# Low-Rank Constraints for Fast Inference in Structured Models

Justin Chiu* and Yuntian Deng* and Alexander Rush

October 14, 2021

# Structured Models

▶ Explicitly model output associations

   ▶ Directly or through <u>latent variables</u>

▶ Focus on combinatorially large latent <u>discrete structures</u>

   ▶ Complementary to distributed representations

# Scaling Structured Models

- ▶ Prior work demonstrated: Size ⬆ Performance ⬆
  - ▶ Hidden Markov Models (HMM)
  - ▶ Probabilistic Context-Free Grammars (PCFG)

- ▶ Prior work scaled via
  - ▶ Sparsity for HMMs[1]
  - ▶ Low-rank tensor decompositions for PCFGs[2]

- ▶ This work: low-rank matrix constraints
  - ▶ More general
  - ▶ Less speedup

[1]Chiu and Rush, *Scaling Hidden Markov Language Models*.
[2]Yang, Zhao, and Tu, 'PCFGs Can Do Better: Inducing Probabilistic Context-Free Grammars with Many Symbols'.

# Inference as Matrix-Vector Products

- Inference: sequence of matrix-vector products

- Speed up via fast mat-vecs

- Applies to a large family of structured models

# Fast Matrix-Vector Products

- Mat-vecs take $O(L^2)$ computation

- Various fast methods
    - Sparsity (nnz entries)
    - Fast Fourier Transform ($L \log L$)
    - <u>Low-Rank factorization</u> ($LR$)

- Connected to work in efficient attention and low-dimensional kernel approximations[3]

---

[3]Choromanski et al., *Rethinking Attention with Performers*; Peng et al., *Random Feature Attention*; Blanc and Rendle, *Adaptive Sampled Softmax with Kernel Based Sampling*.

# Roadmap

▶ Speeding up HMM inference

▶ Speeding up PCFG inference

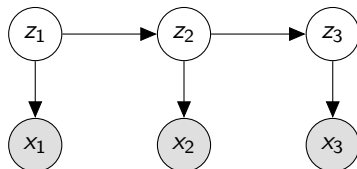▶ Generalization to hypergraph inference

▶ Experiments

# Two Examples

some text here some text here
some text here some text here
some text here

- ▶ Blah

# Hidden Markov Models (HMMs)

For times $t$, model states $z_t \in [Z]$, and tokens $x_t \in [X]$,



with joint distribution

$$p(x, z) = \prod_t p(x_t \mid z_t) p(z_t \mid z_{t-1})$$

## Inference

Given observed $x = (x_1, \ldots, x_T)$ We wish to maximize

$$p(x) = \sum_{z_1} \cdots \sum_{z_T} p(x, z) = \alpha_1^\top \Lambda_2 \Lambda_3 \cdots \Lambda_T \mathbf{1},$$

where we have the

| | |
|---|---|
| start, | $[\alpha_1]_{z_1} = p(x_1 \mid z_1) p(z_1),$ |
| and transition operators, | $[\Lambda_t]_{z_{t-1}, z_t} = p(x_t \mid z_t) p(z_t \mid z_{t-1})$ |

# Inference: Backward Algorithm

▶ Performing multiplications from right to left

$$p(x) = \alpha_1^\top (\Lambda_2(\Lambda_3 \mathbf{1}))$$

▶ Recursively

$$\beta_t = \Lambda_t \beta_{t+1}$$

▶ Requires $O(TZ^2)$ operations in total

## Low-Rank Factorization

Factor matrices $\Lambda \in \mathbb{R}^{Z \times Z}$ into product of $U, V \in \mathbb{R}^{Z \times R}$



resulting in two matrix-vector products of cost $O(ZR)$ each

# Hypergraph Marginalization

►

# Hypergraph Marginalization Algorithm

# Experiments

# Experiments

- Language modeling on PTB

- Feature map $\phi(U) = \exp(UW)$, with learned $W \in \mathbb{R}^{R \times R}$

- Baseline: Softmax HMM

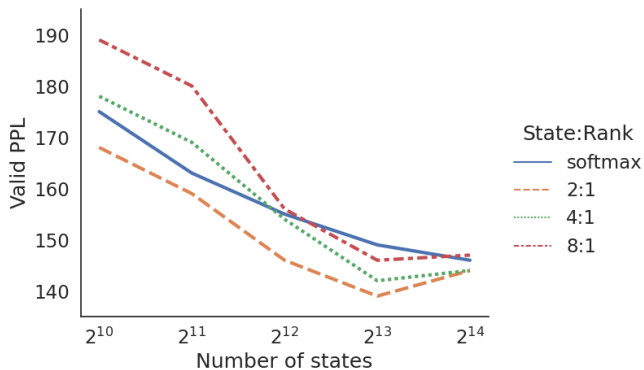# HMM Accuracy

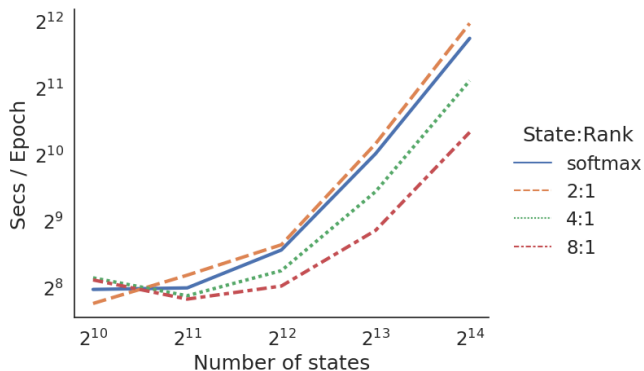| Model | Val | Test |
|---|---|---|
| AWD-LSTM | 60.0 | 57.3 |
| VL-HMM | 128.6 | 119.5 |
| HMM | 144.3 | 136.8 |
| LHMM | 141.4 | 131.8 |

# HMM Speed vs Accuracy Frontier

# HMM Accuracy vs Rank

# HMM Speed vs Rank

# HMM Music Results

| Model | Nott | Piano | Muse | JSB |
|---|---|---|---|---|
| RNN-NADE | 2.31 | **7.05** | **5.6** | 5.19 |
| R-Transformer | **2.24** | 7.44 | 7.00 | 8.26 |
| LSTM | 3.43 | 7.77 | 7.23 | 8.17 |
| LV-RNN | 2.72 | 7.61 | 6.89 | **3.99** |
| SRNN | 2.94 | 8.20 | 6.28 | 4.74 |
| TSBN | 3.67 | 7.89 | 6.81 | 7.48 |
| HMM | 2.43 | 8.51 | 7.34 | 5.74 |
| LHMM | 2.60 | 8.89 | 7.60 | 5.80 |

# HSMM Results

# PCFG Results

# Citations

📄 Blanc, Guy and Steffen Rendle. *Adaptive Sampled Softmax with Kernel Based Sampling*. 2018. arXiv: 1712.00527 [cs.LG].

📄 Chiu, Justin T. and Alexander M. Rush. *Scaling Hidden Markov Language Models*. 2020. arXiv: 2011.04640 [cs.CL].

📄 Choromanski, Krzysztof et al. *Rethinking Attention with Performers*. 2021. arXiv: 2009.14794 [cs.LG].

📄 Peng, Hao et al. *Random Feature Attention*. 2021. arXiv: 2103.02143 [cs.CL].

📄 Yang, Songlin, Yanpeng Zhao, and Kewei Tu. 'PCFGs Can Do Better: Inducing Probabilistic Context-Free Grammars with Many Symbols'. In: *CoRR* abs/2104.13727 (2021). arXiv: 2104.13727. URL: https://arxiv.org/abs/2104.13727.