Neural Constituency Parsing



Dan Klein CS 288



Syntactic Parsing

She enjoys playing tennis.

Syntactic Parsing



Historical Trends



[Slide from Slav Petrov]

Output Correlations







She enjoys playing tennis.



Parsing as Span Classification































Span Classification



Span Classification





Non-Constituents



... But Will We Get a Tree Out?









Does It Work?



Neural parsers no longer have much of the model structure provided to classical parsers.

How do they perform so well without it?

Why don't we need a grammar?

Adjacent tree labels are redundant with LSTM features

If we can predict surrounding tree labels from our LSTM representation of the input, then this information doesn't need to be provided explicitly by grammar production rules

We find that for **92.3%** of spans, the label of the span's parent can predicted from the neural representation of the span



Do we need tree constraints?

Not for F1

Many neural parsers no longer model output correlations with grammar rules, but still use output correlations from tree constraints

Predicting span brackets independently gives **nearly identical performance** on PTB development set F1 and produces valid trees for **94.5%** of sentences



What word representations do we need?

A character LSTM is sufficient

Word Only	91.44
Word and Tag	92.09
Character LSTM Only	92.24
Character LSTM and Word	92.22
Character LSTM, Word, and Tag	92.24

What about lexicon features?

The character LSTM captures the same information

Heavily engineered lexicons used to be critical to good performance, but neural models typically don't use them

Word features from the Berkeley Parser (Petrov and Klein 2007) can be predicted with over **99.7%** accuracy from the character LSTM representation

Do LSTMs introduce useful inductive bias compared to feedforward networks?

Yes!

We compare a truncated LSTM with feedforward architectures that are given the same inputs

The LSTM outperformed the best feedforward by **6.5 F1**



•











What Helps?



F1 (English, dev)





Pre-Training

Problem: Input has more variation than output

Need to handle:

- Rare words not seen during training
- Word forms in morphologically rich languages
- Contextual paraphrase / lexical variation

Historical Trends



[Slide from Slav Petrov]



 Knowledge modularity: Learn domain-general knowledge from one data source and use it solve specific problems elsewhere



Parsing as Span Classification



Pretraining





Encoder Architectures



Encoder Architectures



Results: Multilingual







Does Structure Help?



Figure 1: Labelled bracketing F1 versus minimum span length for the English corpora. F1 scores for the In-Order parser with BERT (orange) and the Chart parser with BERT (cyan) start to diverge for longer spans.

Out of Domain Parsing

	Berkeley		BLLIP		In-Order		Chart	
	F1	Δ Err.						
WSJ Test	90.06	+0.0%	91.48	+0.0%	91.47	+0.0%	93.27	+0.0%
Brown All Genia All EWT All	84.64 79.11 77.38	+54.5% +110.2% +127.6%	85.89 79.63 79.91	+65.6% +139.1% +135.8%	85.60 80.31 79.07	+68.9% +130.9% +145.4%	88.04 82.68 82.22	+77.7% +157.4% +164.2%

Neural parsers improve out-of-domain numbers, but not more than in-domain numbers



Other Neural Constituency Parsers



- Back to at least Henderson 1998!
- Recent directions:
 - Shift-Reduce, eg Cross and Huang 2016
 - SR/Generative, eg Dyer et al 2016 (RNNG)
 - In-Order Generative, eg Liu and Zhang 2017