



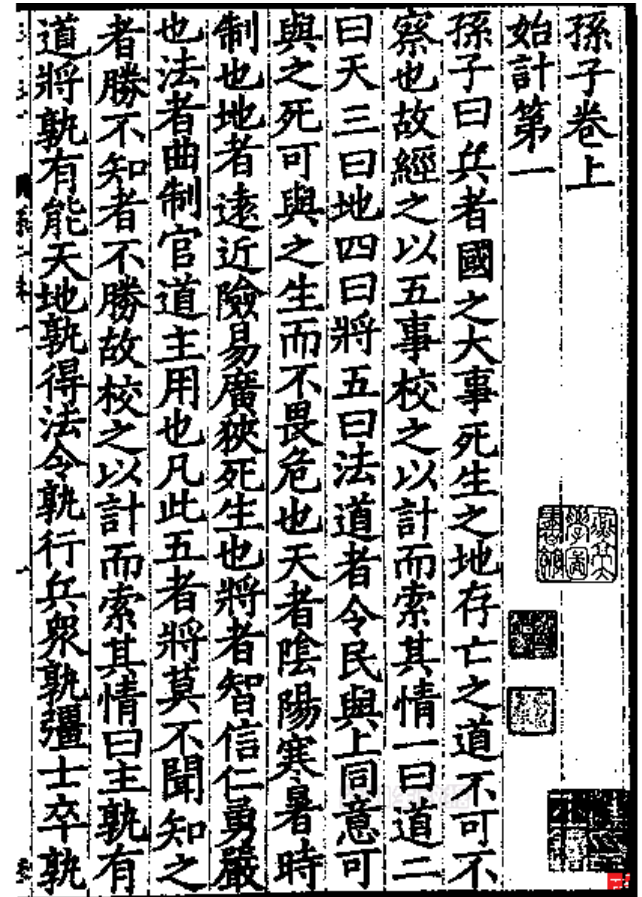
Sentence Segmentation for Ancient Chinese Texts Based on RNN

一种基于循环神经网络的古文断句方法

Boli Wang, Xiaodong Shi, Jinsong Su
Xiamen University

Introduction

- Ancient Chinese texts
 - No punctuations



Introduction

- Ancient Chinese texts
 - No punctuations
- Ju Dou (句读)
 - Rely on language expertise
 - Time-consuming



Introduction

- Ancient Chinese texts
 - No punctuations
- Ju Dou (句读)
 - Rely on language expertise
 - Time-consuming
- Digitalized ancient Chinese texts
 - Segmented: only about 4000
- **Automatic sentence segmentation**

Related work

- Rule-based Methods
 - RegExp (Huang, 2008)
- Statistical Methods
 - N-gram (Chen, 2007)
 - **CRF** (H. Zhang, 2009) (K. Zhang, 2009) (Huang, 2010)
 - Handcrafted local context features
 - State-of-the-art
- NN-based Methods
 - Sequence generation: RNNLM (Wang, 2016)
 - Underperform CRF
 - High time complexity
- **This paper**: seq-to-seq

Our approach

– **Sequence labeling:** $\mathbf{o}_t^{[i]} = P(y_t = i | x_1 x_2 \dots x_N)$

– Punctuations: 。 ? ! , ; :

曰：“朕此衣已三浣矣。”

– 6-tag set (H. Zhang, 2009)

$$T = \{B, M, E3, E2, E, S\}$$

曰/S 朕/B 此/M 衣/M
已/M 三/E3 浣/E2 矣/E

B – the first character of a clause

M – the middle character of a clause

E3 – the last 3rd character of a clause

E2 – the last 2nd character of a clause

E – the last character of a clause

S – a single character forming a clause

Our approach

- Model: RNN

- Input layer

Character embedding: $x_t \rightarrow e_t$

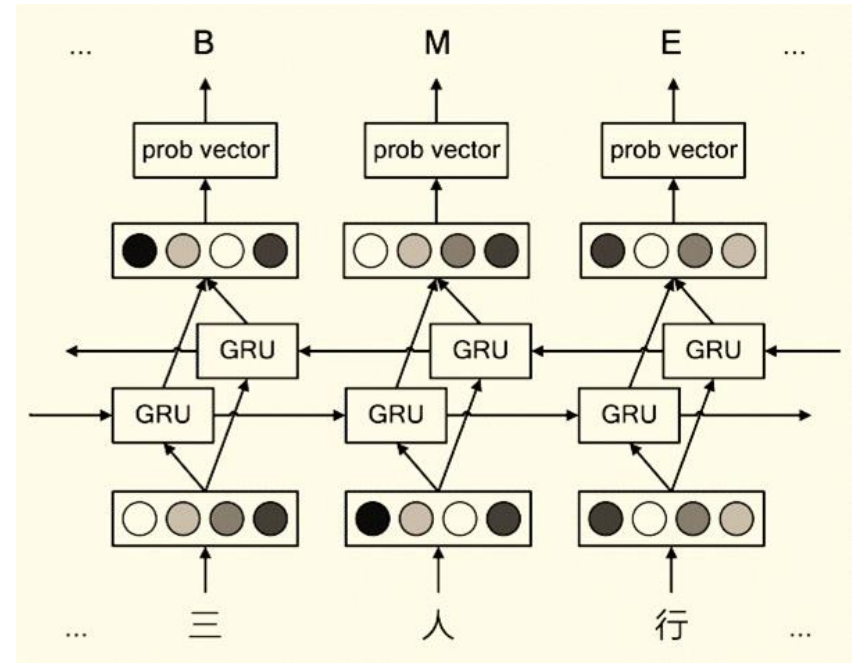
- Hidden layer:

Bi-directional GRU

$$\vec{h}_t = f(e_t, \vec{h}_{t-1}) \quad \overleftarrow{h}_t = f(e_t, \overleftarrow{h}_{t+1})$$

- Output layer: softmax

$$o_t^{[i]} = \frac{\exp(W_{\vec{h}o}^{[i]} \vec{h}_t + W_{\overleftarrow{h}o}^{[i]} \overleftarrow{h}_t + b_o)}{\sum_{j \in T} \exp(W_{\vec{h}o}^{[j]} \vec{h}_t + W_{\overleftarrow{h}o}^{[j]} \overleftarrow{h}_t + b_o)}$$



Our approach

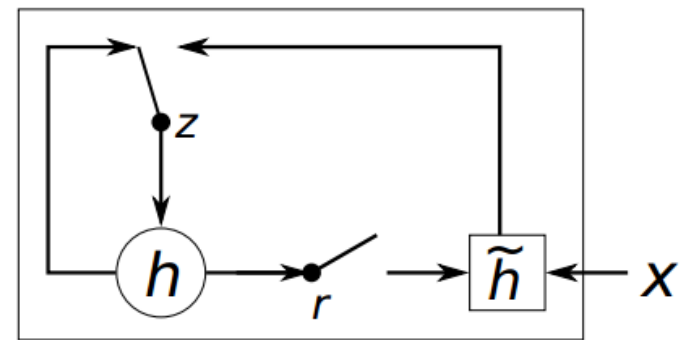
- Gated Recurrent Unit (Cho, 2014)

$$r_t = \text{sigmoid}(W_r x_t + U_r h_{t-1} + b_r),$$

$$\tilde{h}_t = \tanh(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h),$$

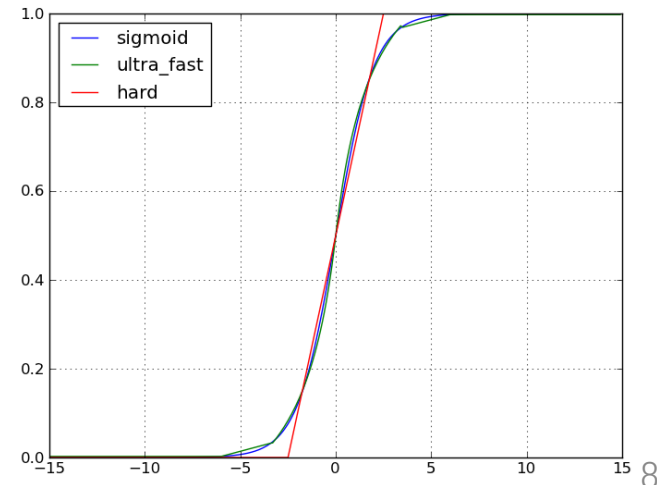
$$z_t = \text{sigmoid}(W_z x_t + U_z h_{t-1} + b_z),$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t,$$



- Our implementation
 - using `hard_sigmoid`

$$\text{hard_sigmoid}(x) = \begin{cases} 0, & x < -2.5 \\ 0.2x + 0.5, & -2.5 \leq x < 2.5 \\ 1, & x \geq 2.5 \end{cases}$$



Our approach

- Training

- Maximum-likelihood estimation

$$Loss(\theta) = - \sum_i \log(P(y_i|x_i, \theta))$$

- Mini-batch gradient descent

- Back-Propagation Through Time (Werbos, 1990)
 - RMSProp (Tieleman, 2012)
 - Gradient Re-normalization (Salimans, 2016)
 - Early Stopping (Yao, 2007)

- Pre-training: using a GRU-RNNLM

Our approach

- Decoding
 - Given a char sequence $x_{1:N}$
 - Goal: find $y_{1:N}$ to maximize $s(x_{1:N}, y_{1:N})$

$$s(x_{1:N}, y_{1:N}) = \sum_{i=1}^N (o_i^{[y_i]} + \alpha \cdot \underline{A_{y_{i-1}y_i}} + \beta \cdot \underline{C(y_{i:N})})$$

state transition probability

length penalty (# of sents)

- Tag inference: beam search

Experiments

– Datasets

Dataset	Number of characters	Charset size	Source
Training Set	237M	23905	STCC
Development Set	0.01M	1890	HDST
Test Set 1	0.32M	6188	4HN
Test Set 2	0.36M	5755	HDST&4HN

– Superfection Traditional Chinese Corpus

http://cloudtranslation.cc/corpus_tc.html

– Baselines

– RNNLM: CLM1 / CLM6 (Wang, 2016)

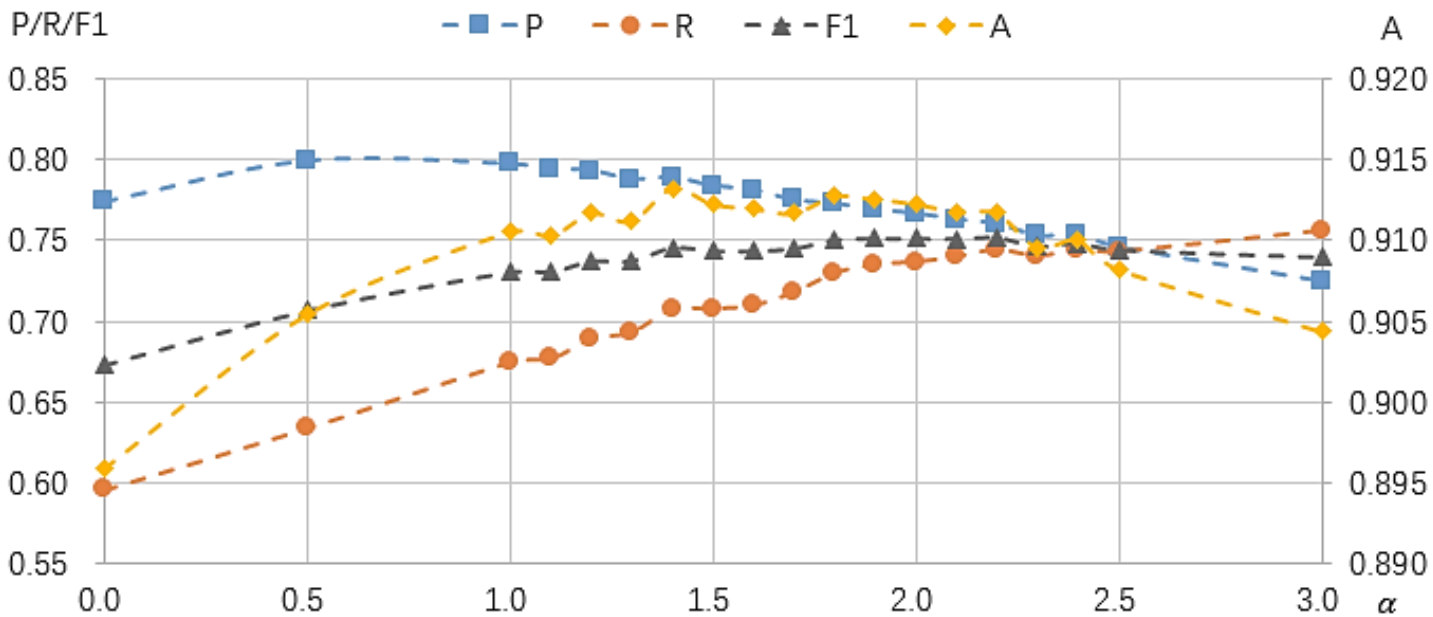
– CRF: reproduce (H. Zhang, 2009) with CRF++

Experiments

- Parameter Selection

- Weight of transition score α

Fix $beam = 50$, $\beta = 0 \Rightarrow$ highest accuracy when $\alpha = 1.4$

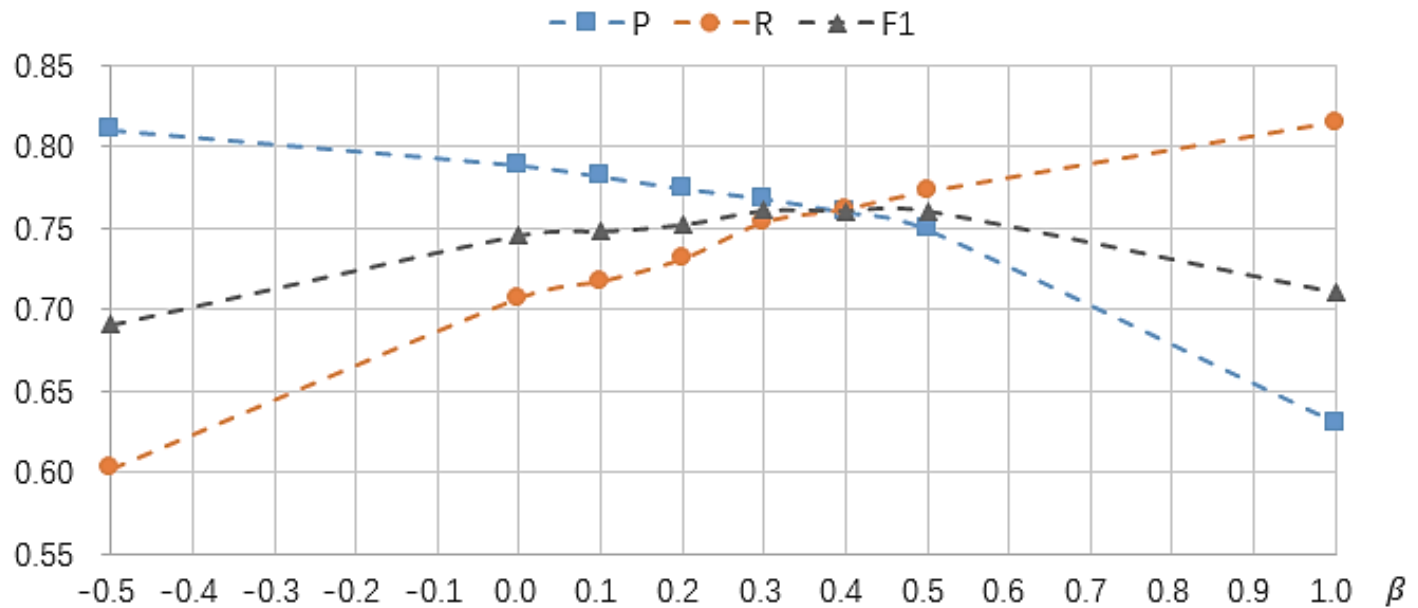


Experiments

– Parameter Selection

– **Weight of length penalty β**

Fix $beam = 50$, $\alpha = 1.4 \Rightarrow$ highest F1 when $\beta = 0.3$

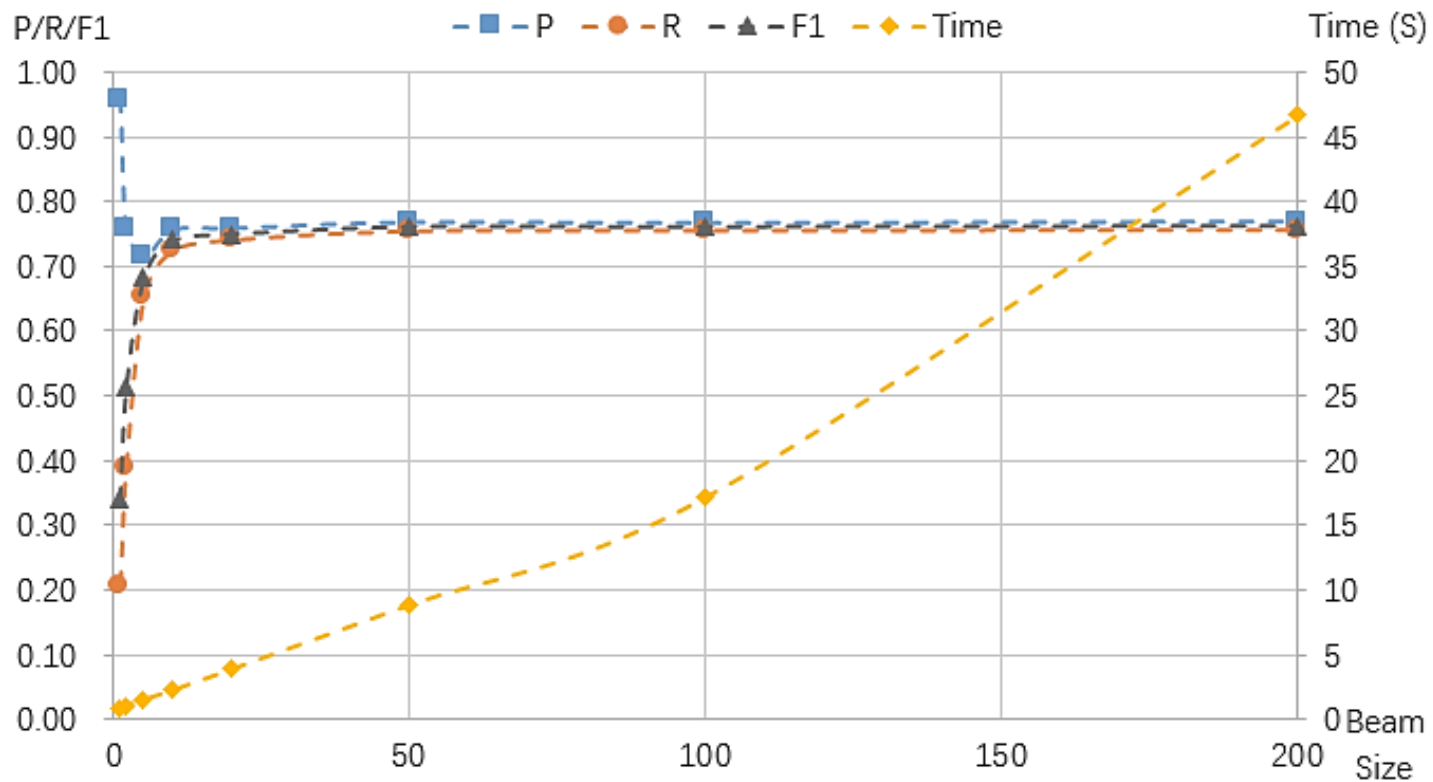


Experiments

– Parameter Selection

– Size of beam

Fix $\alpha = 1.4$, $\beta = 0.3$



Experiments

– Results

Models	Test Set 1			Test Set 2		
	P	R	F1	P	R	F1
CLM6	0.782	0.577	0.664	0.728	0.636	0.679
CLM6+0.65LP	0.749	0.673	0.709	0.686	0.721	0.703
CLM1+0.65LP	0.721	0.733	0.727	0.639	0.769	0.698
CRF	0.816	0.662	0.731	0.786	0.710	0.746
GRU-RNN	0.755	0.677	0.714	0.734	0.633	0.680
GRU-RNN+1.4ST[50]	0.785	0.720	0.751	0.760	0.719	0.739
GRU-RNN+1.4ST+0.3LP[50]	0.761	0.746	0.753	0.735	0.753	0.744
GRU-RNN+1.4ST+0.3LP[200]	0.763	0.748	0.755	0.738	0.756	0.747

Conclusion

- **Bi-directional GRU-RNN** outperform CRF and RNNLM methods.
- Integrations of **transition score** and **length penalty** are effective.
- Future works
 - Classify texts by ages and train sub-models
 - Segmentation → punctuation

Thank you!

by Boli Wang
me@bo-li.wang

