

CS 533: Natural Language Processing

Transition-Based Dependency Parsing

Karl Stratos

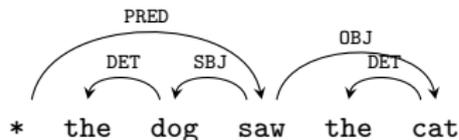


Rutgers University

Transition-Based Dependency Parsing

- ▶ Dependency parsing framed as a sequence of **transitions**

$$c_0 \xrightarrow{t_0} c_1 \xrightarrow{t_1} \dots \xrightarrow{t_{T-1}} c_T$$



- ▶ Runtime **linear** in sentence length!
 - ▶ Major advantage over graph-based dependency parsing

Overview

Dependency Parsing

Transition-Based Framework

- Configuration

- Transitions

Transition Systems

- Arc-Standard

- Arc-Eager

Implementation

- Training

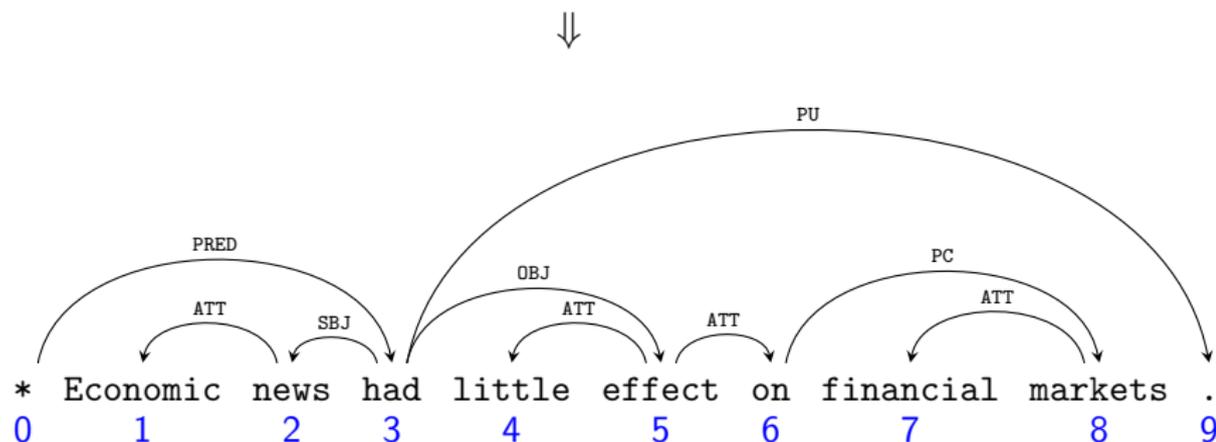
- Greedy Parser

- Beam Search Parser

- Evaluation

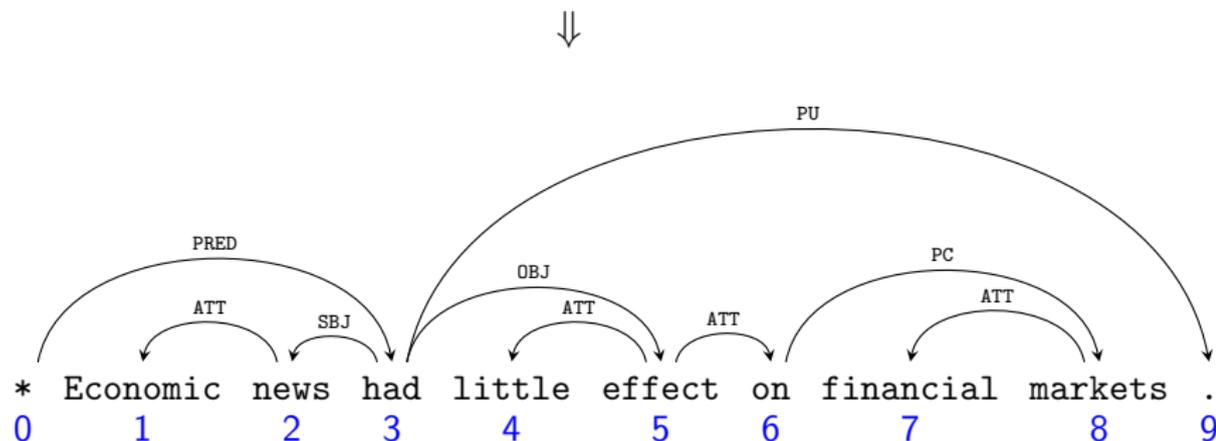
Example Dependency Tree (Nivre 2013)

Economic news had little effect on financial markets.



Example Dependency Tree (Nivre 2013)

Economic news had little effect on financial markets.



$$A = \{(0, \text{PRED}, 3), (3, \text{SBJ}, 2), (2, \text{ATT}, 1), (3, \text{OBJ}, 5), \\ (3, \text{PU}, 9), (5, \text{ATT}, 4), (5, \text{ATT}, 6), (6, \text{PC}, 8), (8, \text{ATT}, 7)\}$$

Dependency Parsing = Arc Finding

- ▶ Sentence: $x_1 \dots x_m$
- ▶ Associated nodes: $\mathcal{N} = \{0, 1, \dots, m\}$
 - ▶ Convention: leftmost root 0
- ▶ Labels: $L = \{\text{PRED}, \text{SBJ}, \dots\}$

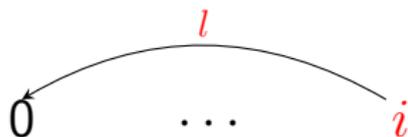
Goal. Find a set of **labeled, directed arcs**

$$A \subseteq \mathcal{N} \times L \times \mathcal{N}$$

that corresponds to a **correct dependency tree** for $x_1 \dots x_m$.

Valid Dependency Tree

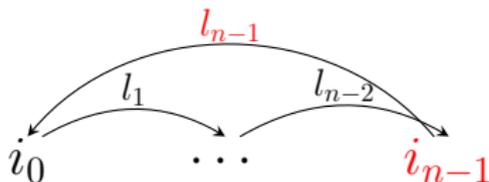
1. (Root): 0 must not have a parent.



2. (Connected): There must be a path from 0 to every $i \in \mathcal{N}$.
3. (Tree): A node must not have more than one parent.

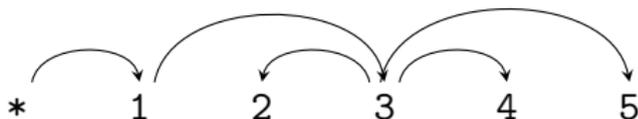


4. (Acyclic): Nodes must not form a cycle.

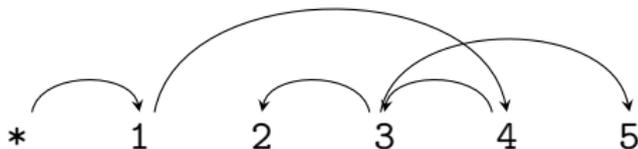


Projective

- ▶ A valid dependency tree is **projective** if for every arc (i, l, j) there is a path from i to k for all $i < k < j$.



- ▶ Valid but non-projective



We will focus on projective trees only!

Overview

Dependency Parsing

Transition-Based Framework

Configuration

Transitions

Transition Systems

Arc-Standard

Arc-Eager

Implementation

Training

Greedy Parser

Beam Search Parser

Evaluation

Parser Configuration

Triple $c = (\sigma, \beta, A)$ where

- ▶ $\sigma = [\dots i]$: “stack” of \mathcal{N} with i at the top
- ▶ $\beta = [i \dots]$: “buffer” of \mathcal{N} with i at the front
- ▶ $A \subseteq \mathcal{N} \times L \times \mathcal{N}$: arcs

Notation

- ▶ \mathcal{C} denotes the space of all possible configurations.
- ▶ $c.\sigma$, $c.\beta$, $c.A$ denote stack, buffer, arcs of $c \in \mathcal{C}$.

Configuration-Based Parsing Scheme

Initial configuration

$$c_0 := ([0], [1 \dots m], \{ \})$$

Apply “transitions” until we reach **terminal** c_T (defined later)

$$c_0 \xrightarrow{t_0} c_1 \xrightarrow{t_1} \dots \xrightarrow{t_{T-1}} c_T$$

and return as a parse

$$c_T.A$$

Overview

Dependency Parsing

Transition-Based Framework

Configuration

Transitions

Transition Systems

Arc-Standard

Arc-Eager

Implementation

Training

Greedy Parser

Beam Search Parser

Evaluation

Shift and Reduce

SHIFT $(\sigma, i|\beta, A) \Rightarrow (\sigma|i, \beta, A)$

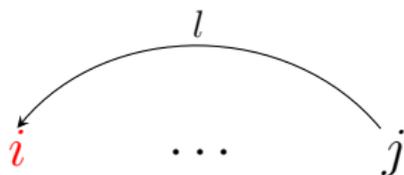
Illegal if β is empty.

REDUCE $(\sigma|i, \beta, A) \Rightarrow (\sigma, \beta, A)$

Illegal if i does not have a parent.

Left-Arc

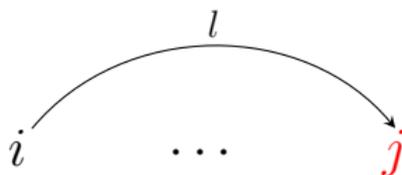
$$\mathbf{LEFT}_l (\sigma|i|j, \beta, A) \Rightarrow (\sigma|j, \beta, A \cup \{(j, l, i)\})$$



Illegal if either $i = 0$ or i already has a parent.

Right-Arc

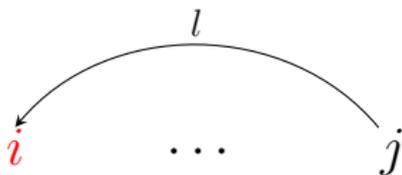
RIGHT_{*l*} $(\sigma|i|j, \beta, A) \Rightarrow (\sigma|i, \beta, A \cup \{(i, l, j)\})$



Illegal if j already has a parent.

“Eager” Left-Arc

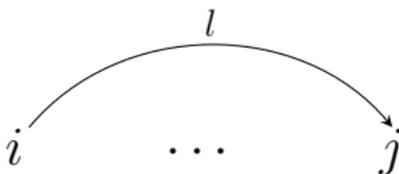
$$\mathbf{LEFT}_l^e \quad (\sigma|i, j|\beta, A) \Rightarrow (\sigma, j|\beta, A \cup \{(j, l, i)\})$$



Illegal if either $i = 0$ or i already has a parent.

“Eager” Right-Arc

$$\mathbf{RIGHT}_l^e (\sigma|i, j|\beta, A) \Rightarrow (\sigma|i|j, \beta, A \cup \{(i, l, j)\})$$



Illegal if j already has a parent.

Legal transitions

- ▶ Certain transitions are illegal depending on $c \in \mathcal{C}$.
- ▶ We will denote the **set of legal actions at** c by $\text{LEGAL}(c)$.

Overview

Dependency Parsing

Transition-Based Framework

- Configuration

- Transitions

Transition Systems

- Arc-Standard

- Arc-Eager

Implementation

- Training

- Greedy Parser

- Beam Search Parser

- Evaluation

Definition

$2|L| + 1$ possible transitions \mathcal{T}^{std}

▶ **SHIFT**: $(\sigma, i|\beta, A) \Rightarrow (\sigma|i, \beta, A)$

▶ **LEFT_l** for each $l \in L$:

$$(\sigma|i|j, \beta, A) \Rightarrow (\sigma|j, \beta, A \cup \{(j, l, i)\})$$

▶ **RIGHT_l** for each $l \in L$:

$$(\sigma|i|j, \beta, A) \Rightarrow (\sigma|i, \beta, A \cup \{(i, l, j)\})$$

Terminal condition: $c.\sigma = [0]$ and $c.\beta = []$

Properties

- ▶ Makes **exactly** $2m$ transitions to parse $x_1 \dots x_m$. Why?
- ▶ **Bottom-up**: a node must collect all its children before getting a parent. Why?
- ▶ **Sound**: if c is terminal, $c.A$ forms a valid projective tree.
- ▶ **Complete**: every valid projective tree A can be produced from c_0 by some sequence of transitions $t_0 \dots t_{T-1} \in \mathcal{T}^{\text{std}}$.

$$t_i = \text{Oracle}^{\text{std}}(c_i)$$
$$c_{i+1} = t_i(c_i)$$

Input: gold arcs A^{gold} , non-terminal configuration $c = (\sigma, \beta, A)$

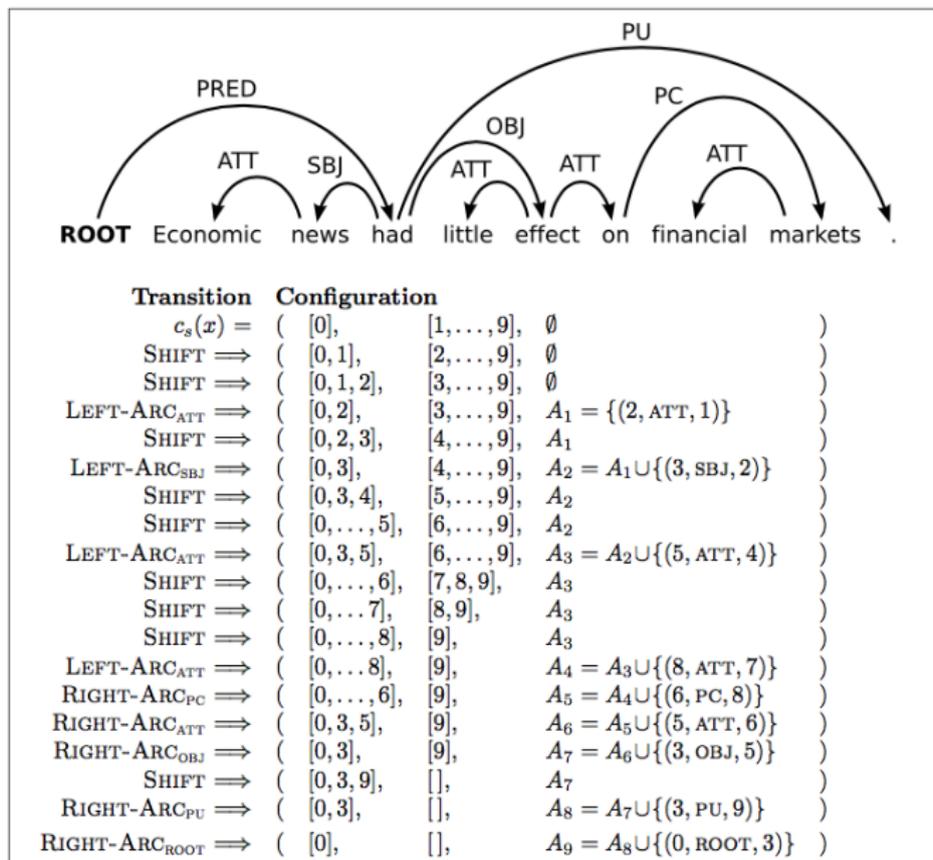
Output: transition $t \in \mathcal{T}^{\text{std}}$ to apply on c

1. Return **SHIFT** if $|\sigma| = 1$.
2. Otherwise $\sigma = [\dots i j]$ for some $i < j$:
 - 2.1 Return **LEFT** _{l} if $(j, l, i) \in A^{\text{gold}}$.
 - 2.2 Return **RIGHT** _{l} if $(i, l, j) \in A^{\text{gold}}$ and for all $l' \in L, j' \in \mathcal{N}$,

$$(j, l', j') \in A^{\text{gold}} \quad \Rightarrow \quad (j, l', j') \in A$$

- 2.3 Return **SHIFT** otherwise.

Example Parse (Nivre 2013)



Overview

Dependency Parsing

Transition-Based Framework

Configuration

Transitions

Transition Systems

Arc-Standard

Arc-Eager

Implementation

Training

Greedy Parser

Beam Search Parser

Evaluation

Definition

$2|L| + 2$ possible transitions \mathcal{T}^{eag}

SHIFT: $(\sigma, i|\beta, A) \Rightarrow (\sigma|i, \beta, A)$

REDUCE: $(\sigma|i, \beta, A) \Rightarrow (\sigma, \beta, A)$

LEFT _{l} ^{e} for each $l \in L$:

$$(\sigma|i, j|\beta, A) \Rightarrow (\sigma, j|\beta, A \cup \{(j, l, i)\})$$

RIGHT _{l} ^{e} for each $l \in L$:

$$(\sigma|i, j|\beta, A) \Rightarrow (\sigma|i|j, \beta, A \cup \{(i, l, j)\})$$

Terminal condition: $c.\beta = []$

- ▶ Stop as soon as the buffer is empty.

Properties

- ▶ Makes **at most** $2m$ transitions to parse $x_1 \dots x_m$. Why?
- ▶ **Partially top-down**: but a node must collect all its **left** children before right children. Why?
- ▶ **Not sound**: even if c is terminal, $c.A$ may form unconnected projective trees (“dependency forest”).
 - ▶ But can be manually corrected by connecting to the root.



- ▶ **Complete**: every valid projective tree A can be produced from c_0 by some sequence of transitions $t_0 \dots t_{T-1} \in \mathcal{T}^{\text{eag}}$.

$$t_i = \text{Oracle}^{\text{eag}}(c_i)$$

$$c_{i+1} = t_i(c_i)$$

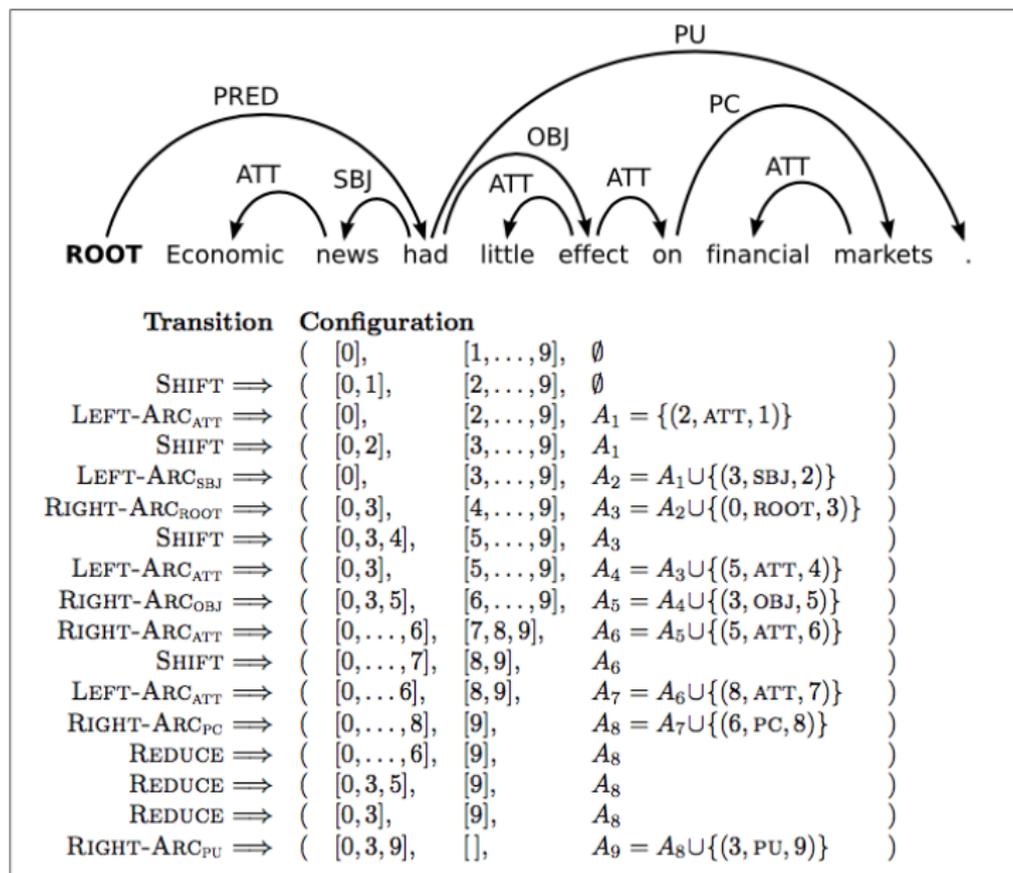
Input: gold arcs A^{gold} , non-terminal configuration

$$c = (\sigma = [\dots i], \beta = [j \dots], A)$$

Output: transition $t \in \mathcal{T}^{\text{eag}}$ to apply on c

1. Return **LEFT** _{l} ^{e} if $(j, l, i) \in A^{\text{gold}}$.
2. Return **RIGHT** _{l} ^{e} if $(i, l, j) \in A^{\text{gold}}$.
3. Return **REDUCE** if there is some $k < i$ such that $(k, l, j) \in A^{\text{gold}}$ or $(j, l, k) \in A^{\text{gold}}$ for some l .
4. Return **SHIFT** otherwise.

Example Parse (Nivre 2013)



Overview

Dependency Parsing

Transition-Based Framework

- Configuration

- Transitions

Transition Systems

- Arc-Standard

- Arc-Eager

Implementation

- Training

- Greedy Parser

- Beam Search Parser

- Evaluation

Getting Training Data

- ▶ **Treebank**: sentence-tree pairs $(x^{(1)}, A^{(1)}) \dots (x^{(M)}, A^{(M)})$
 - ▶ Assume all projective

- ▶ For each $A^{(j)}$, use an oracle to extract

$$(c_0^{(j)}, t_0^{(j)}) \dots (c_{T-1}^{(j)}, t_{T-1}^{(j)})$$

where $t_{T-1}^{(j)}(c_{T-1}^{(j)}) \cdot A = A^{(j)}$.

- ▶ We can now use this to train a **classifier**

$$(x^{(j)}, c_i^{(j)}) \mapsto t_i^{(j)}$$

Linear Classifier

- ▶ Parameters: $w_t \in \mathbb{R}^d$ for each $t \in \mathcal{T}$
- ▶ Each $c \in \mathcal{C}$ for sentence x is “featurized” as $\phi^x(c) \in \mathbb{R}^d$.
 - ▶ Classical approach: **binary features** providing useful signals
 - ▶ Assumes we have access to POS tags of $x_1 \dots x_m$.

$$\phi_{20134}^x(c) := \begin{cases} 1 & \text{if } x_{c.\sigma[0]}.POS = NN \text{ and } x_{c.\beta[0]}.POS = VBD \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_{1988}^x(c) := \begin{cases} 1 & \text{if } x_{c.\sigma[0]}.POS = VBD \text{ with leftmost arc SUBJ} \\ 0 & \text{otherwise} \end{cases}$$

$$\phi_{42}^x(c) := \begin{cases} 1 & \text{if } x_{c.\beta[1]} = \text{cat} \\ 0 & \text{otherwise} \end{cases}$$

Linear Classifier (Continued)

- ▶ **Score** of $t \in \mathcal{T}$ at $c \in \mathcal{C}$ for x :

$$\begin{aligned}\text{score}_x(t|c) &:= w_t \cdot \phi^x(c) \\ &= \sum_{i=1: \phi_i^x(c)=1}^d [w_t]_i\end{aligned}$$

- ▶ From here on, we assume $\{w_t\}_{t \in \mathcal{T}}$ trained from data.

Important Aside

Each c_i is computed from **past decisions** $t_0 \dots t_{i-1}$.

$$c_i = t_{i-1}(t_{i-2}(\dots t_0(c_0)))$$

So the score function on c_i is really a **function of** $t_0 \dots t_{i-1}$.

$$\text{score}_x(t|c) = \text{score}_x(t|t_1 \dots t_{i-1})$$

Will use c_i and $t_0 \dots t_{i-1}$ interchangeably.

Overview

Dependency Parsing

Transition-Based Framework

- Configuration

- Transitions

Transition Systems

- Arc-Standard

- Arc-Eager

Implementation

- Training

- Greedy Parser

- Beam Search Parser

- Evaluation

Greedy

At each configuration c_i , pick

$$t_i \leftarrow \arg \max_{t \in \text{LEGAL}(c_i)} \text{score}_x(t | t_0 \dots t_{i-1})$$

Parsing Algorithm

Input: $\{w_t\}_{t \in \mathcal{T}}$, sentence x of length m

Output: arcs representing a dependency tree for x

1. $c \leftarrow c_0$
2. While $c.\beta \neq []$,

2.1 Select

$$\hat{t} \leftarrow \arg \max_{t \in \text{LEGAL}(c)} \text{score}_x(t|c)$$

2.2 Make a transition: $c \leftarrow \hat{t}(c)$.

3. Return $c.A$.

Overview

Dependency Parsing

Transition-Based Framework

- Configuration

- Transitions

Transition Systems

- Arc-Standard

- Arc-Eager

Implementation

- Training

- Greedy Parser

- Beam Search Parser

- Evaluation

Beam Search

Approximate the **optimal sequence of transitions**:

$$t_0^* \dots t_{T-1}^* = \arg \max_{\substack{t_0 \dots t_{T-1}: \\ t_i \in \text{LEGAL}(c_i) \\ c_T \cdot \beta = []}} \sum_{i=0}^{T-1} \text{score}_x(t_i | t_0 \dots t_{i-1})$$

Parsing Algorithm

Input: $\{w_t\}_{t \in \mathcal{T}}$, sentence x of length m , beam width K

Beam: $\langle c, s \rangle \in \mathcal{C} \times \mathbb{R}$ organized by second argument (score)

Output: arcs representing a dependency tree for x

1. $\mathcal{B} \leftarrow \text{Beam}(\{\langle c_0, 0 \rangle\}, K)$
2. While $c.\beta \neq []$ for some $\langle c, s \rangle \in \mathcal{B}$,
 - 2.1 $\mathcal{B}' \leftarrow \text{Beam}(\{\}, K)$
 - 2.2 For $\langle c, s \rangle \in \mathcal{B}$, for $t \in \text{LEGAL}(c)$,

$$\mathcal{B}'.\text{push}\langle t(c), s + \text{score}_x(t|c) \rangle$$

- 2.3 $\mathcal{B} \leftarrow \mathcal{B}'$
3. Return $c^*.A$ where $c^* \leftarrow \mathcal{B}.\text{pop}()$.

Overview

Dependency Parsing

Transition-Based Framework

Configuration

Transitions

Transition Systems

Arc-Standard

Arc-Eager

Implementation

Training

Greedy Parser

Beam Search Parser

Evaluation

- ▶ Unlabeled Attachment Score (UAS):

$$\frac{\# \text{ words with correct parent}}{\# \text{ words}}$$

- ▶ Labeled Attachment Score (LAS):

$$\frac{\# \text{ words with correct parent and label}}{\# \text{ words}}$$

Current state-of-the-art: 93-95 UAS, 91-93 LAS!

Parting Remarks

- ▶ There are better ways to train model $\{w_t\}_{t \in \mathcal{T}}$.
 - ▶ Online learning, “dynamic” oracles, etc.
- ▶ Today, state-of-the-art parsers are obtained by just replacing

$$\text{score}_x(t|c) = \overbrace{w_t}^{\text{linear}} \cdot \underbrace{\phi^x(c)}_{\text{hand-engineered}}$$

with a neural network (Kiperwasser and Goldberg, 2016).

- ▶ **Graph-based** dependency parsing (Eisner, 1996): similar to CKY in constituency parsing