

# Instructions for EMNLP 2020 Proceedings

## Anonymous EMNLP submission

### Abstract

todo

## 1 Introduction

Hidden Markov models (HMMs) are fundamental to the field of machine learning, providing one of the simplest latent variable models for sequential data. HMMs have been applied to many time-series problems, such as those that can be found in bioinformatics, reinforcement learning, and natural language processing. For the task of language modeling, HMMs were not considered until recently (Krakovna and Doshi-Velez, 2016; Dedieu et al., 2019). Instead, language models modeled the observed words using n-gram models, feedforward neural networks (Bengio et al., 2003), and eventually recurrent neural networks (Mikolov et al., 2010; Zaremba et al., 2014; Merity et al., 2017) or transformers (Radford et al., 2019).

We can make three observations about language modeling progress. The first is that parameterization is important for learning. The feedforward neural network of Bengio et al. (2003) had the same independence assumptions as an n-gram model, yet achieved better performance due to its parameterization. The second is that overparameterization is helpful as well. The work in Zaremba et al. (2014) demonstrated that increasing the size of a recurrent neural network led to marked improvements in performance. The third is that regularization is necessary with small datasets and overparameterization. Experiments by Merity et al. (2017) showed that the performance a recurrent neural network could vary greatly given various training tricks.

In this work, we apply the observations taken from improving fully observed language models to HMMs. We scale HMM training to a large num-

ber of states by introducing the following contributions:

1. A neural parameterization that allows us to compute conditional distributions only as needed while avoiding gradient sparsity.
2. Sparsity constraints that allow for efficient conditional inference in large HMMs.
3. A dropout variant that improves both performance and computational overhead during training.

On two small language modeling datasets and a supervised part of speech task, we found the HMM outperforms n-gram models and performs comparably to neural network counterparts.

## 2 Related work

Prior work has demonstrated the benefits of neural parameterization of conditional distributions. For HMMs, Tran et al. (2016) demonstrated improvements in POS induction with a neural parameterization of an HMM, while Kim et al. (2019) demonstrated improvements in grammar induction with a probabilistic context free grammar. Both of these works used latent variables with relatively small state spaces, as the goal of both was structure induction rather than language modeling itself. We extend the neural parameterization to much larger state space models.

We also draw inspiration from the experiments by Dedieu et al. (2019), who proposed to introduce sparsity constraints in the emission distribution of HMMs in order to make conditional inference tractable in large state spaces. Dedieu et al. (2019) trained a 30k state HMM on character-level language modeling by constraining every state to emit only a single character type. This particular

constraint is problematic for language modeling at the word level, where the vocabulary size is much larger. We build on their work by proposing a sparsity constraint based on Brown clustering (Brown et al., 1992) which allows us to extend their work to vocabularies that are larger than the state space.

Other investigations into improving the performance of HMMs involved relaxing independence or modeling assumptions (Buys et al., 2018) to resemble a recurrent neural network, or through model combination with a recurrent neural network (Krakovna and Doshi-Velez, 2016). There are also extensions of HMMs, such as factorial HMMs Ghahramani and Jordan (1997); Nepal and Yates (2013) and context free grammars (Kim et al., 2019). We leave scaling more expressive models to large state spaces for future work, and focus on scaling the basic HMM.

### 3 Background

We are interested in learning a distribution over sequences of tokens  $\mathbf{x} = \langle x_0, \dots, x_T \rangle$ , with each token  $x_t$  an element of the finite vocabulary  $\mathcal{X}$ .

#### 3.1 Hidden Markov Models

Hidden Markov Models (HMMs) specify a joint distribution over observed tokens  $\mathbf{x}$  and discrete latent states  $\mathbf{z} = \langle z_0, \dots, z_T \rangle$ , with each  $z_t$  from finite set  $\mathcal{Z}$ . The model is defined by the following generative process: For every time step  $t \in \{0, \dots, T\}$ , choose a state given the previous state  $z_t | z_{t-1}$  from the transition distribution  $p(z_t | z_{t-1})$ , where the first state  $z_0$  is chosen from the starting distribution  $p(z_0)$ . Then choose a token given the current state  $x_t | z_t$  from the emission distribution  $p(x_t | z_t)$ . This yields the joint distribution

$$p(\mathbf{x}, \mathbf{z}) = p(z_0)p(x_0 | z_0) \prod_{t=1}^T p(x_t | z_t)p(z_t | z_{t-1}) \quad (1)$$

The distributions are parameterized as follows

$$\begin{aligned} p(z_0) &\propto e^{\phi_{z_0}} \\ p(z_t | z_{t-1}) &\propto e^{\phi_{z_t z_{t-1}}} \\ p(x_t | z_t) &\propto e^{\phi_{x_t z_t}} \end{aligned} \quad (2)$$

where each  $\phi_{z_0}, \phi_{z_t z_{t-1}}, \phi_{x_t z_t} \in \mathbb{R}$  may be parameterized by a scalar or a neural network.

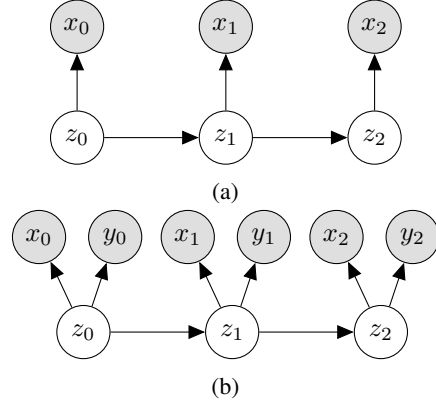


Figure 1: The HMMs for (a) language modeling and (b) part-of-speech tagging. We have the words  $x_t$ , the states  $z_t$ , and the tags  $y_t$ .

A scalar parameterization would result in  $O(|\mathcal{Z}|^2 + |\mathcal{Z}||\mathcal{X}|)$  parameters, since the transition and emission matrices must be explicitly represented. A compact neural parameterization may take as few as  $O(H(|\mathcal{Z}| + |\mathcal{X}|))$  parameters, where  $H$  is the dimension of the state and word embeddings. Such a parameterization represents the states and words with embeddings, and implicitly parameterizes the transition and emission distributions by modeling those matrices as a function of the embeddings. We use a residual network as in Kim et al. (2019) to parameterize conditional distributions.

### 4 Model

For sequence modeling with an HMM, we must marginalize over the latent states  $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$  in order to train the model. This sum can be computed in time  $O(T|\mathcal{Z}|^2)$  via dynamic programming, which becomes prohibitive if the number of latent states  $|\mathcal{Z}|$  is large. However, in order to increase the capacity of HMMs, we would like to consider large  $\mathcal{Z}$ . We present two techniques to manage the computational complexity of marginalization in HMMs.

#### 4.1 Sparse Emission HMMs

We introduce a sparse emission constraint that allows us to efficiently perform conditional inference in large state spaces. Inspired by Cloned HMMs (Dedieu et al., 2019), we constrain each observation  $x$  to be emitted only by the states  $z \in \mathcal{C}_x \subset \mathcal{Z}$ :

$$p(x | z) \propto 1(z \in \mathcal{C}_x) e^{\phi_{zx}} \quad (3)$$

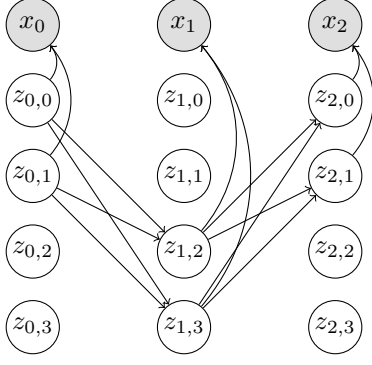


Figure 2: A depiction of the Brown emission constraints. The Brown emission constraint partitions the states such that each cluster is emitted by a disjoint set of states. Above is the trellis after conditioning on the observation sequence  $\mathbf{x}$  with hidden states  $\mathbf{z}$ . The arrows indicate transitions with nonzero probabilities. Rather than considering transition matrices of size  $4 \times 4$ , the emission constraint allows us to use transition matrices of size  $2 \times 2$  after conditioning on  $\mathbf{x}$ .

This allows us to perform marginalization as follows

$$\begin{aligned}
 p(\mathbf{x}) &= \sum_{z_0} p(z_0) p(x_0 | z_0) \sum_{z_1} p(z_1 | z_0) p(x_1 | z_1) \\
 &\quad \cdots \sum_{z_T} p(z_T | z_{T-1}) p(x_T | z_T) \\
 &= \sum_{z_0 \in \mathcal{C}_{x_0}} p(z_0) p(x_0 | z_0) \sum_{z_1 \in \mathcal{C}_{x_1}} p(z_1 | z_0) p(x_1 | z_1) \\
 &\quad \cdots \sum_{z_T \in \mathcal{C}_{x_T}} p(z_T | z_{T-1}) p(x_T | z_T)
 \end{aligned} \tag{4}$$

which takes  $O(T \max_x (|\mathcal{C}_x|)^2)$  computation. We choose sets  $\mathcal{C}_x$  such that  $\forall x, |\mathcal{C}_x| = k$  due to convenience of implementation.

We experiment with two constraint sets, both of which are specified as hyperparameters and remain fixed for a given model. The first constraint set is obtained by sampling subsets of states for each element of  $\mathcal{X}$  i.i.d. from a uniform distribution over subsets. The second constraint set uses Brown Clusters (Brown et al., 1992) to assign all words in a given Brown cluster the same subset of states. This constraint set is advantageous, as it admits further optimizations that we discuss below.

## 4.2 State Dropout

We introduce a variant of dropout called state dropout that operates on the start and transition

distributions. The goal of state dropout is to encourage usage of the full state space in HMMs as well as to reduce the complexity of marginalization in a manner identical to the sparse emission constraints.

State dropout samples a mask over states  $\mathbf{b}_C \in \{0, 1\}^{|\mathcal{C}|}$  for each partition  $\mathcal{C} \subset \mathcal{Z}$ , such that each sampled  $|\mathbf{b}_C| = n$ . We only apply this if the  $\mathcal{C}_x$  are disjoint partitions, as with the Brown constraint set.

Let  $\mathbf{b}$  be the concatenation of the  $\mathbf{b}_C$  for all partitions  $\mathcal{C}$ , such that  $b_z = 1$  if  $b_{Cz} = 1$ , where  $z \in \mathcal{C}$ . The start and transition distributions with dropout are given by

$$\begin{aligned}
 p(z_0) &\propto b_{z_0} e^{\phi_{z_0}} \\
 p(z_t | z_{t-1}) &\propto b_{z_t} b_{z_{t-1}} e^{\phi_{z_t z_{t-1}}}
 \end{aligned} \tag{5}$$

Marginalizing over  $\mathbf{z}$  with state dropout requires  $O(Tn^2)$  computation, where  $n$  is the number of nonzero elements of  $\mathbf{b}_C, \forall \mathcal{C}$ . We utilize state dropout on top of the Brown emission sparsity constraint, subsampling states of size  $n$  for each partition during training.

We also consider unstructured dropout with rate  $p$ , which samples elements of masks  $b_{z_0}, b_{z_t z_{t-1}} \stackrel{iid}{\sim} \text{Bern}(1 - p)$  for all  $z_0, z_t, z_{t-1} \in \mathcal{Z}$ , yielding

$$\begin{aligned}
 p(z_0) &\propto b_{z_0} e^{\phi_{z_0}} \\
 p(z_t | z_{t-1}) &\propto b_{z_t z_{t-1}} e^{\phi_{z_t z_{t-1}}}.
 \end{aligned} \tag{6}$$

Unstructured dropout results in sparsity patterns that are more difficult to take advantage of than state dropout, since there is variance in the number of nonzero elements in a mask. (Fix wording)

We find that models with state dropout generalize better than models with unstructured dropout, and additionally have computational benefits as described above.

## 4.3 Part of speech tagging

We extend the above HMM to part of speech (POS) tagging. In addition to the words  $\mathbf{x}$ , we model the tags  $\mathbf{y} = \langle y_0, \dots, y_T \rangle, y_t \in \mathcal{Y}$ . We introduce the following change to the HMM’s generative process: at every timestep  $t$ , the model chooses a state  $z_t | z_{t-1}$  and emits a word  $x_t | z_t$  and tag  $y_t | z_t$  independently given the state. This yields the fol-

lowing joint distribution

$$p(\mathbf{x}, \mathbf{y}, \mathbf{z}) = p(z_0)p(x_0 | z_0)p(y_0 | z_0) \prod_{t=1}^T p(x_t | z_t)p(y_t | z_t)p(z_t | z_{t-1}) \quad (7)$$

(Give intuition for this model)

#### 4.4 Learning and inference

We optimize the evidence of observed sentences  $\log p(\mathbf{x}) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$  by first marginalizing over latent states  $\mathbf{z}$  then performing gradient ascent. The emission sparsity constraints we introduce as well as state dropout allow us to both perform marginalization and compute the gradient of the evidence efficiently.

(Inference in POS model)

### 5 Experiments

We evaluate the HMM on two language modeling datasets.

#### 5.1 Language Modeling

**Datasets** The two datasets we used are the Penn Treebank (Marcus et al., 1993) and wikitext2 (Merity et al., 2016). Penn Treebank contains 929k tokens in the training corpus, with a vocabulary size of 10k. We use the preprocessing from Mikolov et al. (2011), which lowercases all words and substitutes words outside of the vocabulary with unks. Wikitext2 contains 2M tokens in the training corpus, with a vocabulary size of 33k. Casing is preserved, and all words outside the vocab are unked. Both datasets contain inter-sentence dependencies, due to the use of documents consisting of more than a single sentence.

**Implementation** We train two-layer LSTM recurrent neural networks with 256 or 512 units, as well as two-layer feed-forward neural networks with 256 or 512 units. The HMMs we train follow the sparsity constraints outlined in the previous section with a dropout rate of 0.5, and we vary the total number of states as well as states per word. We optimize all models with AdamW (Loshchilov and Hutter, 2017).

Table 1: Perplexities on the Penn Treebank dataset. The FF model is a 256-dim 2-layer feedforward neural network with a window of 4 previous tokens with 0.3 dropout. The LSTM is a 256-dim 2-layer recurrent neural network with 0.3 dropout. The HMM is a 64k-state HMM with 0.5 state dropout.

| Model | Num Params | Valid PPL | Test PPL |
|-------|------------|-----------|----------|
| FF    | 2.8M       | 164       | -        |
| LSTM  | 3.6M       | 97        | -        |
| HMM   | 19.8M      | 122       | -        |
| HMM   | 7.7M       | 125       | -        |

We experimented with a couple batching strategies: On Penn Treebank, the first strategy discarded the inter-sentence dependencies and shuffled all sentences, and the second treated the corpus as a single flat document without shuffling. On Wikitext2, we either shuffled at the document level or treated the corpus as a single document. Prior work on both corpora treated the corpora as single documents.

See Appendix A for the hyperparameters for all models.

#### 5.2 Results

We report perplexities for Penn Treebank in Table 1 and for wikitext2 in Table 2.

The HMMs in all cases outperform the feedforward models (FF), but underperform the recurrent LSTMs. Since the HMMs require parameters linear in the number of hidden states, we find that the performance of the HMMs scales poorly compared to the other models which only require parameters that scale linearly with the vocabulary size. Although representing and summing over the hidden states allows us to explicitly capture uncertainty in the hidden state, it proves to be a limiting factor in terms of performance.

In the remainder of this section we ablate and analyze the HMMs on Penn Treebank.

**Sparse emission constraint ablation** We ablate the emission sparsity constraints in Table 3, and find that the Brown emission constraints outperforms the uniform emission constraints in all model sizes.

One explanation for the relative benefit of Brown emission constraints over uniform is due to the effect of Brown clusters. The goal of Brown clusters



Table 2: Perplexities on the `wikitext2` dataset. The FF model is a 256-dim 2-layer feedforward neural network with a window of 4 previous tokens with 0.3 dropout. The LSTM is a 256-dim 2-layer recurrent neural network with 0.3 dropout. The HMM is a 32k-state HMM with 0.5 state dropout.

| Model | Num params | Valid PPL | Test PPL |
|-------|------------|-----------|----------|
| FF    |            | 209       | -        |
| LSTM  |            | 125       | -        |
| HMM   |            | 167       | -        |

Table 3: The perplexities for the different emission sparsity constraints in a 1024 state HMM as well as larger HMMs for which exact inference without sparsity is too expensive. The quantities  $|\mathcal{Z}|$  and  $k$  are the number of hidden states and the number of states per word respectively. The HMMs with 1024 states do not have any dropout, while the 8k and 16k state HMMs have unstructured dropout at a rate of 0.5.

| Constraint | $ \mathcal{Z} $ | $k$ | Valid PPL |
|------------|-----------------|-----|-----------|
| Uniform    | 1k              | 32  | -         |
| Brown      | 1k              | 32  | -         |
| None       | 1k              | 1k  | -         |
| Uniform    | 8k              | 128 | 150*      |
| Brown      | 8k              | 128 | 142*      |
| Uniform    | 16k             | 128 | 146*      |
| Brown      | 16k             | 128 | 134*      |

tering is to place two words in the same cluster if they are used in the same context. In Table 4, we observe that the entropy of the emission distribution for models with uniform emission constraints is lower than models with brown constraints, and the entropy of the transition distributions is higher. This implies that the burden of modeling ambiguity is pushed onto the transition distribution, since the uniform models are less likely to place words that appear in similar contexts together.

**Dropout analysis** The effect of the different dropout strategies and rates is shown in Table 5. We found that state dropout outperformed unstructured dropout overall. This in addition to the computational benefits afforded by state dropout led us to prefer this strategy.

Unstructured dropout was sensitive to batching strategies. We observed a larger gap between training and validation perplexities with the unshuffled batching strategy, whereas unstructured dropout

Table 4: The average entropies of the the emission and transition distributions for HMMs with uniform and Brown cluster emission constraints. All models have  $k = 128$  states per word and use unstructured dropout with a rate of  $p = 0.5$ .

| Constraint | $ \mathcal{Z} $ | H(emit) | H(transition) |
|------------|-----------------|---------|---------------|
| Uniform    | 8k              | -       | -             |
| Brown      | 8k              | -       | -             |
| Uniform    | 16k             | -       | -             |
| Brown      | 16k             | -       | -             |

Table 5: State occupancies for the dropout strategies and rates. All models were HMMs with Brown cluster emission constraints, 16k total states, and 128 states per word (and therefore 128 Brown clusters).

| Dropout type | $p$  | Valid PPL | Occupancy |
|--------------|------|-----------|-----------|
| Unstructured | 0.25 | -         | -         |
| Unstructured | 0.5  | -         | -         |
| Unstructured | 0.75 | -         | -         |
| State        | 0.25 | -         | -         |
| State        | 0.5  | -         | -         |
| State        | 0.75 | -         | -         |

appeared to be robust to the batching strategy.

Analysis 1: Discuss what happens with just state dropout, without partitioning. (Too much variance in gradient estimator, would require variance reduction. Get numbers or just handwaive? Consider exact inference (in constrained model) this work. )

**Number of Brown clusters** We next examine the sensitivity of HMMs to the number of Brown clusters in Table 6. We find that models at with 16k or 32k total states are not sensitive to the number of Brown clusters in the range where marginalization is computationally feasible. This contrasts with the observation that the number of Brown clusters influenced performance at 1k total states.

**Weight tying ablation** Analysis 1: Discuss performance using different tying methods (found this to not affect performance) and computational savings in terms of number of parameters. Discuss relative parameter inefficiency compared to LSTM / FF.

**State size ablation** We find that as we increase the number of total states while keeping the Brown

Table 6: The perplexities for HMMs with 16k states and different numbers of Brown clusters for constraining the emission distribution of the HMMs.  $|\mathcal{Z}|$  is the total number of hidden states. All models have 0.5 state dropout.

| $ \mathcal{Z} $ | Num clusters | Valid PPL |
|-----------------|--------------|-----------|
| 16k             | 32           | 136       |
| 16k             | 64           | 137       |
| 16k             | 128          | 133       |
| 16k             | 256          | 137       |

Table 7: The perplexities for HMMs with 128 Brown clusters for constraining the emission distribution of the HMMs.  $|\mathcal{Z}|$  is the total number of hidden states. All models have 0.5 state dropout.

| $ \mathcal{Z} $ | Num clusters | Valid PPL |
|-----------------|--------------|-----------|
| 16k             | 128          | 133       |
| 32k             | 256          | 126       |
| 65k             | 512          | 124       |

clusters fixed, performance improves up until we have 65k states.

Analysis 1: How does state usage change as clusters remain constant but the number of states (and states per word) increases?

**Parameterization ablation** Analysis 1: Is there a performance difference between neural / scalar parameterization, and is it consistent across state sizes?

Discussion 1: Memory savings when working with state dropout to not instantiate the full matrices during training, which makes working with a neural parameterization beneficial.

### 5.3 Error analysis

## 6 Discussion

## Acknowledgments

TBD

## References

- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. 2003. A neural probabilistic language model. *J. Mach. Learn. Res.*, 3(null):1137–1155.
- Peter F. Brown, Peter V. deSouza, Robert L. Mercer, Vincent J. Della Pietra, and Jenifer C. Lai. 1992. Class-based n-gram models of natural language. *Comput. Linguist.*, 18(4):467–479.
- Jan Buys, Yonatan Bisk, and Yejin Choi. 2018. Bridging hmms and rnns through architectural transformations.
- Antoine Dedieu, Nishad Gothoskar, Scott Swingle, Wolfgang Lechach, Miguel Lázaro-Gredilla, and Dileep George. 2019. [Learning higher-order sequential structure with cloned hmms](#).
- Zoubin Ghahramani and Michael I. Jordan. 1997. [Factorial hidden markov models](#). *Mach. Learn.*, 29(2–3):245–273.
- Yoon Kim, Chris Dyer, and Alexander M. Rush. 2019. [Compound probabilistic context-free grammars for grammar induction](#). *CoRR*, abs/1906.10225.
- Viktoriia Krakovna and Finale Doshi-Velez. 2016. [Increasing the interpretability of recurrent neural networks using hidden markov models](#).
- Ilya Loshchilov and Frank Hutter. 2017. [Fixing weight decay regularization in adam](#). *CoRR*, abs/1711.05101.
- Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. [Building a large annotated corpus of English: The Penn Treebank](#). *Computational Linguistics*, 19(2):313–330.
- Stephen Merity, Nitish Shirish Keskar, and Richard Socher. 2017. [Regularizing and optimizing LSTM language models](#). *CoRR*, abs/1708.02182.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. [Pointer sentinel mixture models](#). *CoRR*, abs/1609.07843.
- Tomas Mikolov, Anoop Deoras, Stefan Kombrink, Lukás Burget, and Jan Cernocký. 2011. [Empirical evaluation and combination of advanced language modeling techniques](#). pages 605–608.
- Tomas Mikolov, Martin Karafiát, Lukás Burget, Jan Cernocký, and Sanjeev Khudanpur. 2010. [Recurrent neural network based language model](#). pages 1045–1048.
- Anjan Nepal and Alexander Yates. 2013. Factorial hidden markov models for learning representations of natural language. *CoRR*, abs/1312.6168.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.

Ke M. Tran, Yonatan Bisk, Ashish Vaswani, Daniel  
 Marcu, and Kevin Knight. 2016. [Unsupervised neu-  
 ral hidden markov models](#). *CoRR*, abs/1609.09007.

Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals.  
 2014. [Recurrent neural network regularization](#).  
*CoRR*, abs/1409.2329.

## A Hyperparameters

### LSTM

- 2 layers

## B Supplemental Material