

Trans-dimensional Random Fields (TDRF) for Sequence Modeling

We present the potential of applying random fields for sequence modeling, demonstrated by its success in language modeling.

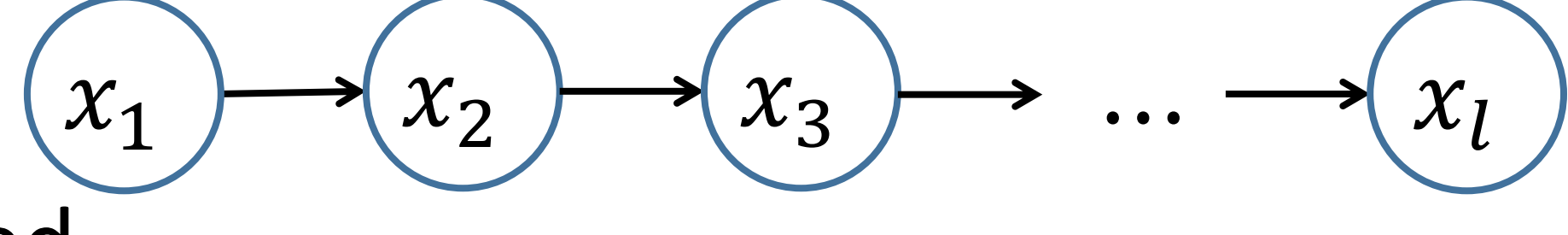
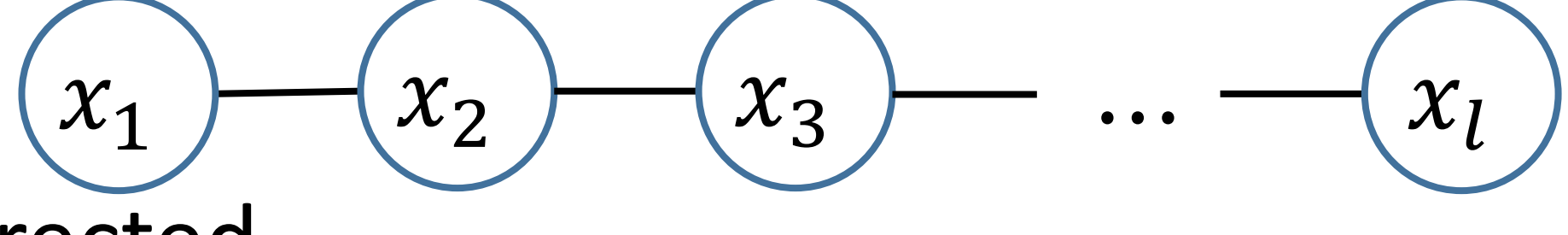
Bin Wang¹

Zhijian Ou¹

Zhiqiang Tan²

¹ Dept. of Electronic Engineering, Tsinghua Univ., China

² Dept. of Statistics, Rutgers Univ., USA

State-of-the-art LMs – Review	TDRF LMs – Motivation
<ul style="list-style-type: none"> Dominant: Conditional approach $p(x_1, x_2, \dots, x_l) = \prod_{i=1}^l p(x_i x_1, \dots, x_{i-1})$ <ul style="list-style-type: none"> N-gram LMs Neural network LMs $p(x_i = k x_1, \dots, x_{i-1}) \approx \frac{w_k^T \phi[x_1, \dots, x_{i-1}]}{\sum_{k=1}^V w_k^T \phi[x_1, \dots, x_{i-1}]}, w_k \in R^h$ <ul style="list-style-type: none"> ⊗ Computational expensive in both training and testing¹ e.g. lexicon size $V = 10k \sim 100k$, embedding dim $h = 250$ <p>¹ Partly alleviated by using un-normalized models, e.g. through noise contrastive estimation training.</p>	$p(x_1, x_2, \dots, x_l) = ?$ <ul style="list-style-type: none"> Dominant: Conditional approach / Directed  <ul style="list-style-type: none"> Alternative: Random field approach / Undirected  <ul style="list-style-type: none"> ⊗ Model training is difficult. ⊙ Capture bidirectional context for language cognition. <p>The cat is <u>on</u> the table. The cat is <u>in</u> the house.</p> <ul style="list-style-type: none"> ⊙ Breakthrough in training with a number of innovations Fixed-dim (e.g. image) -> Trans-dim (sequential modeling)

TDRF LMs – Model Definition	WSME vs TDRF
<ul style="list-style-type: none"> Features $(f_i, i = 1, 2, \dots, F)$ can be defined flexibly. Each feature brings a contribution to the sentence probability. $p(x; \lambda) = \frac{1}{Z(\lambda)} \exp\left(\sum_{i=1}^F \lambda_i f_i(x)\right), x \triangleq (x_1, x_2, \dots, x_l)$ $f_i(x) = \begin{cases} 1, & \text{'meeting on DAY-OF-WEEK' appears in } x \Rightarrow \lambda_i \text{ is activated} \\ 0, & \text{Otherwise} \Rightarrow \lambda_i \text{ is removed} \end{cases}$ <ul style="list-style-type: none"> ⊙ More flexible features, beyond the n-gram features, can be well supported in TDRF LMs. ⊙ Computational efficient in computing sentence probability for testing. <p>Jelinek 1995: put language back into language modeling</p>	<ul style="list-style-type: none"> Whole-sentence maximum entropy (WSME) (Rosenfeld, Chen, Zhu 2001) $p(l, x^l; \lambda) = \frac{1}{Z(\lambda)} \exp[\lambda^T f(x^l)], x \triangleq (l, x^l), x^l \triangleq (x_1, x_2, \dots, x_l)$ $= \frac{Z_l(\lambda)}{Z(\lambda)} \cdot \frac{1}{Z_l(\lambda)} \cdot \exp[\lambda^T f(x^l)], Z_l(\lambda) = \sum_{x^l} \exp[\lambda^T f(x^l)]$ <p>A mixture distribution with unknown weights, which differ from each other greatly, e.g. 10^{40} ! Poor sampling → poor estimation of gradient → poor fitting</p> <ul style="list-style-type: none"> Trans-dimensional RF (TDRF) model $p(l, x^l; \lambda) = \pi_l \cdot \frac{1}{Z_l(\lambda)} \cdot \exp[\lambda^T f(x^l)], l = 1, \dots, m$ <p>Empirical length probabilities in the training data Serve as a control device to improve sampling from multiple distributions!</p>

TDRF LMs – Model Estimation	Experiments																																																																												
<ul style="list-style-type: none"> Maximum-likelihood training $\frac{\partial \text{LogLikelihood}}{\partial \lambda} = E_{\tilde{p}(x)}[f_i(x)] - E_{p(x; \lambda)}[f_i(x)] = 0$ <div style="display: flex; justify-content: space-around;"> <div style="border: 1px solid black; padding: 5px;">Expectation under empirical distribution $\tilde{p}(x)$</div> <div style="border: 1px solid black; padding: 5px;">Expectation under model distribution $p(x; \lambda)$</div> </div> <ul style="list-style-type: none"> Consider $p(l, x^l; \lambda, \zeta) \propto \pi_l \cdot \frac{1}{e^{\zeta_l}} \cdot \exp[\lambda^T f(x^l)]$ where ζ_l is hypothesized values of the true $\zeta_l^*(\lambda) = \log Z_l(\lambda)$. The marginal probability of length l is: $p(l; \lambda, \zeta) = \frac{\pi_l e^{-\zeta_l + \zeta_l^*(\lambda)}}{\sum_j \pi_j e^{-\zeta_j + \zeta_j^*(\lambda)}}$ Joint SA is used to find $\zeta_l^* = \zeta_l^*(\lambda^*)$ and λ^* that solves $\begin{cases} \pi_l = p(l; \lambda, \zeta), & l = 1, \dots, m \\ 0 = E_{\tilde{p}(x)}[f_i(x)] - E_{p(l, x^l; \lambda, \zeta)}[f_i(x)] \end{cases}$ 	<p>LM Training — Penn Treebank portion of WSJ corpus Test speech — WSJ'92 set, by rescoreing of 1000-best lists</p> <table border="1"> <thead> <tr> <th>model</th> <th>WER</th> <th>PPL (\pm std. dev.)</th> <th>#feat</th> </tr> </thead> <tbody> <tr> <td>KN4</td> <td>8.71</td> <td>295.41</td> <td>1.6M</td> </tr> <tr> <td>RNN</td> <td>7.96</td> <td>256.15</td> <td>5.1M</td> </tr> <tr> <td colspan="4">WSMEs (200c)</td> </tr> <tr> <td>w+c+ws+cs</td> <td>8.87</td> <td>$\approx 2.8 \times 10^{12}$</td> <td>5.2M</td> </tr> <tr> <td>w+c+ws+cs+cpw</td> <td>8.82</td> <td>$\approx 6.7 \times 10^{12}$</td> <td>6.4M</td> </tr> <tr> <td colspan="4">TDRFs (100c)</td> </tr> <tr> <td>w+c</td> <td>8.56</td> <td>268.25\pm3.52</td> <td>2.2M</td> </tr> <tr> <td>w+c+ws+cs</td> <td>8.16</td> <td>265.81\pm4.30</td> <td>4.5M</td> </tr> <tr> <td>w+c+ws+cs+cpw</td> <td>8.05</td> <td>265.63\pm7.93</td> <td>5.6M</td> </tr> <tr> <td>w+c+ws+cs+wsh+csh</td> <td>8.03</td> <td>276.90\pm5.00</td> <td>5.2M</td> </tr> <tr> <td colspan="4">TDRFs (200c)</td> </tr> <tr> <td>w+c</td> <td>8.46</td> <td>257.78\pm3.13</td> <td>2.5M</td> </tr> <tr> <td>w+c+ws+cs</td> <td>8.05</td> <td>257.80\pm4.29</td> <td>5.2M</td> </tr> <tr> <td>w+c+ws+cs+cpw</td> <td>7.92</td> <td>264.86\pm8.55</td> <td>6.4M</td> </tr> <tr> <td>w+c+ws+cs+wsh+csh</td> <td>7.94</td> <td>266.42\pm7.48</td> <td>5.9M</td> </tr> <tr> <td colspan="4">TDRFs (500c)</td> </tr> <tr> <td>w+c</td> <td>8.72</td> <td>261.02\pm2.94</td> <td>2.8M</td> </tr> <tr> <td>w+c+ws+cs</td> <td>8.29</td> <td>266.34\pm6.13</td> <td>5.9M</td> </tr> </tbody> </table>	model	WER	PPL (\pm std. dev.)	#feat	KN4	8.71	295.41	1.6M	RNN	7.96	256.15	5.1M	WSMEs (200c)				w+c+ws+cs	8.87	$\approx 2.8 \times 10^{12}$	5.2M	w+c+ws+cs+cpw	8.82	$\approx 6.7 \times 10^{12}$	6.4M	TDRFs (100c)				w+c	8.56	268.25 \pm 3.52	2.2M	w+c+ws+cs	8.16	265.81 \pm 4.30	4.5M	w+c+ws+cs+cpw	8.05	265.63 \pm 7.93	5.6M	w+c+ws+cs+wsh+csh	8.03	276.90 \pm 5.00	5.2M	TDRFs (200c)				w+c	8.46	257.78 \pm 3.13	2.5M	w+c+ws+cs	8.05	257.80 \pm 4.29	5.2M	w+c+ws+cs+cpw	7.92	264.86 \pm 8.55	6.4M	w+c+ws+cs+wsh+csh	7.94	266.42 \pm 7.48	5.9M	TDRFs (500c)				w+c	8.72	261.02 \pm 2.94	2.8M	w+c+ws+cs	8.29	266.34 \pm 6.13	5.9M
model	WER	PPL (\pm std. dev.)	#feat																																																																										
KN4	8.71	295.41	1.6M																																																																										
RNN	7.96	256.15	5.1M																																																																										
WSMEs (200c)																																																																													
w+c+ws+cs	8.87	$\approx 2.8 \times 10^{12}$	5.2M																																																																										
w+c+ws+cs+cpw	8.82	$\approx 6.7 \times 10^{12}$	6.4M																																																																										
TDRFs (100c)																																																																													
w+c	8.56	268.25 \pm 3.52	2.2M																																																																										
w+c+ws+cs	8.16	265.81 \pm 4.30	4.5M																																																																										
w+c+ws+cs+cpw	8.05	265.63 \pm 7.93	5.6M																																																																										
w+c+ws+cs+wsh+csh	8.03	276.90 \pm 5.00	5.2M																																																																										
TDRFs (200c)																																																																													
w+c	8.46	257.78 \pm 3.13	2.5M																																																																										
w+c+ws+cs	8.05	257.80 \pm 4.29	5.2M																																																																										
w+c+ws+cs+cpw	7.92	264.86 \pm 8.55	6.4M																																																																										
w+c+ws+cs+wsh+csh	7.94	266.42 \pm 7.48	5.9M																																																																										
TDRFs (500c)																																																																													
w+c	8.72	261.02 \pm 2.94	2.8M																																																																										
w+c+ws+cs	8.29	266.34 \pm 6.13	5.9M																																																																										

Comparison	Computation efficient in training	Computation efficient in testing	Bidirectional context	Flexible features	Performance
N-gram LMs	✓	✓	✗	✗	✗
Neural network LMs	✗	✗	✗	✓	✓
TDRF LMs	✗	✓	✓	✓	✓