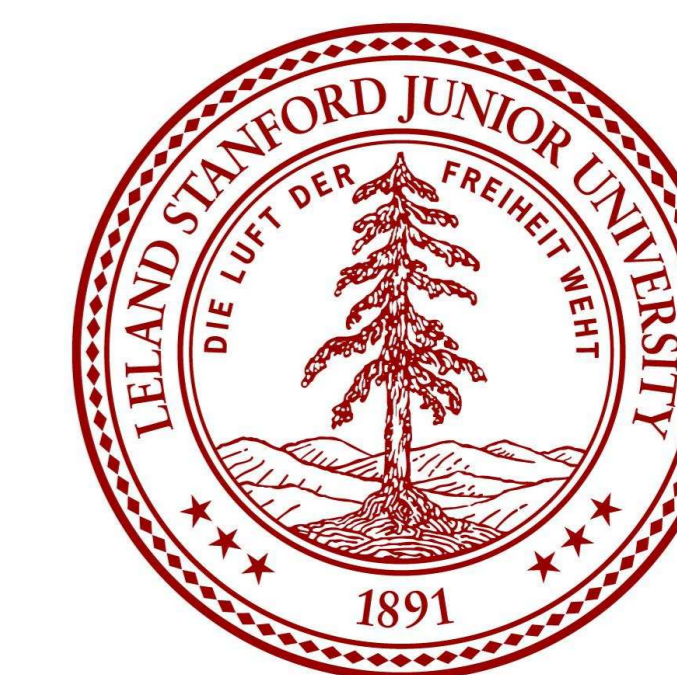




Three Dependency-and-Boundary Models for Grammar Induction

VALENTIN I. SPITKOVSKY, HIYAN ALSHAWI AND DANIEL JURAFSKY

vals@stanford.edu, hiyan@google.com and jurafsky@stanford.edu



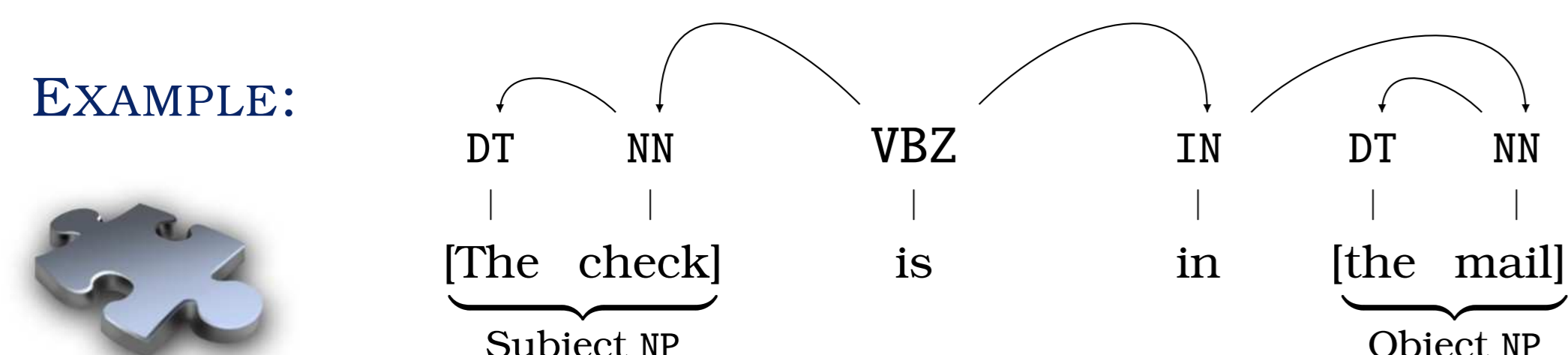
IDEA

Use boundary cues in head-driven dependency grammars.

INTUITION

Induce structure by working inwards from edges.

EXAMPLE:



- learn from left fringe (determiner DT) how to parse object NP
- based on right fringe (noun NN), correctly parse subject NP
- between the two, glean make-up of larger phrases (e.g., VP)

MOTIVATION

MACHINE LEARNING: FOCUS ON OBSERVABLE, COMPLEMENTARY FEATURES

- weak equivalence of phrase representations (Xia and Palmer, 2001)
- redundant views of data ease learning (Blum and Mitchell, 1998)

GRAMMAR INDUCTION: A BOUNTY OF CONSTITUENT BOUNDARY MARKERS

- at sentence beginnings and ends (Hänig et al., 2008; Hänig, 2010)
- around function words (Berant et al., 2006)
- around punctuation marks (Seginer, 2007; Ponvert et al., 2010; 2011; Spitkovsky et al., 2011)
- at capitalization (or script) change-points (Spitkovsky et al., 2012)
- at web markup bracketing end-points (Spitkovsky et al., 2010)
- around other semantic annotations (Naseem and Barzilay, 2011)

LANGUAGE ACQUISITION: WORD BOUNDARIES MATTER TO HUMAN LEARNING

- importance of exposure to isolated words (Brent and Siskind, 2001)

DBM-1

Stop generation based on fringe words of partial yields.

$$\mathbb{P}_{\text{STOP}}(\cdot \mid \text{dir}; \text{adj}, \mathbf{c}_e)$$

EXAMPLE:

(The check is in the mail.)

$$\mathbb{P} = \underbrace{(1 - \mathbb{P}_{\text{STOP}}(\diamond \mid \text{L}; \text{T}))}_0 \times \mathbb{P}_{\text{ATTACH}}(\text{VBZ} \mid \diamond; \text{L})$$

$$\times (1 - \mathbb{P}_{\text{STOP}}(\cdot \mid \text{L}; \text{T}, \text{VBZ})) \times \mathbb{P}_{\text{ATTACH}}(\text{NN} \mid \text{VBZ}; \text{L})$$

$$\times (1 - \mathbb{P}_{\text{STOP}}(\cdot \mid \text{R}; \text{T}, \text{VBZ})) \times \mathbb{P}_{\text{ATTACH}}(\text{IN} \mid \text{VBZ}; \text{R})$$

$$\times \mathbb{P}_{\text{STOP}}(\cdot \mid \text{L}; \text{F}, \text{DT}) \text{ // VBZ} \times \mathbb{P}_{\text{STOP}}(\cdot \mid \text{R}; \text{F}, \text{NN}) \text{ // VBZ}$$

$$\times (1 - \mathbb{P}_{\text{STOP}}(\cdot \mid \text{L}; \text{T}, \text{NN}))^2 \times \mathbb{P}_{\text{ATTACH}}(\text{DT} \mid \text{NN}; \text{L})$$

$$\times (1 - \mathbb{P}_{\text{STOP}}(\cdot \mid \text{R}; \text{T}, \text{IN})) \times \mathbb{P}_{\text{ATTACH}}(\text{NN} \mid \text{IN}; \text{R})$$

$$\times \mathbb{P}_{\text{STOP}}^2(\cdot \mid \text{R}; \text{T}, \text{NN}) \times \mathbb{P}_{\text{STOP}}^2(\cdot \mid \text{L}; \text{F}, \text{DT}) \text{ // NN}$$

$$\times \mathbb{P}_{\text{STOP}}(\cdot \mid \text{L}; \text{T}, \text{IN}) \times \mathbb{P}_{\text{STOP}}(\cdot \mid \text{R}; \text{F}, \text{NN}) \text{ // IN}$$

$$\times \mathbb{P}_{\text{STOP}}^2(\cdot \mid \text{L}; \text{T}, \text{DT}) \times \mathbb{P}_{\text{STOP}}^2(\cdot \mid \text{R}; \text{T}, \text{DT})$$

$$\times \mathbb{P}_{\text{STOP}}(\diamond \mid \text{L}; \text{F}) \times \mathbb{P}_{\text{STOP}}(\diamond \mid \text{R}; \text{T}) \text{ // 1}$$

- truly head-outward model (Alshawi, 1996)
- still split-head, hence efficient (Eisner and Satta, 1999)
- conditions on more observable state — left and right words of phrases being constructed — than hidden head words

→ well-suited to unsupervised learning

DBM-2

Models incomplete inputs based on boundary punctuation.

$$\mathbb{P}_{\text{ATTACH}}(c_r \mid \diamond; \text{L}, \text{comp}) \text{ and } \mathbb{P}_{\text{STOP}}(\cdot \mid \text{dir}; \text{adj}, c_e, \text{comp})$$

EXAMPLES:

(Ungrammatical news-style fragments.)

Odds and Ends captions and headlines
George Morton proper noun phrases
Revenue: \$3.57 billion monetary values
c - Domestic car line items
1:11am date and time expressions

- incomplete fragments are uncharacteristically short
 - roots of fragments are generally not verbs or modals
 - have multiple overlapping grammars coexist in a model
- avoid pitfalls, like inducing nouns as sentence heads

DBM-3

Incorporates sentence-internal punctuation boundaries.

$$\mathbb{P}_{\text{ATTACH}}(c_d \mid c_h; \text{dir}, \text{cross})$$

EXAMPLE: *Continental*s believe that the strongest growth area will be southern Europe.

- punctuation-crossing vets common remote constructions
 - e.g., subordinating conjunctions (IN) and their dependent modal verbs (MD), which are, on average, 4.8 tokens apart
 - avoids bad long distance relations (e.g., far-off DT-NN pairs)
- learn to piece together inter-punctuation fragments

SUMMARY

Unsupervised split-head dependency grammars:

GB (Paskin, 2001); DMV (Klein and Manning, 2004); EVG (Headden et al., 2009).

	$\mathbb{P}_{\text{ATTACH}}$	$\mathbb{P}_{\text{ATTACH}}$	\mathbb{P}_{STOP}
GB	$1 / \{w\} $	$d \mid h; \text{dir}$	$1 / 2$
DMV	$c_r \mid \diamond; \text{L}$	$c_d \mid c_h; \text{dir}$	$\cdot \mid \text{dir}; \text{adj}, c_h$
EVG	$c_r \mid \diamond; \text{L}$	$c_d \mid c_h; \text{dir}, \text{adj}$	$\cdot \mid \text{dir}; \text{adj}, c_h$
DBM-1	$c_r \mid \diamond; \text{L}$	$c_d \mid c_h; \text{dir}$	$\cdot \mid \text{dir}; \text{adj}, c_e$
DBM-2	$c_r \mid \diamond; \text{L}, \text{comp}$	$c_d \mid c_h; \text{dir}$	$\cdot \mid \text{dir}; \text{adj}, c_e, \text{comp}$
DBM-3	$c_r \mid \diamond; \text{L}, \text{comp}$	$c_d \mid c_h; \text{dir}, \text{cross}$	$\cdot \mid \text{dir}; \text{adj}, c_e, \text{comp}$
	(head-root)	(dependent-head)	(direction, adjacency)

RESULTS

Previous state-of-the-art: 38.2% directed dependency accuracy.

— average over all 19 languages of the 2006/7 CoNLL sets (Spitkovsky et al., 2011)

- DBM-1: 40.7% (uniform initialization, no predefined input length cutoff)
- DBM-1 → DBM-2 → DBM-3: **42.9%** (staged curriculum training)

ACKNOWLEDGMENTS

Partially funded by Defense Advanced Research Projects Agency (DARPA) Machine Reading Program, under Air Force Research Laboratory (AFRL) prime contract no. FA8750-09-C-0181. We thank Roi Reichart and Marta Recasens — as well as Marie-Catherine de Marneffe, Roy Schwartz and Mengqiu Wang — for helpful comments on draft versions of our paper.