Expected F-Measure Training for Shift-Reduce Parsing with Recurrent Neural Networks

Wenduan Xu

Cambridge University

12th June, 2016

Another Title...

Training a global shift-reduce model that is also optimized towards the final evaluation metric



step	stack (s_n, \ldots, s_1, s_0)	queue $(q_0, q_1 \ldots, q_n)$	action
0		Ms. Haag plays Elianti	

step	stack (s_n, \ldots, s_1, s_0)	queue $(q_0, q_1 \dots, q_n)$	action
0		Ms. Haag plays Elianti	
1		Haag plays Elianti	SHIFT
2		plays Elianti	SHIFT

stack (s_n, \ldots, s_1, s_0)	queue $(q_0, q_1 \dots, q_n)$	action
	Ms. Haag plays Elianti	
	Haag plays Elianti	SHIFT
	plays Elianti	SHIFT
	plays Elianti	REDUCE
	stack (s_n, \ldots, s_1, s_0)	stack (s_n, \ldots, s_1, s_0) queue $(q_0, q_1 \ldots, q_n)$ Ms. Haag plays Elianti Haag plays Elianti plays Elianti plays EliantiJoint

step	stack (s_n, \ldots, s_1, s_0)	queue $(q_0, q_1 \ldots, q_n)$	action
0		Ms. Haag plays Elianti	
1		Haag plays Elianti	SHIFT
2		plays Elianti	SHIFT
3		plays Elianti	REDUCE
4		plays Elianti	UNARY
5		Elianti	SHIFT
6			SHIFT
7			UNARY

step	stack (s_n, \ldots, s_1, s_0)	queue $(q_0, q_1 \ldots, q_n)$	action
0		Ms. Haag plays Elianti	
1		Haag plays Elianti	SHIFT
2		plays Elianti	SHIFT
3		plays Elianti	REDUCE
4		plays Elianti	UNARY
5		Elianti	SHIFT
6			SHIFT
7			UNARY
8	212		REDUCE
9			REDUCE







SHIFT SHIFT SHIFT REDUCE SHIFT SHIFT ...

• Linear model (perceptron, SVM etc.)

 $score(t_i) = \mathbf{f}(\langle s, q \rangle, t_i) \cdot \mathbf{w}$

• Linear model (perceptron, SVM etc.)

 $score(t_i) = \mathbf{f}(\langle s, q \rangle, t_i) \cdot \mathbf{w}$ $score(t) = \sum_i score(t_i)$

• Linear model (perceptron, SVM etc.)

 $score(t_i) = \mathbf{f}(\langle s, q \rangle, t_i) \cdot \mathbf{w}$ $score(t) = \sum_i score(t_i)$ $t^* = \arg \max_t score(t)$

• Linear model (perceptron, SVM etc.)

 $score(t_i) = \mathbf{f}(\langle s, q \rangle, t_i) \cdot \mathbf{w}$ $score(t) = \sum_i score(t_i)$ $t^* = \arg \max_t score(t)$

· Great flexibility in defining the feature functions

• Linear model (perceptron, SVM etc.)

 $score(t_i) = \mathbf{f}(\langle s, q \rangle, t_i) \cdot \mathbf{w}$ $score(t) = \sum_i score(t_i)$ $t^* = \arg \max_t score(t)$

- Great flexibility in defining the feature functions
 - results in millions of sparse indicator features

• Linear model (perceptron, SVM etc.)

 $score(t_i) = \mathbf{f}(\langle s, q \rangle, t_i) \cdot \mathbf{w}$ $score(t) = \sum_i score(t_i)$ $t^* = \arg \max_t score(t)$

- Great flexibility in defining the feature functions
 - results in millions of sparse indicator features
- Used in many parsers
 - (e.g., Yamada and Matsumoto, 2003; Huang and Sagae, 2010;
 Zhang and Clark, 2011; Zhang and Nivre, 2011; Goldberg and Nivre, 2012; Bohnet et al., 2013; Zhu et al., 2013;)

Shift-Reduce Parsing Models

_	feature templates
1	S_0 wp, S_0 c, S_0 pc, S_0 wc,
	S_1 wp, S_1 c, S_1 pc, S_1 wc,
	S_2pc , S_2wc ,
	S ₃ pc, S ₃ wc,
2	Q_0 wp, Q_1 wp, Q_2 wp, Q_3 wp,
3	S ₀ Lpc, S ₀ Lwc, S ₀ Rpc, S ₀ Rwc,
	S ₀ Upc, S ₀ Uwc,
	S ₁ Lpc, S ₁ Lwc, S ₁ Rpc, S ₁ Rwc,
	S_1 Upc, S_1 Uwc,
4	S_0wcS_1wc , S_0cS_1w , S_0wS_1c , S_0cS_1c ,
	S_0wcQ_0wp , S_0cQ_0wp , S_0wcQ_0p , S_0cQ_0p ,
	S_1wcQ_0wp , S_1cQ_0wp , S_1wcQ_0p , S_1cQ_0p ,
5	$S_0wcS_1cQ_0p, S_0cS_1wcQ_0p, S_0cS_1cQ_0wp,$
	$S_0 c S_1 c Q_0 p$, $S_0 p S_1 p Q_0 p$,
	$S_0wcQ_0pQ_1p$, $S_0cQ_0wpQ_1p$, $S_0cQ_0pQ_1wp$,
	$S_0 c Q_0 p Q_1 p$, $S_0 p Q_0 p Q_1 p$,
	$S_0wcS_1cS_2c$, $S_0cS_1wcS_2c$, $S_0cS_1cS_2wc$,
	$S_0cS_1cS_2c$, $S_0pS_1pS_2p$,
6	S ₀ cS ₀ HcS ₀ Lc, S ₀ cS ₀ HcS ₀ Rc,
	$S_1 c S_1 H c S_1 R c$,
	$S_0 c S_0 R c Q_0 p$, $S_0 c S_0 R c Q_0 w$,
	$S_0cS_0LcS_1c$, $S_0cS_0LcS_1w$,
	$S_0 cS_1 cS_1 Rc$, $S_0 wS_1 cS_1 Rc$.

Table 1: Feature templates.

(Zhang and Clark, 2011)



Figure 2: Our neural network architecture.

(Chen and Manning, 2014)

NNBeam (Train)Beam (Test)globalC&M, 2014✓X✓X

	NN	Beam (Train)	Beam (Test)	global
C&M, 2014	\checkmark	×	\checkmark	X
this work	\checkmark	\checkmark	\checkmark	1

	NN	Beam (Train)	Beam (Test)	global
C&M, 2014	\checkmark	×	\checkmark	X
this work	\checkmark	\checkmark	\checkmark	\checkmark

At the same time, the model is optimized towards the final evaluation metric \checkmark

Global Shift-Reduce NN Models



(Weiss et al., 2015)

Related Work

- Watanabe and Sumita, 2015
 - max-margin based objective
 - max-violation updates (Huang et al., 2012)
- Zhou et al., 2015
 - based on Chen and Manning, 2014
 - contrastive learning (Hinton, 2002; LeCun and Huang, 2005; Liang and Jordan, 2008)
- Andor et al., 2016
 - based on Chen and Manning, 2014 and Weiss et al., 2015
 - CRF (Bottou et al., 1997; Le Cun et al., 1998; Lafferty et al., 2001)

Related Work

- Watanabe and Sumita, 2015
 - max-margin based objective
 - max-violation updates (Huang et al., 2012)
- Zhou et al., 2015
 - based on Chen and Manning, 2014
 - contrastive learning (Hinton, 2002; LeCun and Huang, 2005; Liang and Jordan, 2008)
- Andor et al., 2016
 - based on Chen and Manning, 2014 and Weiss et al., 2015
 - CRF (Bottou et al., 1997; Le Cun et al., 1998; Lafferty et al., 2001)
- Optimizing task-specific metrics for parsing
 - (e.g., Goodman, 1996; Smith and Eisner, 2006; Auli and Lopez, 2011)

• Train a baseline model using a cross-entropy loss (pretraining)



$$L(\theta) = -\sum_{k}^{T_i} p(t_k)$$

• Train a baseline model using a cross-entropy loss (pretraining)



















$$J(\theta) = -\mathsf{xF1}(\theta) = -\sum_{y_i \in \Lambda(x_n)} p(y_i|\theta)\mathsf{F1}(\Delta_{y_i}, \Delta_{x_n}^{\mathcal{G}}),$$

$$p(y_i|\theta) = \frac{\exp\{\omega(y_i)\}}{\sum_{y \in \Lambda(x_n)} \exp\{\omega(y)\}},$$

$$J(\theta) = -\mathbf{x}\mathsf{F1}(\theta) = -\sum_{y_i \in \Lambda(x_n)} \frac{\mathsf{p}(y_i|\theta)\mathsf{F1}(\Delta_{y_i}, \Delta_{x_n}^{\mathsf{G}}),$$

$$p(y_i|\theta) = \frac{\exp\{\omega(y_i)\}}{\sum_{y \in \Lambda(x_n)} \exp\{\omega(y)\}},$$

$$J(\theta) = -\mathsf{xF1}(\theta) = -\sum_{y_i \in \Lambda(x_n)} p(y_i|\theta) \mathsf{F1}(\Delta_{y_i}, \Delta_{x_n}^{\mathcal{G}}),$$

$$p(y_i|\theta) = \frac{\exp\{\omega(y_i)\}}{\sum_{y \in \Lambda(x_n)} \exp\{\omega(y)\}},$$

$$egin{aligned} rac{\partial J(heta)}{\partial heta} &= -\sum_{y_i \in \Lambda(x_n)} \sum_{y_{ij} \in y_i} rac{\partial J(heta)}{\partial s_ heta(y_{ij})} rac{\partial s_ heta(y_{ij})}{\partial heta} \ &= -\sum_{y_i \in \Lambda(x_n)} \sum_{y_{ij} \in y_i} \delta_{y_{ij}} rac{\partial s_ heta(y_{ij})}{\partial heta}. \end{aligned}$$

$$J(\theta) = -\mathsf{xF1}(\theta) = -\sum_{y_i \in \Lambda(x_n)} p(y_i|\theta) \mathsf{F1}(\Delta_{y_i}, \Delta_{x_n}^{\mathcal{G}}),$$

$$p(y_i|\theta) = \frac{\exp\{\omega(y_i)\}}{\sum_{y \in \Lambda(x_n)} \exp\{\omega(y)\}},$$

$$rac{\partial J(heta)}{\partial heta} = -\sum_{\mathbf{y}_i \in \Lambda(\mathbf{x}_n)} \sum_{\mathbf{y}_{ij} \in \mathbf{y}_i} rac{\partial J(heta)}{\partial s_{ heta}(\mathbf{y}_{ij})} rac{\partial s_{ heta}(y_{ij})}{\partial heta} \\ = -\sum_{\mathbf{y}_i \in \Lambda(\mathbf{x}_n)} \sum_{\mathbf{y}_{ij} \in \mathbf{y}_i} \delta_{\mathbf{y}_{ij}} rac{\partial s_{ heta}(y_{ij})}{\partial heta}.$$

output	action sequence	$\omega(y_i)$	F1
<i>y</i> ₁	<i>Y</i> ₁₁ <i>Y</i> ₁₂ <i>Y</i> _{1<i>i</i>}	-0.60	0.67
<i>y</i> ₂	Y 21 Y 22 · · · Y 2j	-1.5	0.81
<i>y</i> ₃	y 31 y 32 · · · y 3k	-4.96	0.90

output	action sequence	$\omega(y_i)$	F1
<i>y</i> ₁	<i>Y</i> ₁₁ <i>Y</i> ₁₂ <i>Y</i> _{1<i>i</i>}	-0.60	0.67
y 2	Y 21 Y 22 · · · Y 2j	-1.5	0.81
<i>y</i> ₃	y 31 y 32 · · · y 3k	-4.96	0.90

$$J(\theta) = -\mathsf{x}\mathsf{F1}(\theta) = -\sum_{y_i \in \Lambda(x_n)} p(y_i|\theta)\mathsf{F1}(\Delta_{y_i}, \Delta_{x_n}^{\mathsf{G}}) = 71.00$$

output	action sequence	$\omega(y_i)$	F1
<i>y</i> ₁	<i>Y</i> ₁₁ <i>Y</i> ₁₂ <i>Y</i> _{1<i>i</i>}	-0.60	0.67
<i>y</i> ₂	y 21 y 22 y 2j	-1.5	0.81
<i>y</i> ₃	y 31 y 32 · · · y 3k	-4.96	0.90

$$J(\theta) = -\mathsf{x}\mathsf{F1}(\theta) = -\sum_{y_i \in \Lambda(x_n)} p(y_i|\theta)\mathsf{F1}(\Delta_{y_i}, \Delta_{x_n}^G) = 71.00$$

output	action sequence	$\omega(y_i)$	F1
<i>z</i> ₁	$Z_{11} \ Z_{12} \ldots \ Z_{1i}$	-0.90	0.75
z_2	$Z_{21} \ Z_{22} \ldots \ Z_{2j}$	-0.99	0.85
Z 3	$Z_{31} Z_{32} \ldots Z_{3k}$	-3.76	0.95

output	action sequence	$\omega(y_i)$	F1
<i>y</i> ₁	<i>Y</i> ₁₁ <i>Y</i> ₁₂ <i>Y</i> _{1<i>i</i>}	-0.60	0.67
<i>y</i> ₂	y 21 y 22 y 2j	-1.5	0.81
<i>y</i> ₃	y 31 y 32 · · · y 3k	-4.96	0.90

$$J(\theta) = -\mathsf{x}\mathsf{F1}(\theta) = -\sum_{y_i \in \Lambda(x_n)} p(y_i|\theta)\mathsf{F1}(\Delta_{y_i}, \Delta_{x_n}^G) = 71.00$$

output	action sequence	$\omega(y_i)$	F1
z_1	$Z_{11} \ Z_{12} \ldots \ Z_{1i}$	-0.90	0.75
Z 2	$Z_{21} \ Z_{22} \ldots \ Z_{2j}$	-0.99	0.85
Z 3	$Z_{31} Z_{32} \ldots Z_{3k}$	-3.76	0.95

$$J(\theta) = -\mathsf{x}\mathsf{F1}(\theta) = -\sum_{z_i \in \Lambda(\mathsf{x}_n)} p(z_i|\theta)\mathsf{F1}(\Delta_{z_i}, \Delta_{\mathsf{x}_n}^G) = 80.20$$

Experiments on CCGBank

Supertagger	Dev	Test	
C&C	91.50	92.02	
RNN	93.07	93.00	
BRNN	93.49	93.52	

• A bidirectional RNN supertagger extending the unidirectional one in Xu et al., 2015, using the same model and training parameters

Eval: F1 over Labeled, Directed CCG Deps



 $\langle the, NP/N_1, 1, books, \rangle$ $\langle likes, (S \setminus NP_1)/NP_2, 1, John \rangle$ $\langle which, (NP/NP_1)/(S/NP)_2, 2, likes \rangle$ $\langle which, (NP/NP_1)/(S/NP)_2, 1, books \rangle$ $\langle likes, (S \setminus NP_1)/NP_2, 2, books \rangle$

The Greedy Model and Beam Search (Dev)

beam	F1
b = 1	84.61
<i>b</i> = 2	84.94
<i>b</i> = 4	85.01
<i>b</i> = 6	85.02
<i>b</i> = 8	85.02
b = 16	85.01

 $b \in \{6, 8\}$ gives +0.41% F1 over b = 1





Test Set Parsing Results

Model	LP	LR	LF	CAT	Speed
C&C (normal)	85.45	83.97	84.70	92.83	97.90
C&C (hybrid)	86.24	84.17	85.19	93.00	95.25
Zhang11 ($b = 16$)	87.04	84.14	85.56	92.95	49.54
Xu14 ($b = 128$)	87.03	85.08	86.04	93.10	12.85
Am16 $(b = 1)$	-	-	83.27	91.89	350.00
Am16 $(b = 16)$	-	-	85.57	92.86	10.00
RNN-greedy $(b = 1)$	88.53	81.65	84.95	93.57	337.45
RNN-greedy $(b = 6)$	88.54	82.77	85.56	93.68	96.04
RNN-xF1 $(b = 8)$	88.74	84.22	86.42	93.87	67.65

- Zhang11 = Zhang and Clark, 2011*, Xu14 = Xu et al., 2014; AM16 = Ambati et al., 2016 (NN + Struct. Percep (Weiss et al., 2015))
- The xF1 model improves LR by 2.57% and LF by 1.47% over RNN-greedy (b = 1)

Test Set Parsing Results

Model	LP	LR	LF	CAT	Speed
C&C (normal)	85.45	83.97	84.70	92.83	97.90
C&C (hybrid)	86.24	84.17	85.19	93.00	95.25
Zhang11 ($b = 16$)	87.04	84.14	85.56	92.95	49.54
Xu14 ($b = 128$)	87.03	85.08	86.04	93.10	12.85
Am16 $(b = 1)$	-	-	83.27	91.89	350.00
Am16 $(b = 16)$	-	-	85.57	92.86	10.00
RNN-greedy $(b = 1)$	88.53	81.65	84.95	93.57	337.45
RNN-greedy $(b = 6)$	88.54	82.77	85.56	93.68	96.04
RNN-xF1 ($b = 8$)	88.74	84.22	86.42	93.87	67.65

- Zhang11 = Zhang and Clark, 2011*, Xu14 = Xu et al., 2014; AM16 = Ambati et al., 2016 (NN + Struct. Percep (Weiss et al., 2015))
- The xF1 model improves LR by 2.57% and LF by 1.47% over RNN-greedy (b = 1)
- Auli and Lopez, 2011 uses a softmax-margin objective (Gimpel and Smith, 2010) on the C&C parser

The End: Questions?

Model	LP	LR	LF	CAT	Speed
C&C (normal)	85.45	83.97	84.70	92.83	97.90
C&C (hybrid)	86.24	84.17	85.19	93.00	95.25
Zhang11 ($b = 16$)	87.04	84.14	85.56	92.95	49.54
Xu14 ($b = 128$)	87.03	85.08	86.04	93.10	12.85
Am16 $(b = 1)$	-	-	83.27	91.89	350.00
Am16 $(b = 16)$	-	-	85.57	92.86	10.00
RNN-greedy $(b = 1)$	88.53	81.65	84.95	93.57	337.45
RNN-greedy $(b = 6)$	88.54	82.77	85.56	93.68	96.04
RNN-xF1 $(b = 8)$	88.74	84.22	86.42	93.87	67.65

- Zhang11 = Zhang and Clark, 2011*, Xu14 = Xu et al., 2014, AM16 = Ambati et al., 2016 (NN + Struct. Percep (Weiss et al., 2015))
- The xF1 model improves LR by 2.57% and LF by 1.47% over RNN-greedy (b = 1)
- Auli and Lopez, 2011 uses a softmax-margin objective (Gimpel and Smith, 2010) on the C&C parser