

Pron Verb z_3 - - -
 I like to eat - - -
 Pron Verb - - -
 She likes to work out - -
 z_1 z_2 - - -
 There is - - -

$$T_{10} = \left[\begin{array}{c} \cdot \end{array} \right] \rightarrow T_{10}(\text{ProN}) = \frac{\text{Count}(\text{ProN in 1st places})}{\# \text{ Sentences}}$$

Natural Language Processing CSE 325/425



Sihong Xie

$$A_0 = \left[\begin{array}{c} \cdot \end{array} \right]$$

$$T_1, A, B$$

$$A_0^{\text{ProN} \rightarrow \text{Verb}}$$

Lecture 20:

- Project 2 discussion
- Neural CYK

$$A \leftarrow \mu A_L + (1-\mu) A_0$$

$\nearrow \text{arg of } \text{em}()$
 \parallel
 A_0

$$= \frac{\text{Count}(\text{ProN} \rightarrow \text{Verb})}{\text{Count}(\text{ProN})}$$

$$= \frac{2}{2} = 1$$

Not-so-close-example

$$\begin{bmatrix} 3 & 2 \\ 3 & 6 \end{bmatrix} = 1 \times A_L = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} + 2 \times A_0 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Neural CYK

2018 ~ 2019

- Motivation: don't construct a CFG, but learn to predict a parse tree directly from a sentence.

$$P(T|S) = f(T; S)$$

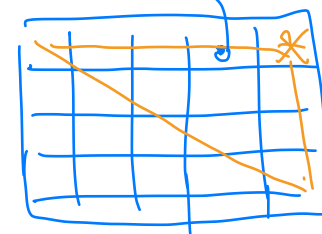
- Similar to CRF: no need to learn the A, B, and π matrices in HMM, but learn a parametric model to predict a POS-tag sequence from a sentence.

$$P(\alpha, o) = \frac{1}{Z} \exp \left\{ \sum_{t=1}^{T-1} \underbrace{\vec{w}^T \vec{f}(o_t, q_t, q_{t+1})}_{\text{can be replaced with NN output}} \right\}$$

- Discriminative
- Similar to RNN/LSTM: use neural networks to combine rich information in a data-driven way
 - No hand-designed features.

Neural CYK

$\delta_i(p, q) = \text{score for } (p, q) \text{ w/ } i$



- Recall: what we need to find the optimal parse tree in CYK.
 - the scores for each span (p, q) for each non-terminal N^i
 - the scores are computed based on a model (PCFG) and solutions to subproblems with a smaller span (p, r) and $(r+1, q)$

$$\psi_i(p, q) = \arg \max_{(j, k, r)} P(N^i \rightarrow N^j N^k) \delta_j(p, r) \delta_k(r+1, q)$$

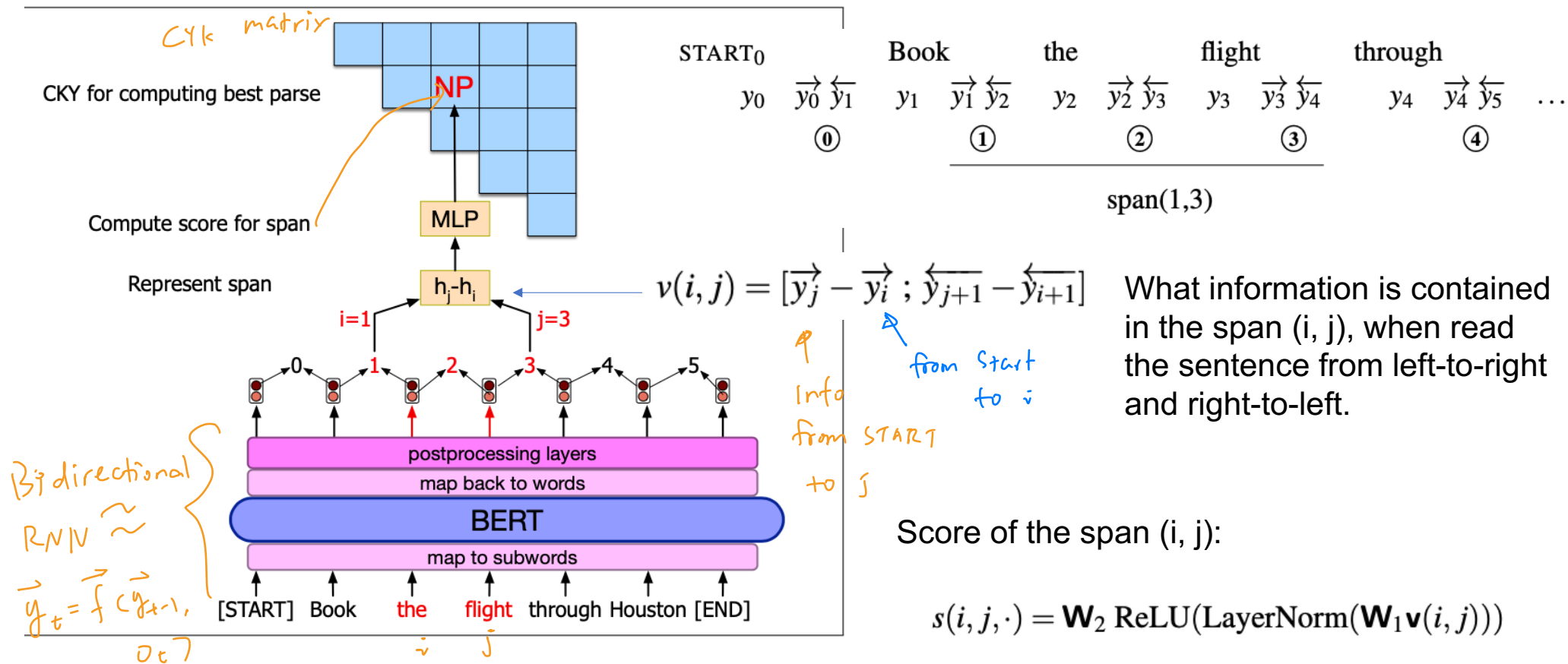
$$\delta_i(p, q) = \max_{\substack{1 \leq j, k \leq n \\ p \leq r < q}} P(N^i \rightarrow N^j N^k) \delta_j(p, r) \delta_k(r+1, q)$$

$$P(\hat{f}) = \delta_1(1, m)$$

- Replace the probability $\Pr(N^i \rightarrow N^j N^k)$ with some scores predicted by a Neural network.

Prob CYK inside prob: $\beta_{NP}(i, j)$

Neural CYK



Neural CYK

No CFG rules

- Finding the optimal tree
 - Assume that the scores of all spans are independent.
 - Related to the PCFG assumptions

A parse tree $T = \{(i_t, j_t, l_t) : t = 1, \dots, |T|\}$
 "clique" in CRF

A optimal tree $\hat{T} = \underset{T}{\operatorname{argmax}} s(T)$

Score of a tree $s(T) = \sum_{(i,j,l) \in T} s(i,j,l)$

$$\begin{aligned}
 & P \left(\begin{array}{c} {}^1S \\ \swarrow \quad \searrow \\ {}^2NP \quad {}^3VP \\ \swarrow \quad \searrow \quad | \\ the \quad man \quad snores \end{array} \right) \\
 &= P({}^1S_{13} \rightarrow {}^2NP_{12} \quad {}^3VP_{33}, {}^2NP_{12} \rightarrow the_1 \quad man_2, {}^3VP_{33} \rightarrow snores_3) \\
 &= P({}^1S_{13} \rightarrow {}^2NP_{12} \quad {}^3VP_{33}) P({}^2NP_{12} \rightarrow the_1 \quad man_2 | {}^1S_{13} \rightarrow {}^2NP_{12} \quad {}^3VP_{33}) \\
 &\quad P({}^3VP_{33} \rightarrow snores_3 | {}^1S_{13} \rightarrow {}^2NP_{12} \quad {}^3VP_{33}, {}^2NP_{12} \rightarrow the_1 \quad man_2) \\
 &= P({}^1S_{13} \rightarrow {}^2NP_{12} \quad {}^3VP_{33}) P({}^2NP_{12} \rightarrow the_1 \quad man_2) P({}^3VP_{33} \rightarrow snores_3) \\
 &= P(S \rightarrow NP \quad VP) P(NP \rightarrow the \quad man) P(VP \rightarrow snores)
 \end{aligned}$$

CYK using the scores:

on the diagonal

$$s_{\text{best}}(i, i+1) = \max_l s(i, i+1, l)$$

on other cells $s_{\text{best}}(i, j) = \max_l s(i, j, l)$

$$+ \max_k [s_{\text{best}}(i, k) + s_{\text{best}}(k, j)]$$

