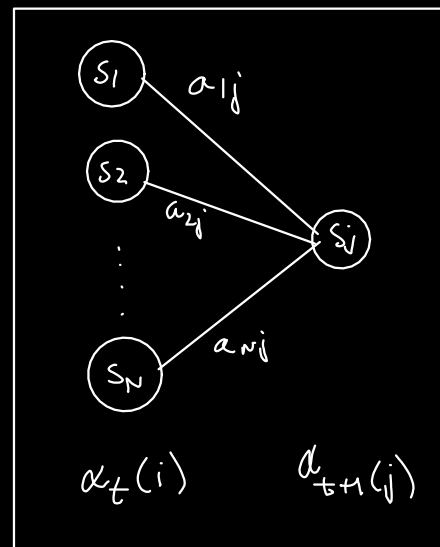
$$\alpha_1(i) = \pi_i b_i(o_1)$$

2. Induction

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) \cdot a_{ij} \right] b_j(o_{t+1})$$

3. Termination

$$P(O|\lambda) = \sum_{i=1}^N \alpha_T(i)$$



Backward Procedure

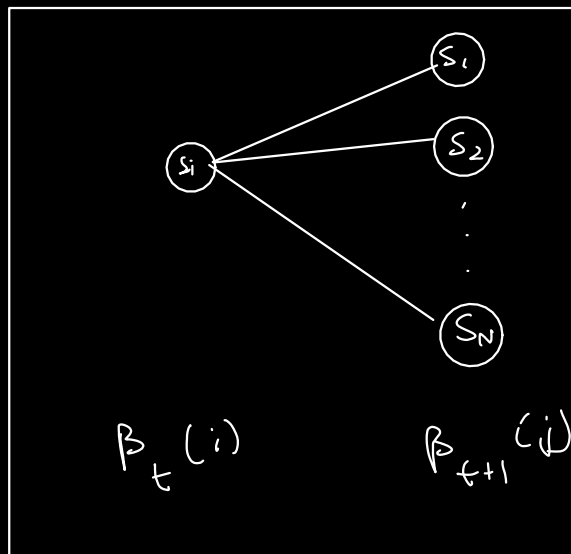
$$\beta_t(i) = P(o_{t+1} \dots o_T | q_t = s_i, \lambda)$$

1. Init

$$\beta_T(i) = 1$$

2. Induction

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)$$



$$O(N^2T)$$

2. Given observation sequence $O = \{o_1, \dots, o_T\}$
Get best $Q = \{q_1, \dots, q_T\}$

Solution:

Best state individually likely at a position i
Best state given all the previously observed states and observations

□ Viterbi Algorithm

Viterbi Algorithm

$$\delta_t(i) = \max_{q_1 \dots q_{t-1}} P(q_1 q_2 \dots q_t = i, o_1 o_2 \dots o_t | \lambda)$$

$$\delta_{t+1}(j) = \left[\max_i \delta_t(i) a_{ij} \right] b_j(o_{t+1})$$

Init

$$\delta_1(i) = \pi_i b_i(o_1)$$

$$\psi(i) = 0$$

Recursion

$$\delta_t(j) = \max \left[\delta_{t-1}(i) a_{ij} \right] b_j(o_t)$$

$$\psi_t(j) = \operatorname{argmax}_{1 \leq i \leq N} \left[\delta_{t-1}(i) a_{ij} \right]$$

Termination

$$p^* = \max_{1 \leq i \leq N} (\delta_T(i))$$

$$q_T^* = \operatorname{argmax}_{1 \leq i \leq N} (\delta_T(i))$$

Backtracking

$$q_t^* = \psi_{t+1}(q_{t+1}^*)$$

Tokenisation, LM, Smoothing

Tokenization

- Challenges:
- hyphenation
 - number
 - dates
 - abbreviations
 - punctuation
 - URLs
 - suffixification

$$\text{N-grams} \quad P(w | s(w)) = \frac{c(w, s(w) \dots)}{c(s(w) \dots)}$$

Language Model

A module that computes $P(w)$ or $P(w_n | w_1, w_2, \dots, w_{n-1})$

Perplexity and Entropy

Entropy over RV: X

$$H(X) = - \sum_{x=1}^n P(x) \log P(x)$$

Perplexity:

$$PP(w) = 2^{H(w)}$$

Perplexity = 2^{entropy} = $2^{\text{avg. \# of bits}}$
(of our model p)

avg. \# of words \rightarrow x can be any RV. In NLP, it's each word.

$$\begin{aligned} &= 2^{-\sum_{i=1}^N \frac{1}{N} P(x_i) \cdot \log_2 q(x_i)} \\ &\quad \text{read dist.} \quad \text{our prediction!} \\ &= e^{-\sum_{i=1}^N \frac{1}{N} P(x_i) \cdot \ln q(x_i)} \\ &\quad \text{We assume all words have the same frequency.} \\ &= e^{-\sum_{i=1}^N \frac{1}{N} \cdot \ln q(x_i)} \\ &= q(x_1)^{-\frac{1}{N}} \cdot q(x_2)^{-\frac{1}{N}} \dots q(x_N)^{-\frac{1}{N}} \\ &= \prod_{i=1}^N q(x_i)^{-\frac{1}{N}} = \sqrt[N]{\frac{1}{q(x_1) q(x_2) \dots q(x_N)}} \end{aligned}$$

Smoothing

Laplace: Add 1 to all counts.

Linear Interpolation:

$$\begin{aligned} \hat{P}(w_n | w_{n-1} w_{n-2}) &= \lambda_1 P(w_n | w_{n-1} w_{n-2}) \\ &\quad + \lambda_2 P(w_n | w_{n-1}) \\ &\quad + \lambda_3 P(w_n) \\ \sum_i \lambda_i &= 1 \end{aligned}$$

Weights from held-out corpus.

Katz

$$P_{katz}(w_n | w_{n-N+1}^{n-1}) = \begin{cases} P^*(w_n | w_{n-N+1}^{n-1}) & \text{if } C(w_{n-N+1}^{n-1}) > 0 \\ \alpha(w_{n-N+1}^{n-1}) P_{katz}(w_n | w_{n-N+2}^{n-1}) & \text{otherwise} \end{cases}$$

The w_i, w_{i-1}, w_{i-2} for clarity are referred as a sequence x, y, z .

Katz method incorporates discounting as integral part of the algorithm.

$$P_{katz}(z | x, y) = \begin{cases} P^*(z | x, y) & \text{if } C(x, y, z) > 0 \\ \alpha(x, y) P_{katz}(z | y) & \text{else if } C(x, y) > 0 \\ P^*(z) & \text{otherwise} \end{cases}$$

$$P_{katz}(z | y) = \begin{cases} P^*(z | y) & \text{if } C(y, z) > 0 \\ \alpha(y) P_{katz}(z) & \text{otherwise} \end{cases}$$

$$P^*(w_n | w_{n-N+1}^{n-1}) = \frac{c^*(w_{n-N+1}^{n-1})}{c(w_{n-N+1}^{n-1})}$$

Kneser-Ney

- Re-estimated counts c^* for greater than 1 counts could be estimated pretty well by just subtracting 0.75 from the MLE count c .
- Absolute discounting method formalizes this intuition by subtracting a fixed (absolute) discount d from each count.
 - The rational is that we have good estimates already for the high counts, and a small discount d won't affect them much.
 - The affected are only the smaller counts for which we do not necessarily trust the estimate anyhow.
- The equation for absolute discounting applied to bigrams (assuming a proper coefficient α on the backoff to make everything sum to one) is:

$$P_{absolute}(w_i | w_{i-1}) = \begin{cases} \frac{c(w_{i-1}w_i) - D}{c(w_{i-1})}, & \text{if } c(w_{i-1}w_i) > 0 \\ \alpha(w_i) P_{absolute}(w_i) & \text{otherwise} \end{cases}$$

- In practice distinct discount values d for the 0 and 1 counts are computed.
- Kneser-Ney discounting augments absolute discounting with a more sophisticated way to handle the backoff distribution. Consider the job of predicting the next word in the sentence, assuming we are backing off to a unigram model:
 - I can't see without my reading XXXXXX.
 - The word "glasses" seem much more likely to follow than the word "Francisco".
 - But "Francisco" is in fact more common, and thus a unigram model will prefer it to "glasses".
 - 1. Thus we would like to capture that although "Francisco" is frequent, it is only frequent after the word "San".
 - 2. The word "glasses" has a much wider distribution.

Thus the idea is instead of backing off to the unigram MLE count (the number of times the word w has been seen), we want to use a completely different backoff distribution!

- We want a heuristic that more accurately estimates the number of times we might expect to see word w in a new unseen context.
- The Kneser-Ney intuition is to base our estimate on the number of **different contexts word w has appeared in**.
- Words that have appeared in more contexts are more likely to appear in some new context as well.
- New backoff probability can be expressed as the "**continuation probability**" presented in following expression:

- Continuation Probability:

$$P_{\text{continuation}}(w_i) = \frac{|\{w_{i-1} : c(w_{i-1}w_i) > 0\}|}{\sum_{w_i} |\{w_{i-1} : c(w_{i-1}w_i) > 0\}|}$$

- Kneser-Ney backoff is formalized as follows assuming proper coefficient α on the backoff to make everything sum to one:

$$P_{KN}(w_i | w_{i-1}) = \begin{cases} \frac{c(w_{i-1}w_i) - D}{c(w_{i-1})}, & \text{if } c(w_{i-1}w_i) > 0 \\ \alpha(w_i) \frac{|\{w_{i-1} : c(w_{i-1}w_i) > 0\}|}{\sum_{w_i} |\{w_{i-1} : c(w_{i-1}w_i) > 0\}|} & \text{otherwise} \end{cases}$$

- Kneser-Ney **backoff** algorithm was shown to be less superior to its **interpolated** version. **Interpolated Kneser-Ney** discounting can be computed with an equation like the following (omitting the computation of β):

$$P_{KN}(w_i | w_{i-1}) = \frac{c(w_{i-1}w_i) - D}{c(w_{i-1})} + \beta(w_i) \frac{|\{w_{i-1} : c(w_{i-1}w_i) > 0\}|}{\sum_{w_i} |\{w_{i-1} : c(w_{i-1}w_i) > 0\}|}$$

- Practical note – it turns out that any interpolation model can be represented as a backoff model, and hence stored in ARPA backoff format. The interpolation is done when the model is built, thus the ‘bigram’ probability stored in the backoff format is really ‘bigram already interpolated with unigram’.

WITTEN BELL

Witten-Bell Smoothing

HYDERABAD

- Let's consider bi-grams
 - $T(w_{i-1})$ is the number of different words (types) that occur to the right of w_{i-1}
 - $N(w_{i-1})$ is the number of all word occurrences (tokens) to the right of w_{i-1}
 - $Z(w_{i-1})$ is the number of bigrams in the current data set starting with w_{i-1} that do not occur in the training data

- If $c(w_{i-1}, w_i) = 0$

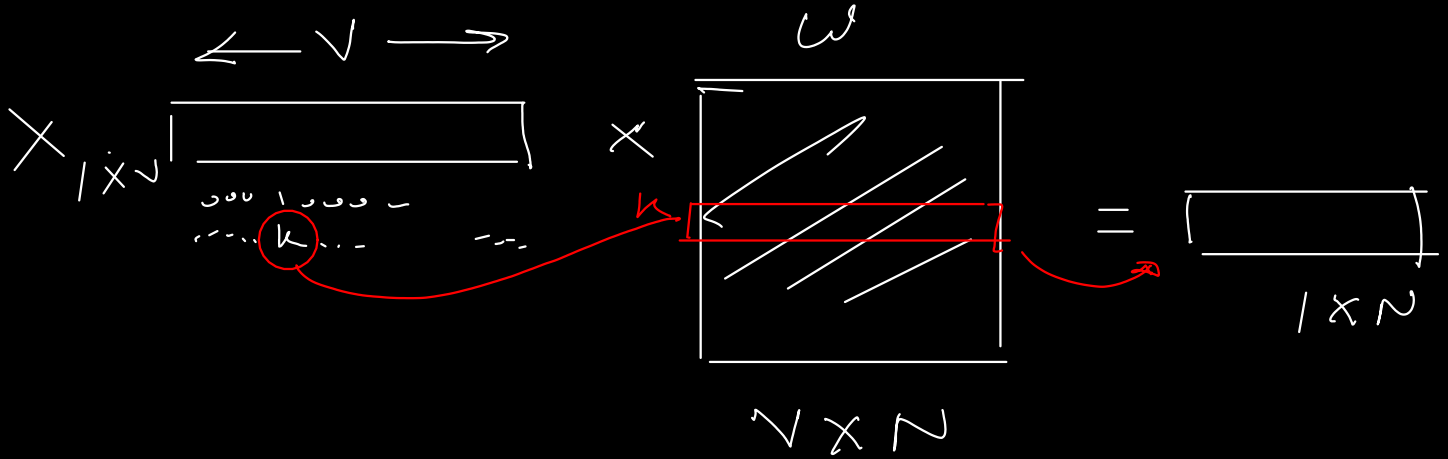
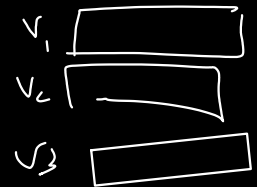
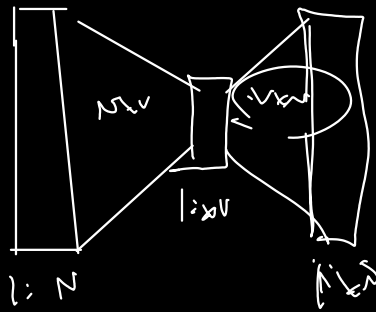
- If $c(w_{i-1}, w_i) > 0$

$$P^{WB}(w_i | w_{i-1}) = \frac{T(w_{i-1})}{Z(w_{i-1})(N + T(w_{i-1}))}$$

$$P^{WB}(w_i | w_{i-1}) = \frac{c(w_{i-1}w_i)}{N(w_{i-1}) + T(w_{i-1})}$$

Word2Vec

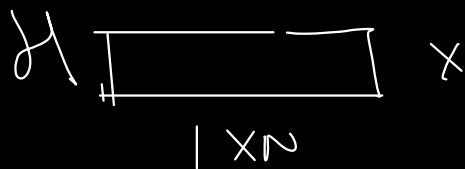
Input → hidden layer



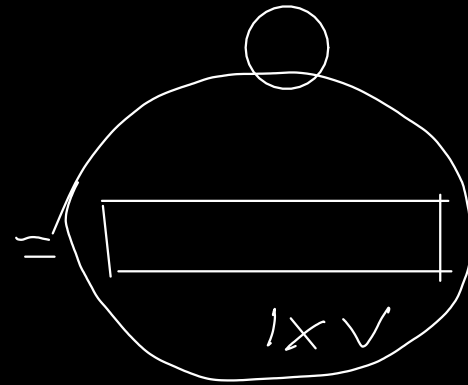
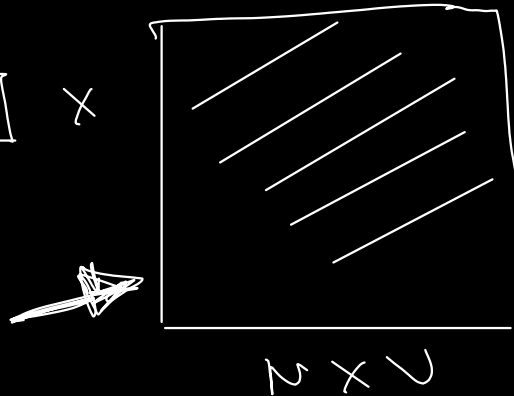
$$x_{1 \times V} \times W_{V \times N} = h_{1 \times N}$$

Hidden → Output layer

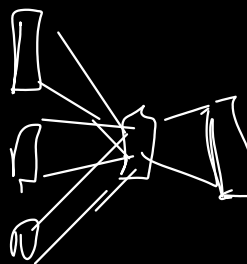
W'



Hidden layer



$$h = xw$$



$$u = \frac{1}{\sigma} (h_1 \dots h_N) w$$

GloVe

- uses global count statistics.
- uses co-occurrence matrix

We have P ratios,

As you can see, words that are related to the nature of "ice" and "steam" ("solid" and "gas" respectively) occur far more often with their corresponding words than the non-corresponding word. In contrast, words like "water" and "fashion" which are not particularly related to either have a probability ratio near 1. Note that the

$$\frac{P(k|ice)}{P(k|steam)}$$

$$k = \text{solid} : 8.9$$

$$k = \text{gas} : 8.8 \times 10^{-2}$$

$$k = \text{fashion, water} \sim 1$$

Building co-occurrence can be done during
so does not affect complexity.

vocabulary building,

General Word Embeddings

Matrix methods: LSA(SVD), HAP, CoMS, Hologger PCA

Neural approaches: NNLM, HLBL, RNN, Skip-gram (GloVe/VeB2)

Matrix

✓

Neural

Non compositionality

- stat
- fast train
- scales poor

- messy stat
- slow train
- scales with size,

Similarity Datasets

- similarity

- Manhattan sum,

MFN

3000 words

Rev

2032 " "

→ Disprop to word count

Evaluation

- NM
- Analogy
- IR
- Sense disambig.

Karakas

k1: karta

k2: karna

k3: karan (instrument)

k4: recipient

k6: source

k7: →

t → time

p → place

Thematic Roles

- Agent

sentient (ate → long)

- Instrument

(opened → key)

- Cause

(killed → epidemic)

- Experiencer

- Recipient

- Path

- Location

- Measure

- Theme

The boy plucked the flower

det n v det n
1 2 3 4 5

Shift-Reduce Parsing

Rules :

- $R_1 \quad S \rightarrow NP VP$
 $R_2 \quad NP \rightarrow \text{det } n$
 $R_3 \quad PP \rightarrow \text{prep } NP$
 $R_{4a} \quad VP \rightarrow v$
 $R_{4b} \quad VP \rightarrow v NP$
 $R_{4c} \quad VP \rightarrow v NP NP$
 $R_{4d} \quad VP \rightarrow VP PP$

Sl. No	stack	Input rem.	Action	Rules rem.
1	—	1	nil	—
1'	det	2	shift	R_1
1''	det n	3	shift	R_1
2	NP	3	reduce	R_3
2'	NP v	4	shift	R_1
3	NP VP	4	reduce	R_{4b}
4	S	4	reduce	R_2
4'	S det	5	shift	R_1
4''	S det n	null	shift	R_1

Sl. No

5 S NP null reduce R3

Error. ∴ Backtrack

3¹ NP VP det 5 shift R1

3¹¹ NP VP det n null shift R1

4 NP VP NP null reduce R3

Error. ∴ Backtrack

2¹¹ NP v det 5 shift R1

2¹¹¹ NP v det n null shift R1

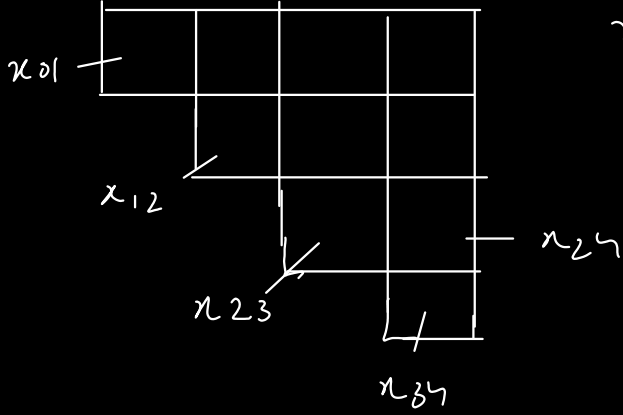
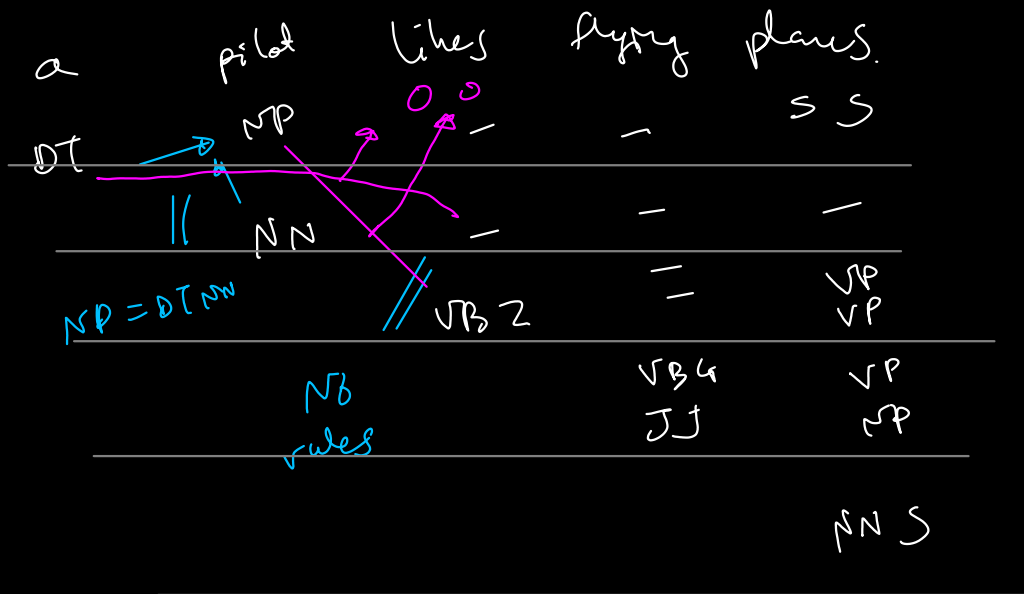
3 NP v NP null reduce R3

4 NP VP null reduce R4c

5 S null reduce R2

Accept

CYK Parsing

$$\begin{array}{c|c|c|c|c|c|c} \omega_1 & \omega_2 & \omega_3 & \dots & \omega_{n-1} & \omega_n \\ \hline 0 & 1 & 2 & 3 & & n \end{array}$$

$$x_{ij} : w_{i+1} \dots w_j$$


05 = 01 15
02 25
03 35
04 45

-- The mother goes - - -

mtl | null
senc | senc.

↳ [spec: det
root: mother
attr: [n: sing
genc: fem]]

goes: [root: go
attr: [T: present
A: simple
N: sing.]]