

Scaling Hidden Markov Language Models

March 3, 2021

Abstract

Modern methods for language models based purely on neural networks are opaque and difficult to analyze, while alternative methods based on probabilistic graphical models are interpretable but not performant. We hypothesize that probabilistic methods are not as performant neural methods because they have neither been scaled to similar sizes as their modern neural counterparts nor taken advantage of modern parallel hardware. A significant barrier to scaling in classical methods is inference, which typically scales superlinearly in the size of the model. In order to make probabilistic models performant while maintaining interpretability, we present techniques for circumventing the heavy computational costs of inference. We first focus on scaling one of the simplest models, hidden Markov models, then show our techniques generalize to more complex models, such as probabilistic context-free grammars.

1 Introduction

Neural network-based generative models have led to progress in difficult tasks such as language modeling and machine translation. However, this progress comes at the cost of interpretability. Neural networks are flexible function approximators, but that flexibility results in opaque models that are difficult to analyze. As opposed to post-hoc analysis of models [8], we instead propose to explore the space of models that are interpretable from the start. In particular, we focus on probabilistic graphical models, which allow for the execution of arbitrary probabilistic queries. These queries include the ability to calculate conditional and marginal probabilities given evidence. This interpretability comes at a cost: Performing these queries via inference is expensive in both time and space. The computational complexity of inference is a major impediment to scaling graphical models to large sizes, which we hypothesize prevents them from reaching the performance of neural networks. We seek to answer the question of whether we can have graphical models that are both performant and interpretable.

We use language modeling as a benchmark task, as the complex and sometimes long-range phenomena in natural language provides a good testbed. Language model benchmarks are currently dominated by autoregressive neural models, which makes it difficult to analyze what properties of language the models are actually capturing.

In this work, we investigate scaling probabilistic graphical models which explicitly reason about latent variables. We focus on the one of the simplest latent variable models, the hidden Markov model (HMM). Despite the simplicity of HMMs, they are still computationally expensive to scale due to the cost of inference. In order to scale HMMs, we utilize a kernel-based generalized softmax that greatly reduces the complexity of inference. Before diving into the generalized softmax, we first review HMMs.

2 Hidden Markov Models for Language Modeling

Hidden Markov models (HMMs) have a rich history in natural language processing. They have been used for a series of tasks, including speech recognition [5], part of speech tagging [3], and word alignment in machine translation [9]. They have also been used as components in neural models, such as in attention [7]. We focus on applying HMMs to the task of language modeling, where the goal is to model the sequence of tokens in a sentence $x = (x_1, \dots, x_T)$.

HMMs are one of the simplest latent variable models for language modeling with a per-token latent variable. The simplicity is a result of the strong conditional independence assumptions: Every timestep has a single discrete state variable used to represent context, which is then the only information used to emit a token. Additionally, the next state is a function of only the previous state. These independence assumptions result in a bottleneck where all information must flow through the

single state variable at each timestep, causing performance to be limited by the number of states. Limiting dependencies in this way also gives the model interpretability.

Formally, HMMs have the following generative process: for each $t \in [T]$ timestep, first choose a latent state $z_t \in \mathcal{Z}$, where $|\mathcal{Z}|$ is the number of latent states, then choose a token to emit $x_t \in \mathcal{X}$. This defines the joint distribution:

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

with transitions $p(z_t | z_{t-1})$, emissions $p(x_t | z_t)$, and a distinguished start state $z_0 = S$.

Training an HMM requires marginalizing over the unobserved sequence $z = (z_1, \dots, z_T)$ in order to obtain the evidence, $p(x) = \sum_z p(x, z)$, via the forward algorithm. The forward algorithm can be written as a sequence of matrix-vector multiplications: Let the transition operators $\Lambda_t \in \mathbb{R}^{|\mathcal{Z}| \times |\mathcal{Z}|}$ for $t \in [2, \dots, T]$, with entries

$$[\Lambda_t]_{z_{t-1}, z_t} = p(x_t | z_t) p(z_t | z_{t-1}), \quad (2)$$

and the start vector, $\boldsymbol{\pi} \in \mathbb{R}^{|\mathcal{Z}|}$, given by

$$[\boldsymbol{\pi}]_{z_1} = p(x_1 | z_1) p(z_1 | z_0 = S). \quad (3)$$

The evidence is then given by

$$p(x) = \boldsymbol{\pi}^\top \Lambda_2 \cdots \Lambda_T \mathbf{1}, \quad (4)$$

where $\mathbf{1} \in \mathbb{R}^{|\mathcal{Z}| \times 1}$ is the column vector of all ones.

On a serial machine, the cost of computing each matrix-vector product in the above equation takes time $O(|\mathcal{Z}|^2)$, resulting in a total running time of $O(T|\mathcal{Z}|^2)$ for the forward algorithm. The quadratic dependence on the number of states precludes scaling to extremely large state spaces.

In the following sections we will introduce techniques that aim to reduce the dependence of the time and space complexity of inference on the size of the state space.

3 Generalized Softmax with Kernels

Parameterization of the probability distributions in a model affects both computational performance as well as generalization. Recent work on efficient attention in neural models shows that careful choice of parameterization of the softmax kernel, the main component of attention, results in large computational gains at little cost in accuracy [1, 4]. In this section we will cover generalized softmax with kernels and its use in efficient inference.

Softmax is commonly used to parameterize conditional distributions in neural nets, and has origins in smooth approximations of argmax as well as the conditional max entropy problem []. Focusing on the transition distribution $p(z_t | z_{t-1})$, the softmax parameterization is as follows:

$$p(z_t | z_{t-1}) = \frac{\exp(\mathbf{u}_{z_{t-1}}^\top \mathbf{v}_{z_t})}{\sum_{z'} \exp(\mathbf{u}_{z_{t-1}}^\top \mathbf{v}_{z'})},$$

with $\mathbf{u}_{z_{t-1}}, \mathbf{v}_{z_t} \in \mathbb{R}^d$ vector-space embeddings of values of z_{t-1}, z_t .¹ The numerator, $\exp(\mathbf{u}_{z_{t-1}}^\top \mathbf{v}_{z_t})$, is the exponential kernel [6].

Although the parameterization of softmax is specific to the exponential kernel, softmax can be generalized with arbitrary nonnegative kernels $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$. Replacing the exponential kernel with the nonnegative kernel K yields generalized softmax,

$$p(z_t | z_{t-1}) = \frac{K(\mathbf{u}_{z_{t-1}}, \mathbf{v}_{z_t})}{\sum_{z'} K(\mathbf{u}_{z_{t-1}}, \mathbf{v}_{z'})}, \quad (5)$$

where we again have embeddings $\mathbf{u}_{z_{t-1}}, \mathbf{v}_{z_t} \in \mathbb{R}^d$.²

Kernels have a long history in machine learning, where they are used to extend linear classification models to nonlinear feature spaces, such as in support vector machines or regression. In those settings, kernels are intuitively used to measure the similarity between two data points, hinging on a connection to an inner product in a reproducing kernel Hilbert space. In the generalized softmax setting, kernels instead parameterize conditional distributions by measuring the propensity of one random variable taking a particular value given the value of another random variable. However, the connection to inner products in reproducing kernel Hilbert spaces is still useful. In particular, the

¹ Softmax also has a temperature parameter that rescales the dot product, which we omit.

² The term generalized softmax is inspired by generalized linear models, which generalized linear models to non-gaussian errors with a link function, similar to the use of kernels in generalized softmax.

connection guarantees the existence of a feature map $\phi : \mathbb{R}^d \rightarrow \mathbb{F}$, a map from the input space \mathbb{R}^d to the feature space \mathbb{F} , which allows us to write generalized softmax in matrix form:

$$p(z_t \mid z_{t-1}) = [\text{diag}(\mathbf{d})\phi(U)\phi(V)^\top]_{z_{t-1}, z_t}, \quad (6)$$

where we have the stacked embeddings $\phi(U), \phi(V) \in \mathbb{R}^{|\mathcal{Z}| \times \dim(\mathbb{F})}$ and the normalizing constants $\mathbf{d} \in \mathbb{R}^{|\mathcal{Z}|}$,

$$\begin{aligned} \phi(U) &= [\phi(\mathbf{u}_1)^\top, \phi(\mathbf{u}_2)^\top, \dots] \in \mathbb{R}^{|\mathcal{Z}| \times \dim(\mathbb{F})}, \\ \phi(V) &= [\phi(\mathbf{v}_1)^\top, \phi(\mathbf{v}_2)^\top, \dots] \in \mathbb{R}^{|\mathcal{Z}| \times \dim(\mathbb{F})}, \\ [\mathbf{d}]_{z_{t-1}} &= \phi(\mathbf{u}_{z_{t-1}})^\top \sum_z \phi(\mathbf{v}_z). \end{aligned}$$

We will rely on the dimension of the feature space being small relative to the number of states, $\dim(\mathbb{F}) < |\mathcal{Z}|$, in order to improve the time and space complexity of inference.

4 Inference with Generalized Softmax

Unfortunately, in the case of softmax, which relies on the exponential kernel, the feature space \mathbb{F}_{exp} has infinite dimension. This can be seen from the Taylor expansion:

$$\exp(\mathbf{u}^\top \mathbf{v}) = \sum_{n=0}^{\infty} \frac{(\mathbf{u}^\top \mathbf{v})^n}{n!} = \sum_{n=0}^{\infty} \frac{(\sum_{i=1}^d [\mathbf{u}]_i [\mathbf{v}]_i)^n}{n!} = \sum_{n=0}^{\infty} \frac{\sum_{\mathbf{j} \in [d]^n} (\prod_{i=1}^n [\mathbf{u}]_{[\mathbf{j}]_i}) (\prod_{i=1}^n [\mathbf{v}]_{[\mathbf{j}]_i})}{n!},$$

where we have expanded the terms of the numerator by summing over all subsets \mathbf{j} of size n elements of $[d]$, including repeats and different orderings. This yields a feature map of

$$[\phi(\mathbf{u})]_{n, \mathbf{j}} = \frac{\prod_{i=1}^n [\mathbf{u}]_{[\mathbf{j}]_i}}{\sqrt{n!}}$$

for all $n \in \mathbb{N}$.³ We can therefore express the exponential kernel as an inner product of infinite dimensional vectors [2],

$$\exp(\mathbf{u}^\top \mathbf{v}) = \sum_{n=0}^{\infty} \sum_{\mathbf{j} \in [d]^n} [\phi(\mathbf{u})]_{n, \mathbf{j}} [\phi(\mathbf{v})]_{n, \mathbf{j}} = \langle \phi(\mathbf{u}), \phi(\mathbf{v}) \rangle,$$

where the dimension of the feature space $\dim(\mathbb{F}_{\text{exp}}) = \dim(\phi(\mathbf{u})) = \infty$. Since the infinite dimensional feature space of the exponential kernel is larger than the finite $|\mathcal{Z}|$, recent work has instead turned to approximations of the exponential kernel based on random Fourier features [1, 4, 6]. These approaches use a finite-dimensional unbiased estimator of the exponential kernel.

We propose to take a different strategy, where we directly learn a feature map for the generalized softmax while disregarding the explicit functional form of the kernel. This allows us to use the matrix representation of generalized softmax to efficiently perform inference.

Previously, Equation 4 shows how to compute the evidence, $p(x)$, with a series of matrix-vector products. This was done by first defining the transition operator, Λ_t as in Equations 2 and 3. We can perform a similar procedure with the transition distribution using generalized softmax. With a generalized softmax parameterization of the transition distribution and a feature space \mathbb{F}_ϕ with feature map ϕ and finite dimension $\dim(\mathbb{F})$, we have the transition operators for all $t \in [2, \dots, T]$,

$$\Lambda_t = \underbrace{\text{diag}(\mathbf{d})\phi(U)}_{L_t} \underbrace{\phi(V)^\top \text{diag}(\mathbf{o}_t)}_{R_t}, \quad (7)$$

where the entries of $\mathbf{o}_t \in \mathbb{R}^{|\mathcal{Z}|}$ are given by

$$[\mathbf{o}_t]_{z_t} = p(x_t \mid z_t),$$

the emission probability and $\mathbf{d} \in \mathbb{R}^{|\mathcal{Z}|}$ contains the normalizing constant for each state. We can interpret $L_t \in \mathbb{R}^{|\mathcal{Z}| \times \dim \mathbb{F}}$ as a projection into feature space, and $R_t \in \mathbb{R}^{\dim(\mathbb{F}) \times |\mathcal{Z}|}$ as a projection back into state space. The evidence can then be computed using L_t and R_t as follows:

$$p(x) = \boldsymbol{\pi} \Lambda_2 \cdots \Lambda_T \mathbf{1} = \boldsymbol{\pi} (L_2 R_2) \cdots (L_T R_T) \mathbf{1} = \boldsymbol{\pi} L_2 R_2 \cdots L_T R_T \mathbf{1}, \quad (8)$$

where we have substituted Equation 7 into the computation of the evidence and applied the associative property of matrix multiplication. When performed from left to right, the matrix vector

³ While this is one possible feature map, there may exist multiple valid feature maps.

products now take $O(|\mathcal{Z}| \dim(\mathbb{F}))$ time, rather than $O(|\mathcal{Z}|^2)$, with the overall complexity of inference now linear in the number of states, $O(T|\mathcal{Z}| \dim(\mathbb{F}))$, with generalized softmax.

We parameterize the feature map $\phi(\mathbf{u}) = W_\phi \mathbf{u}$, with a linear projection $W_\phi \in \mathbb{R}^{\dim(\mathbb{F}) \times d}$, and train all parameters, including W_ϕ and the parameters of the transition and emission distributions,⁴ by optimizing the evidence with gradient ascent.

4.1 Connection to Kernel Mean Embeddings and Kernel Belief Propagation

(Not done yet, tbd) Many probabilistic queries can be written as conditional expectations. In particular, we can write marginalization in an HMM as a conditional expectation:

$$p(x) = \sum_z p(x, z) = \sum_z p(z) p(x | z) = \mathbb{E}_z [p(x | z)].$$

Decomposing the expectation over time and examining a single timestep, we have

$$p(x_t | x_{<t}) = \sum_{z_t} p(x_t | z_t) p(z_t | x_{<t}).$$

The second term, $p(z_t | x_{<t})$, can be further decomposed using the generative process of the HMM to

$$p(z_t | x_{<t}) = \sum_{z_{t-1}} p(z_t | z_{t-1}) p(z_{t-1} | x_{<t}),$$

altogether yielding the equation

$$p(x_t | x_{<t}) = \sum_{z_t} p(x_t | z_t) \sum_{z_{t-1}} p(z_t | z_{t-1}) p(z_{t-1} | x_{<t}) = \mathbb{E}_{z_{t-1}, z_t | z_{t-1}} [p(x_t | z_t)] p(z_{t-1} | x_{<t}).$$

We can use the matrix form of generalized softmax to speed up this computation. Let $[\mathbf{f}_t]_{z_t} = p(x_t | z_t)$, $[\boldsymbol{\alpha}_t]_{z_{t-1}} = p(z_{t-1} | x_{<t})$ in

$$\mathbb{E}_{z_t | x_{<t}} [p(x_t | z_t)] = (\boldsymbol{\alpha}_t \circ \mathbf{d})^\top \phi(U) \phi(V)^\top \mathbf{f}_t,$$

where $\mathbf{d} = \frac{1}{\phi(U) \sum_z \phi(\mathbf{v}_z)} \in \mathbb{R}^{|\mathcal{Z}|}$.

5 Kernel Approximations

Move to appendix

Taylor approximation, generating functions
limitations?

5.1 Bochner Theorem

5.2 Nystrom

6 Kernelized Inference

7 Generalization: Kernelized Belief Propagation

8 Spectral Methods?

References

- [1] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy Colwell, and Adrian Weller. Rethinking attention with performers, 2020.
- [2] Andrew Cotter, Joseph Keshet, and Nathan Srebro. Explicit approximations of the gaussian kernel. *CoRR*, abs/1109.4603, 2011. URL <http://arxiv.org/abs/1109.4603>.

⁴ We give the full parameterization of all distributions in the appendix.

- [3] Bernard Merialdo. Tagging English text with a probabilistic model. *Computational Linguistics*, 20(2):155–171, 1994. URL <https://www.aclweb.org/anthology/J94-2001>.
- [4] Hao Peng, Nikolaos Pappas, Dani Yogatama, Roy Schwartz, Noah Smith, and Lingpeng Kong. Random feature attention. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=QtTKTdVrFBB>.
- [5] Lawrence R. Rabiner. *A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition*, page 267–296. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1990. ISBN 1558601244.
- [6] Ankit Singh Rawat, Jiecao Chen, Felix Xinnan X Yu, Ananda Theertha Suresh, and Sanjiv Kumar. Sampled softmax with random fourier features. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Álché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32, pages 13857–13867. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/e43739bba7cdb577e9e3e4e42447f5a5-Paper.pdf>.
- [7] Shiv Shankar and Sunita Sarawagi. Posterior attention models for sequence to sequence learning. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bk1tNhC9FX>.
- [8] Ian Tenney, Dipanjan Das, and Ellie Pavlick. Bert rediscovers the classical nlp pipeline. In *Association for Computational Linguistics*, 2019. URL <https://arxiv.org/abs/1905.05950>.
- [9] Stephan Vogel, Hermann Ney, and Christoph Tillmann. Hmm-based word alignment in statistical translation. In *Proceedings of the 16th Conference on Computational Linguistics - Volume 2, COLING ’96*, page 836–841, USA, 1996. Association for Computational Linguistics. doi: 10.3115/993268.993313. URL <https://doi.org/10.3115/993268.993313>.