# Lecture 6: Transforming Dependency to Context-free Grammars

Jelle Zuidema
ILLC, Universiteit van Amsterdam

Unsupervised Language Learning, 2014

# Plan for today
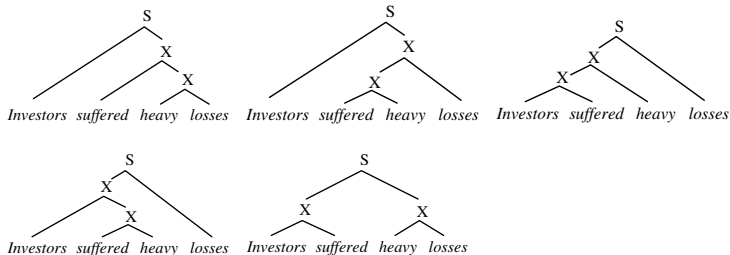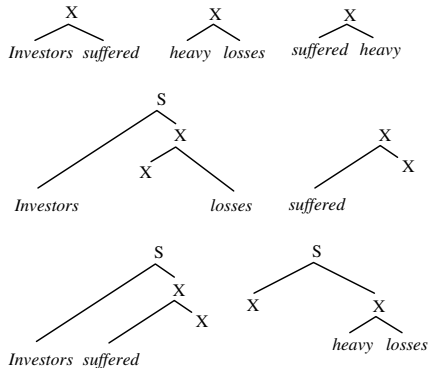
(based on slides from Mark Johnson)

# How Does U-DOP Operate?

**1. Assign *all* possible binary trees to strings where each root node is labeled *S* and other nodes labeled *X,* and store them in a parse forest**

E.g., for WSJ sentence *Investors suffered heavy losses*:

**2. Convert the set of all trees into all subtrees. For instance:**



=> Note that some subtrees contain discontiguous yields

**3. Compute most probable tree among shortest derivations for new string (as in DOP):**

Probability of…

    a *subtree t* :

$$P(t) \; = \; \frac{|t|}{\sum_{t' \,:\; root(t')=root(t)} |t'|}$$

    a *derivation*  $d = t_1 \circ ... \circ t_n$ :

$$P(t_1 \circ ... \circ t_n) \; = \; \prod_i P(t_i)$$

    a *parse tree*  $T$ :

$$P(T) \; = \; \sum_d \prod_i P(t_{id})$$

## U-DOP compared to other models on WSJ-10

(using 7422 sentences up to **10** words, as in Klein and Manning 2004)

| Model | F-score on WSJ-10 |
|---|---|
| CCM | 71.9 |
| DMV | 52.1 |
| DMV+CCM | 77.6 |
| U-DOP | **82.7** |
| U-DOP without discontiguous subtrees | **72.1** |

CCM:   Klein and Manning (2002) based on all linear
       (contiguous) contexts *without* holes
DMV:   Klein and Manning (2004) using
       dependency structures
U-DOP: equivalent to CCM *plus* discontiguous contexts *with*
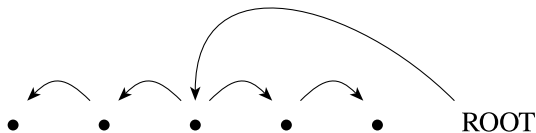       holes: *11% improvement in F-score*

# Shortcomings of U-DOP

- Viewed from the statistical inference perspective, the model relies much on heuristics: initialization, training & stopping

- Results with UML-DOP (Bod,06) suggests it is approximately Maximum Likelihood...

- ... but not over the entire PTSG space, as there are exponentially many subtrees, and exponentially many trees for a sentence!

- Implementation must somehow restrict space; efficiency remains the achilles heel.

Introduction
Constituent Context Model
**Dependency Model**

**Dependencies**
Dependency Model with Valence
Combined Model
Evaluation

## Definitions

Definitions:

- a dependency $d$ is a pair $\langle h, a \rangle$, where $h$ is the head and $a$ the argument, both are words in sentence $s$;
- a dependency structure $D$ is set of dependencies which form a planar, acyclic graph rooted in ROOT;
- the skeleton $G$ of a dependency structure specifies the arrows, but not the words; $G$ and $s$ together fully determine the dependency structure.

Introduction
Constituent Context Model
**Dependency Model**

Dependencies
**Dependency Model with Valence**
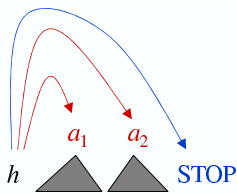Combined Model
Evaluation

## Dependency Model with Valence

Klein and Manning (2004) propose a model that generates dependencies outwards from the head:

- generate a set of arguments on one side of the head, then a STOP argument to terminate;
- the do the same thing on the other side;
- terminate with probability $P_{STOP}$, if not STOP then choose another argument with probability $P_{CHOOSE}$.

# Extended Chomsky Hierarchy

| language | grammar | rules |
|---|---|---|
| $\{a, b, cbabb\}$ | Set | $\in$ |
| $(ab)^n$ | ngram | $\langle a, b \rangle, \langle b, a \rangle, \langle ab, a \rangle$ |
| $a^n ba^m$ | Left-linear | $S \to AB, B \to bA$ |
| $a^n b^n$ | Context-free | $S \to aSb, S \to ab$ |
| $a^n b^n c^n d^n \vert 1 \leq n$ | Range Concatenation | $S[abc] \to A[a, c]B[b]$ |
| | Unrestricted | |

# Probabilistic Extensions

| grammar | probabilistic grammar |
|---------|----------------------|
| Set | Probability distribution |
| ngram | Markov model |
| Left-linear | Hidden Markov (HMM) |
| Context-free | PCFG |
| Range Concatenation | PLCRS |
| Unrestricted | |

# CFG encoding of Dependency Grammars

- Given that dependency grammars must be somewhere on the
  CH, presumable below contextfree, can we reuse the
  technology we developed for context-free grammars (rule
  extraction, CYK, Inside algorithm, Inside-outside) for
  dependency grammars?
- Yes!- at least for some kind of dependency grammars and
  given the right preprocessing.

# Plan for today

## Projective Bilexical Dependency Grammars

► Projective Bilexical Dependency Grammar (PBDG)

$$0 \overset{\frown}{\ } \text{gave} \qquad \text{Sandy} \overset{\frown}{\ } \text{gave}$$
$$\text{gave} \overset{\frown}{\ } \text{dog} \qquad \text{the} \overset{\frown}{\ } \text{dog}$$
$$\text{gave} \overset{\frown}{\ } \text{bone} \qquad \text{a} \overset{\frown}{\ } \text{bone}$$

► A dependency parse generated by the PBDG

0  Sandy  gave  the  dog  a  bone

► Weights can be attached to dependencies (and preserved in CFG transforms)
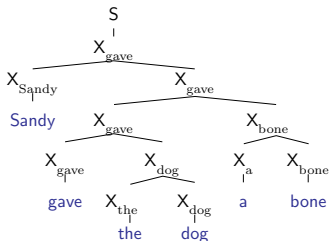
## A naive encoding of PBDGs as CFGs

$$
\begin{aligned}
S &\rightarrow X_u && \text{where } 0 \overset{\frown}{\phantom{}} u \\
X_u &\rightarrow u && \\
X_u &\rightarrow X_v X_u && \text{where } v \overset{\frown}{\phantom{}} u \\
X_u &\rightarrow X_u X_v && \text{where } u \overset{\frown}{\phantom{}} v
\end{aligned}
$$

## Spurious ambiguity in naive encoding

▶ Naive encoding allows dependencies on different sides of head to be freely reordered

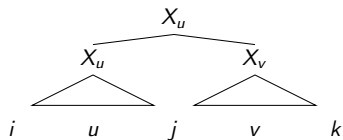⇒ Spurious ambiguity in CFG parses (not present in PBDG parses)

# Parsing naive CFG encoding takes $O(n^5)$ time

► A production schema such as

$$X_u \rightarrow X_u X_v$$

has 5 variables, and so can match input in $O(n^5)$ different ways
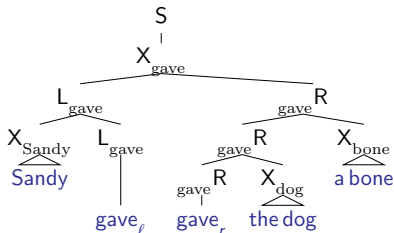
# Plan for today

## Simple split-head encoding

▶ Replace input word $u$ with a *left variant* $u_\ell$ and a *right variant* $u_r$
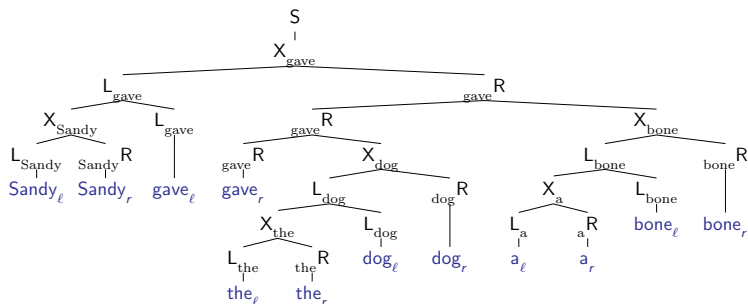  (can be avoided in practice with fancy book-keeping)

Sandy gave the dog a bone
$$\Downarrow$$
$\text{Sandy}_\ell \ \text{Sandy}_r \ \text{gave}_\ell \ \text{gave}_r \ \text{the}_\ell \ \text{the}_r \ \text{dog}_\ell \ \text{dog}_r \ a_\ell \ a_r \ \text{bone}_\ell \ \text{bone}_r$

▶ PCFG separately collects left dependencies and right dependencies


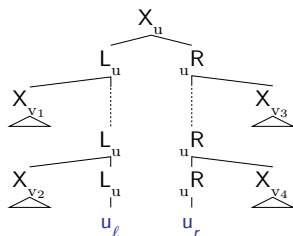
$$
\begin{aligned}
S &\rightarrow X_u & \text{where } 0 \curvearrowright u \\
X_u &\rightarrow L_u \ {}_u R & \text{where } u \in \Sigma \\
L_u &\rightarrow u_l \\
L_u &\rightarrow X_v \ L_u & \text{where } v \curvearrowleft u \\
{}_u R &\rightarrow u_r \\
{}_u R &\rightarrow {}_u R \ X_v & \text{where } u \curvearrowright v
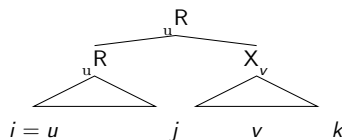\end{aligned}
$$

9 / 22

## Simple split-head CFG parse

# $L_u$ and $_uR$ heads are phrase-peripheral $\Rightarrow O(n^4)$

▸ Heads of $L_u$ and $_uR$ are always at right (left) edge



$$
\begin{aligned}
S &\rightarrow X_u && \text{where } 0 \overset{\frown}{\rightarrow} u \\
X_u &\rightarrow L_u \; _uR && \text{where } u \in \Sigma \\
L_u &\rightarrow u_l \\
L_u &\rightarrow X_v \; L_u && \text{where } v \overset{\frown}{\leftarrow} u \\
_uR &\rightarrow u_r \\
_uR &\rightarrow \; _uR \; X_v && \text{where } u \overset{\frown}{\rightarrow} v
\end{aligned}
$$

▸ $X_u \rightarrow L_u \; _uR$ take $O(n^3)$
▸ $_uR \rightarrow \; _uR \; X_v$ take $O(n^4)$

# Plan for today

## The Unfold-Fold transform

- ▶ Unfold-fold originally proposed for transforming recursive programs; used here to transform CFGs into new CFGs

- ▶ *Unfolding* a nonterminal replaces it with its expansion

$$
\begin{array}{ll}
A \to \alpha \, B \, \gamma & \begin{array}{l} A \to \alpha \, \beta_1 \, \gamma \\ A \to \alpha \, \beta_2 \, \gamma \end{array} \\
B \to \beta_1 & \\
B \to \beta_2 \qquad \Rightarrow & B \to \beta_1 \\
\qquad \cdots & B \to \beta_2 \\
& \qquad \cdots
\end{array}
$$

- ▶ *Folding* is the inverse of unfolding (replace RHS with nonterminal)

$$
\begin{array}{ll}
A \to \alpha \, \beta \, \gamma & A \to \alpha \, B \, \gamma \\
B \to \beta \qquad \Rightarrow & B \to \beta \\
\qquad \cdots & \qquad \cdots
\end{array}
$$

- ▶ Transformed grammar generates same language (Sato 1992)

# Unfold-fold converts $O(n^4)$ to $O(n^3)$ grammar

▶ Unfold $X_v$ responsible for $O(n^4)$ parse time

$$
\begin{array}{lll}
L_u & \to & u_l \\
L_u & \to & X_v \, L_u \\
X_v & \to & L_v \, {_v}R
\end{array}
\quad \Rightarrow \quad
\begin{array}{lll}
L_u & \to & u_l \\
L_u & \to & L_v \, {_v}R \, L_u
\end{array}
$$

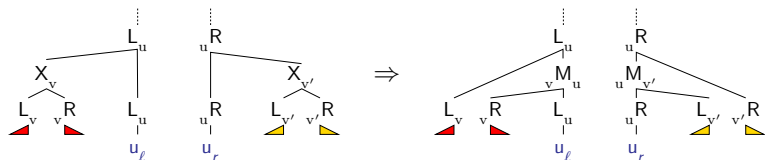▶ Introduce new non-terminals ${_x}M_y$ (doesn't change language)

$$
{_x}M_y \; \to \; {_x}R \, L_y
$$

▶ Fold two children of $L_u$ into ${_x}M_y$

$$
\begin{array}{lll}
L_u & \to & u_l \\
L_u & \to & L_v \, {_v}R \, L_u \\
{_x}M_y & \to & {_x}R \, L_y
\end{array}
\quad \Rightarrow \quad
\begin{array}{lll}
L_u & \to & u_l \\
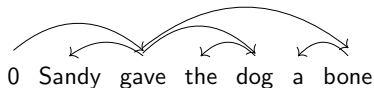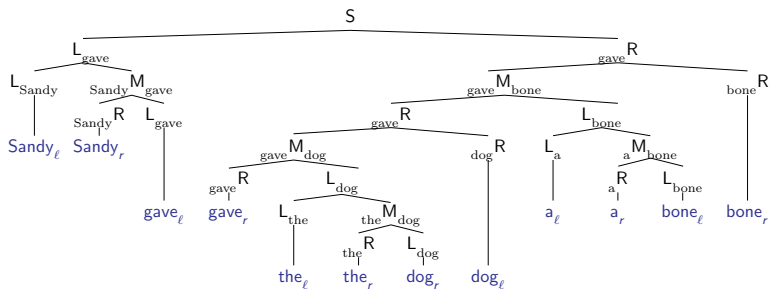L_u & \to & L_v \, {_v}M_u \\
{_x}M_y & \to & {_x}R \, L_y
\end{array}
$$

# Transformed grammar collects left and right dependencies separately



- $X_v$ constituents (which cause $O(n^4)$ parse time) no longer used
- Head annotations now all phrase peripheral $\Rightarrow O(n^3)$ parse time
- Dependencies can be recovered from parse tree
- Basically same as Eisner and Satta $O(n^3)$ algorithm
  - explains why Inside-Outside sanity check fails for Eisner/Satta
  - two copies of each terminal $\Rightarrow$ each terminals' Outside probability is *double* the Inside sentence probability

# Parse using $O(n^3)$ transformed split-head grammar

## Parsing time of CFG encodings of same PBDG

| CFG schemata | sentences parsed / second |
|---|---|
| Naive $O(n^5)$ CFG | 45.4 |
| $O(n^4)$ simple split-head CFG | 406.2 |
| $O(n^3)$ transformed split-head CFG | 3580.0 |

- ▶ Weighted PBDG; all pairs of heads have some dependency weight
- ▶ Dependency weights precomputed before parsing begins
- ▶ Timing results on a 3.6GHz Pentium 4 machine parsing section 24 of the PTB
- ▶ CKY parsers with grammars hard-coded in C (no rule lookup)
- ▶ Dependency accuracy of Viterbi parses = 0.8918 for all grammars
- ▶ *Feature extraction is much slower than even naive CFG*