CS388: Natural Language Processing

Lecture 13: Dependency II





- Dependency syntax: syntactic structure is defined by dependencies Head (parent, governor) connected to dependent (child, modifier) Each word has exactly one parent except for the ROOT symbol Dependencies must form a directed acyclic graph



Recall: Dependencies



Projective <-> no "crossing" arcs



dogs in houses and cats

Crossing arcs:



Recall: Projectivity



the dog ran to the house

credit: Language Log







Recall: Eisner's Algorithm

- Left and right children are built independently, heads are edges of spans
- Complete item: all children are attached, head is at the "tall end"
- Incomplete item: arc from "tall end" to "short end", may still expect children







Neural CRFs for dependency parsing: let c = LSTM embedding of i, p = LSTM embedding of parent(i). score(i, parent(i), \mathbf{x}) = $p^{T}Uc$ *score*(tree) = sum of edge scores $H^{(arc-dep)} \oplus 1$ $H^{(arc-head)}$ $S^{(arc)}$ $U^{(arc)}$



Recall: Biaffine Neural Parsing



- UAS: unlabeled attachment score. Accuracy of choosing each word's parent (*n* decisions per sentence)
- LAS: additionally consider label for each edge
- Log-linear CRF parser, decoding with Eisner algorithm: 91 UAS
- Higher-order features from Koo parser: 93 UAS
- Best English results with neural CRFs (Dozat and Manning): 95-96 UAS

Evaluating Dependency Parsing





Very little work on HPSG in NLP since no real treebank exists

HPSG



Joint model of constituency and dependency combining ideas from Dozat + Manning and Stern et al. Zhou and Zhao (2019)

Parsing with "HPSG"

Parsing with "HPSG"

 Slightly stronger results than Dozat + Manning, significantly better results on Chinese Model

Chen an Andor e Zhang e Cheng e Kuncord Ma and Dozat a Li et al. Ma et al Our (Div Our (Joi Our (Di Our (Joi

	English		Chinese	
	UAS	LAS	UAS	LAS
nd Manning (2014)	91.8	89.6	83.9	82.4
et al. (2016)	94.61	92.79	_	_
et al. (2016)	93.42	91.29	87.65	86.17
et al. (2016)	94.10	91.49	88.1	85.7
o et al. (2016)	94.26	92.06	88.87	87.30
Hovy (2017)	94.88	92.98	89.05	87.74
nd Manning (2017)	95.74	94.08	89.30	88.23
(2018a)	94.11	92.08	88.78	86.23
l. (2018)	95.87	94.19	90.59	89.29
vision)	94.32	93.09	89.14	87.31
int)	96.09	94.68	91.21	89.15
vision*)	-	-	91.69	90.54
int*)	-	-	93.24	91.95
	_			

Zhou and Zhao (2019)

Transition-based (shift-reduce) dependency parsing

Approximate, greedy inference — fast, but a little bit weird!

Shift-Reduce Parsing

- Similar to deterministic parsers for compilers Also called transition-based parsing
- A tree is built from a sequence of incremental decisions moving left to right through the sentence
- Stack containing partially-built tree, buffer containing rest of sentence
- Shifts consume the buffer, reduces build a tree on the stack

Shift-Reduce Parsing

Shift: top of buffer -> top of stack

- Initial state: Stack: [ROOT] Buffer: [I ate some spaghetti bolognese]
 - Shift 1: Stack: [ROOT I] Buffer: [ate some spaghetti bolognese]
 - Shift 2: Stack: [ROOT | ate] Buffer: [some spaghetti bolognese]

- State: Stack: [ROOT | ate] Buffer: [some spaghetti bolognese]
- Left-arc (reduce): Let σ denote the stack, $\sigma | w_{-1} =$ stack ending in w_{-1}
 - Pop two elements, add an arc, put them back on the stack"

$$\sigma|w_{-2}, w_{-1} \to \sigma|w_{-1} \quad u$$

State: Stack: [ROOT ate]

- v_{-2} is now a child of w_{-1}
- Buffer: [some spaghetti bolognese]

- Start: stack contains [ROOT], buffer contains [I ate some spaghetti bolognese] Arc-standard system: three operations
 - Shift: top of buffer -> top of stack

Left-Arc:
$$\sigma | w_{-2}, w_{-1} \rightarrow \sigma |$$

- Fight-Arc $\sigma | w_{-2}, w_{-1}
 ightarrow \sigma | w_{-2}$, w_{-1} is now a child of $| w_{-2}
 angle$
- End: stack contains [ROOT], buffer is empty []

How many transitions do we need if we have n words in a sentence?

Arc-Standard Parsing

- w_{-2} is now a child of w_{-1} w_{-1}

- Could do the left arc later! But no reason to wait
- Can't attach ROOT <- ate yet even though this is a correct dependency!</p>

Arc-Standard Parsing

- top of buffer -> top of stack S pop two, left arc between them LA RA pop two, right arc between them
- [I ate some spaghetti bolognese]
- [ate some spaghetti bolognese]
- [some spaghetti bolognese]
- [some spaghetti bolognese]

Arc-Standard Parsing

top of buffer -> top of stack S pop two, left arc between them LA RA pop two, right arc between them

[some spaghetti bolognese]

[bolognese]

[bolognese]

Arc-Standard Parsing

top of buffer -> top of stack S pop two, left arc between them LA RA pop two, right arc between them

Stack consists of all words that are still waiting for right children, end with a bunch of right-arc ops

spagnetti

some

bolognese

Final state:

- Arc-eager (Nivre, 2004): lets you add right arcs sooner and keeps items on stack, separate reduce action that clears out the stack
- Arc-swift (Qi and Manning, 2017): explicitly choose a parent from what's on the stack
- Many ways to decompose these, which one works best depends on the language and features (nonprojective variants too!)

Other Systems

[ROOT]

- How do we make the right decision in this case?
- Only one legal move (shift)

[ROOT ate some spaghetti]

- How do we make the right decision in this case? (all three actions legal)
- Multi-way classification problem: shift, left-arc, or right-arc?

 $\operatorname{argmax}_{a \in \{S, LA, RA\}} w^{\top} f(\operatorname{stack}, \operatorname{buffer}, a)$

Building Shift-Reduce Parsers

[I ate some spaghetti bolognese]

[bolognese]

[ROOT ate some spaghetti]

Features to know this should left-arc?

- One of the harder feature design tasks!
- leftmost and rightmost children of top items on the stack

Features for Shift-Reduce Parsing

[bolognese]

In this case: the stack tag sequence VBD - DT - NN is pretty informative — looks like a verb taking a direct object which has a determiner in it

Things to look at: top words/POS of buffer, top words/POS of stack,

Training a Greedy Model

[ROOT ate some spaghetti] [bolognese]

- $\operatorname{argmax}_{a \in \{S, LA, RA\}} w^{\top} f(\operatorname{stack}, \operatorname{buffer}, a)$
- Can turn a tree into a decision sequence a by building an oracle
- Train a classifier to predict the right decision using these as training data
- Training data assumes you made correct decisions up to this point and teaches you to make the correct decision, but what if you screwed up...

Greedy training State space Gold end state

making bad decisions but don't condition on bad decisions

Many early-2000s constituency parsers were ~5 sentences/sec

Using S-R used to mean taking a performance hit compared to graph-based, that's no longer (quite as) true

Speed Tradeoffs

De	ev	Test		Speed
UAS	LAS	UAS	LAS	(sent/s)
89.9	88.7	89.7	88.3	51
90.3	89.2	89.9	88.6	63
90.0	88.8	89.9	88.5	560
90.1	88.9	90.1	88.7	535
92.1	90.8	92.0	90.5	12
92.2	91.0	92.0	90.7	1013

Chen and Manning (2014)

Global Decoding

[dinner] [ROOT gave] him

[ROOT gave dinner] him

Global Decoding

- [ROOT gave him] [dinner]

Global Decoding: A Cartoon

Both wrong! Also both probably low scoring!

Correct, high scoring option

Global Decoding: A Cartoon

Lookahead can help us avoid getting stuck in bad spots

- Global model: maximize sum of scores over all decisions
- Similar to how Viterbi works: we maintain uncertainty over the current state so that if another one looks more optimal going forward, we can use that one

[ROOT gave him] [dinner]

Global Shift-Reduce Parsing

[ROOT gave him] [dinner]

• Greedy: repeatedly execute $a_{\text{best}} \leftarrow \operatorname{argmax}_a w^{\top} f(s, a)$ $s \leftarrow a_{\text{best}}(s)$

Can we do search exactly? How many states s are there?

No! Use beam search

Global:

$$\operatorname{argmax}_{\mathbf{s},\mathbf{a}} w^{\top} f(\mathbf{s},\mathbf{a}) = \sum_{i=1}^{2n} w^{\top} f(s_i)$$
$$s_{i+1} = a_i(s_i)$$
ny states s are there?

THERSITY OF THE ASSIDILATE CHUITALING

Beam Search

Maintain a beam of k plausible states at the current timestep, expand each and only keep top k best new ones

How good is beam search?

- k=1: greedy search
- Choosing beam size:
 - 2 is usually better than 1
 - Usually don't use larger than 50
 - Depends on problem structure

Global Shift-Reduce Parsing

Global Training

- If using global inference, should train the parser in a global fashion as well: use structured perceptron / structured SVM
- Model treats an entire derivation as something to featurize
- No algorithm like Viterbi for doing efficient parsing, so use beam search

State-of-the-art Transition-Based Parsers

- 2005: Eisner algorithm graph-based parser was SOTA (~91 UAS)
- 2010: Koo's 3rd-order parser was SOTA for graph-based (~93 UAS)
- 2012: Maltparser was SOTA was for transition-based (~90 UAS)
- 2014: Chen and Manning got 92 UAS with transition-based neural model
- 2016: Improvements to Chen and Manning

Dependency Parsers

State-of-the-art Parsers

Chen and Manning (2014)

- Close to state-of-the-art, released by Google publicly
- 94.61 UAS on the Penn Treebank using a global transition-based system with early updating (compared to 95.8 for Dozat, 93.7 for Koo in 2009)
 - Additional data harvested via "tri-training", form of self-training
- Feedforward neural nets looking at words and POS associated with words in the stack / those words' children / words in the buffer
- Feature set pioneered by Chen and Manning (2014), Google fine-tuned it
- Shift-reduce parsers are often playing "catch-up", hard to really push the SOTA with shift-reduce because it's harder to design models
 - Andor et al. (2016)

Shift-Reduce Constituency

combine with no label for ternary rules

Can do shift-reduce for constituency as well, reduce operation builds constituents

ural action	label action	stack after	bracket
PRP)	label-NP	0_1	$_0 NP_1$
/MD)	nolabel	$0 \frown 1 \frown 2$	
e/VBP)	nolabel	$0 \square 1 \square 2 \square 3$	
	nolabel	$0 \frown 1 \frown 3$	
ting/VBG)	nolabel	$0 \frown 1 \frown 3 \frown 4$	
h/NN)	label-NP	$0 \bigcirc 1 \bigcirc 3 \bigcirc 4 \bigcirc 5$	$_4NP_5$
	label-S-VP	$0 \square 1 \square 3 \square 5$	$_{3}S_{5}, _{3}VP_{5}$
	label-VP	$0 \frown 1 \frown 5$	$_{1}\mathrm{VP}_{5}$
	label-S	$0 \frown 5$	$_{0}S_{5}$

(b) static oracle actions

Cross and Huang (2016)

- Shift-reduce parsing can work nearly as well as graph-based
- Arc-standard system for transition-based parsing
- Purely greedy or more "global" approaches
- Next time: semantic parsing