
Doubly Sparse: Sparse Mixture of Sparse Experts for Efficient Softmax Inference

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Computations for the softmax function in neural network models are expensive
2 when the number of output classes is large. This can become a significant issue
3 in both training and inference for such models. In this paper, we present Doubly
4 Sparse Softmax (DS-Softmax), *Sparse Mixture of Sparse of Sparse Experts*, to
5 improve the efficiency for softmax inference. During training, our method learns
6 a two-level class hierarchy by dividing entire output class space into several par-
7 tially overlapping experts. Each expert is responsible for a learned subset of the
8 output class space and each output class only belongs to a small number of those
9 experts. During inference, our method quickly locates the most probable expert
10 to compute small-scale softmax. Our method is learning-based and requires no
11 knowledge of the output class partition space a priori. We empirically evaluate our
12 method on several real-world tasks and demonstrate that we can achieve significant
13 computation reductions without loss of performance.

14 1 Introduction

15 Deep learning models have demonstrated impressive performance in many classification problems (Le-
16 Cun et al., 2015). In many of these models, softmax function/layer is commonly used to produce
17 categorical distributions over the output space. Due to its linear complexity, computation for softmax
18 layer can become a bottleneck with large output dimensions, such as language modelling (Bengio
19 et al., 2003), neural machine translation (Bahdanau et al., 2014) and face recognition (Sun et al.,
20 2014). In some models, softmax contributes to more than 95% computation. This becomes more of
21 an issue when computational resource is limited, like mobile devices (Howard et al., 2017). Many
22 methods have been proposed to reduce softmax complexity for both training and inference phases.
23 In terms of inference, our goal is not to computing the exact categorical distribution over the whole
24 vocabulary, but rather to search for top-K classes accurately and efficiently.

25 Our work aims to improve the inference efficiency of the softmax layer. We propose a novel Doubly
26 Sparse softmax (DS-Softmax) layer. The proposed method is motivated by (Shazeer et al., 2017),
27 and it learns a two-level overlapping hierarchy using *sparse* mixture of *sparse* experts. Each expert
28 is trained to only contain a small subset of entire output class space, while each class is permitted
29 to belong to more than one expert. Given a set of experts and an input vector, the DS-Softmax first
30 selects the top expert that is most related to the input (in contrast to a dense mixture of experts), and
31 then the chosen expert could return a scored list of most probable classes in its sparse subset. This
32 method can reduce the linear complexity in original softmax significantly since it does not need to
33 consider the whole vocabulary.

34 We conduct experiments in different real tasks, ranging from language modeling to neural machine
35 translation. We demonstrate our method can reduce softmax computation dramatically without loss of
36 prediction performance. For example, we achieved more than 23x speedup in language modelling and

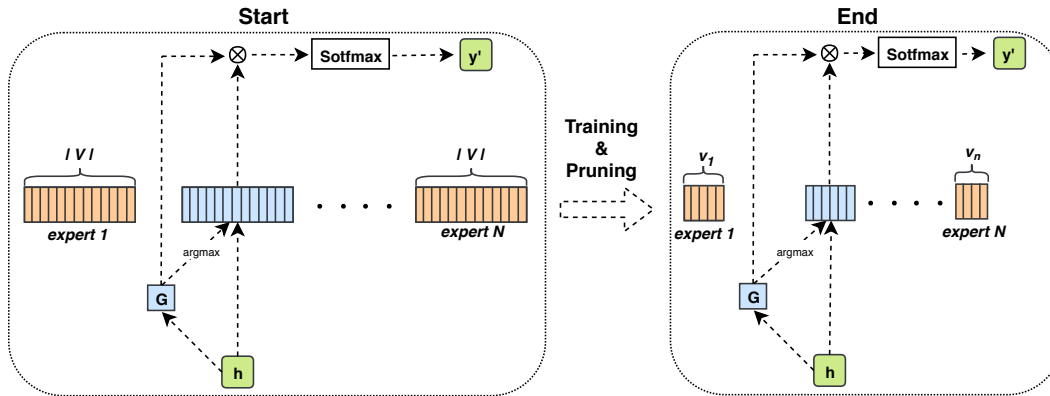


Figure 1: Overview of DS-Softmax. Initial model is similar to sparsely gating mixture of experts model. After pruning, each expert will only consists partial outputs v_n instead of $|V|$.

37 15x speedup in translation with similar performances. In real device, our method also demonstrates
 38 similar speedup with theoretic one.

39 2 DS-Softmax: Sparse Mixture of Sparse Experts

40 Goodman (2001) studied a two-level hierarchy for language modeling, where each word belongs to
 41 one unique cluster. (A “cluster” here refers to a cluster of words.) From this perspective, our method
 42 can be as an extension of their method to allow overlapping hierarchy. This is because, in language
 43 modeling, it is often difficult to exactly assign a word to a single cluster. For example, if we want to
 44 predict next word of “I want to eat ___” and one possible correct answer is “cookie”, we can quickly
 45 notice that possible answer belongs to something eatable. So if we only search right answer inside
 46 words with the eatable property, we can dramatically increase the efficiency. On the other, though
 47 “cookie” is one of the correct answers, it might also like appear under some non-eatable context,
 48 such as “a piece of data” in computer science. Thus, a two-level overlapping hierarchy can naturally
 49 accommodate word homonyms like this by allowing each word to belong to more than one cluster.
 50 We believe this observation is likely to be true in other applications besides language modeling.

51 **Overview** Doubly Sparse softmax (DS-Softmax) is designed to capture such overlapped two-level
 52 hierarchy among output classes. In DS-Softmax, the first level is the *sparse* mixture and second level
 53 contains several *sparse* experts. (Here an expert can be thought as a similar concept as cluster.) The
 54 sparse mixture is to choose the right expert/cluster while sparse experts are responsible to separate
 55 full output space into multiple, overlapped and small class clusters. The design of mixture gating is
 56 inspired by Shazeer et al. (2017) but each expert in their model needs to search whole output space,
 57 while DS-Softmax only searches a small subset. This becomes much faster given large output space.

58 The first level of sparsification is a sparse gating mechanism inspired by Sparsely-Gated Mixture of
 59 Experts (Shazeer et al., 2017), where only partial experts are activated. For faster inference purpose,
 60 only top-one expert is chosen, which is corresponding to choose the right experts. The second level
 61 sparsification is the sparse experts, which output a categorical distribution for only a subset output
 62 classes. To sparsify each expert, we apply group lasso loss to restrain the weights inside softmax.
 63 Furthermore, the utilization of each experts is balanced with additional losses. The detail of our
 64 method can be found in Appendix.

65 3 Experiments

66 We evaluate the proposed method on both synthetic and real tasks. For the synthetic task, our goal is
 67 to demonstrate that our learning method could discover the hidden two-level hierarchy automatically.
 68 We also evaluate both theoretical speedup (FLOPs) and real device speedup (latency on CPU) on three
 69 different real tasks: natural language modelling, neural machine translation and Chinese handwritten

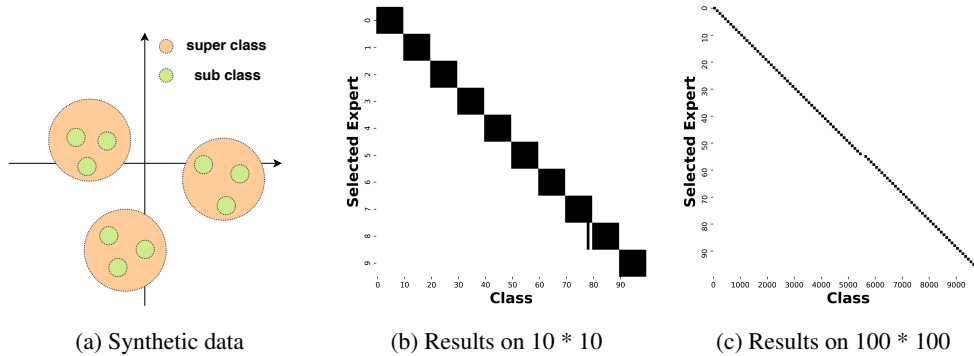


Figure 2: (a) Illustration of data generation. (b) and (c) Results on discovered sparse experts on 10x10 and 100x100 datasets. The x-axis indicates class and y-axis shows the selected expert for handling this class. The order of x-axis is arranged through their super class information. For example, each 10 sub classes are belonged to one super classes in (b).

70 character recognition. In those real tasks, all layers except the DS-Softmax layer are pre-trained in all
 71 tasks.

72 3.1 Synthetic task

73 One two-level hierarchy synthetic dataset is illustrated Fig 2a. Each super class contain multiple
 74 sub classes. Two different sizes are evaluated, 10x10 (super classes x sub classes) and 100x100.
 75 The result is illustrated in Fig. 2b and Fig. 2c. We found our DS-Softmax can perfectly capture the
 76 hierarchy. For sanity check and visualization purposes, the ground-truth two hierachy in the synthetic
 77 data does not have overlappings.

78 3.2 Comparisons on FLOPs and Real Device

79 Real device experiments were conducted on machine with Two Intel(R) Xeon(R) CPU @ 2.20GHz
 80 and 16G memory. All tested models are re-implemented using numpy. Two configurations of
 81 SVD-Softmax Shim et al. (2017) are evaluated, SVD-5 and SVD-10. They use top 5% and 10%
 82 dimension for final evaluation in their preview window and window width is 16. Latency of each
 83 sample is shown in table 1. According to the result, our DS-Softmax can achieve not only better
 84 FLOPs speedup but also much better performance on latency.

Task	Full		DS-64 (Ours)			SVD-5			SVD-10		
	Value	ms	Value	FLOPs	ms	Value	FLOPs	ms	Value	FLOPs	ms
PTB	0.252	0.73	0.258	15.99x	0.05	0.249	6.67x	0.12	0.251	5.00x	0.18
Wiki-2	0.257	3.07	0.259	23.86x	0.12	0.253	7.35x	0.43	0.255	5.38x	0.63
En-Ve	25.2	1.91	25.0	15.08x	0.12	25.0	6.77x	0.39	25.1	5.06x	0.42
CASIA	0.906	1.61	0.901	6.91x	0.25	0.899	3.00x	0.59	0.902	2.61x	0.68

Table 1: Comparison with SVD-softmax on real device latency. The “ms” indicates the latency in microseconds. Bold fonts indicate better results.

85 4 Conclusion

86 In this paper, we present *doubly sparse: sparse mixture of sparse experts* for efficient softmax
 87 inference. Our method is trained end-to-end. It learns a two-level overlapping class hierarchy. Each
 88 expert is learned to be only responsible for a small subset of the output class space. During inference,
 89 our method first identifies the responsible expert and then perform a small scale softmax computation
 90 just for that expert. Our experiments on several real-world tasks have demonstrated the efficacy of
 91 our proposed method.

92 **References**

- 93 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly
94 learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- 95 Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. A neural probabilistic
96 language model. *Journal of machine learning research*, 3(Feb):1137–1155, 2003.
- 97 Joshua Goodman. Classes for fast maximum entropy training. In *Acoustics, Speech, and Signal
98 Processing, 2001. Proceedings.(ICASSP'01). 2001 IEEE International Conference on*, volume 1,
99 pp. 561–564. IEEE, 2001.
- 100 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv
101 preprint arXiv:1503.02531*, 2015.
- 102 Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand,
103 Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for
104 mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- 105 Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436, 2015.
- 106 Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and
107 Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv
108 preprint arXiv:1701.06538*, 2017.
- 109 Kyuhong Shim, Minjae Lee, Iksoo Choi, Yoonho Boo, and Wonyong Sung. Svd-softmax: Fast
110 softmax approximation on large vocabulary neural networks. In *Advances in Neural Information
111 Processing Systems*, pp. 5463–5473, 2017.
- 112 Yi Sun, Xiaogang Wang, and Xiaoou Tang. Deep learning face representation from predicting 10,000
113 classes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
114 1891–1898, 2014.

115 **A Detail of Methods**

116 **Sparse gating.** The first level of sparsification is a sparse gating mechanism inspired by Shazeer
 117 et al. (2017), which is to design to choose the right experts. The sparse gating outputs a sparse
 118 activation over a set of experts. For faster inference purpose, only the top-one expert is chosen
 119 here. One major difference comparing to Shazeer et al. (2017) is described as follows. Suppose
 120 we have K experts. Given input activation vector $h \in \mathbb{R}^d$, gating values $G_k(h)$, $k = 1, \dots, K$, are
 121 normalized prior to the selection as shown in Eq. 1 and then we choose the gate with the largest value
 122 $g_k = \max_i G_i(h)$ and set all other gates to be zero. Also, corresponding k -th expert is chosen.

$$G_k(h) = \frac{\exp(W_k^g h)}{\sum_{k'} \exp(W_{k'}^g h)}, \quad (1)$$

$$g_k = \begin{cases} G_k(h), & \text{if } k = \arg \max_i G_i(h), \\ 0, & \text{otherwise.} \end{cases}$$

123 This allows gradient to be back-propagated to whole W^g instead of W_k^g only, $W^g \in \mathbb{R}^{K \times d}$. In
 124 Shazeer et al. (2017), normalization is done after top-K experts are selected. We can not do that since
 125 we only choose top-1 expert since it will carry no gradient information since it becomes constant 1.
 126 Given the sparse gate, we compute the probability of class c as,

$$O(h) = p(c|h) = \frac{\exp(\sum_k g_k W_{(c,k)}^e h)}{\sum_{c'} \exp(\sum_k g_k W_{(c',k)}^e h)}, \quad (2)$$

127 where $W_{(c,k)}^e \in \mathbb{R}^d$ is softmax embedding weight vector for class c in expert k . Note that only one g_k
 128 (the chosen expert) is nonzero in the formulation above. The gating values can be interpreted as an
 129 inverse temperature term for final categorical distribution produced by the chosen expert k Hinton et al.
 130 (2015), shown in Eq. 2. A smaller g_k gives a more uniform distribution and larger g_k corresponds to
 131 a sharper one.

132 **Sparse experts with group lasso.** The second level sparsification is the sparse experts, which
 133 output a categorical distribution for only a subset output classes. To sparsify each expert, we apply
 134 group lasso loss to restrain the $W_{(c,k)}^e$, shown in Eq. 3. Then, pruning is carried out for $W_{(c,k)}^e$ during
 135 training with γ is a lasso threshold according to Eq. 4.

$$\mathcal{L}_{lasso} = \sum_k \sum_c \|W_{(c,k)}^e\|_2, \quad (3)$$

$$W_{(c,k)}^e = \begin{cases} W_{(c,k)}^e, & \text{if } \|W_{(c,k)}^e\|_2 > \gamma, \\ 0, & \text{otherwise.} \end{cases} \quad (4)$$

136 **Loading Balance.** We denote the sparsity percentage out of full softmax in k -th expert as sparsity $_k$
 137 and proportion of k -th expert activated as utilization $_k$. Then, the overall speedup compared to the full
 138 softmax can be calculated as $1 / \sum_k (\text{utilization}_k * \text{sparsity}_k)$. Thus, better utilization is essential for
 139 speedup as well. For example, there is no speedup if the expert with full output space is always chosen.
 140 We borrow a similar loading balance function from Shazeer et al. (2017) in Eq. 5. It encourages the
 141 utilization percentage of each expert to be balanced by maximizing the coefficient of variation (CV)
 142 for gating outputs. In addition, to encourage the exclusiveness of classes, we incorporate group lasso
 143 loss on expert level where each class should only exist in only one expert as shown in Eq. 6.

$$\mathcal{L}_{load} = -\text{CV} \left(\sum_{h \in H(x)} G(h) \right), \quad (5)$$

$$\mathcal{L}_{expert} = \sum_k \sqrt{\sum_c \|W_{(c,k)}^e\|_2^2}. \quad (6)$$

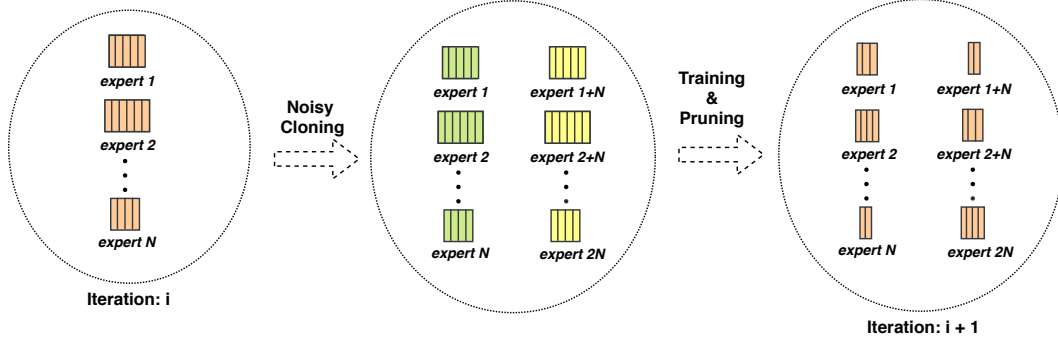


Figure 3: The mitosis training scheme: the sparsity is inherited when parent experts produce offspring, reducing the memory requirements for training with more experts.

Algorithm 1

- 1: **Initialization:** Let x be the input, y be the corresponding label, H be the pretrained function, V be the output dimension and $D(y', y)$ be an arbitrarily distance function. Set $W^e \leftarrow$ parameters for experts and $W^g \leftarrow$ parameters for the gating network. The hyper-parameter t denotes target performance.
 - 2: **while** $epoch < \text{Max}$ **do**
 - 3: $epoch = epoch + 1$
 - 4: $\mathcal{L}_{task} = D(O(H(x)), y)$
 - 5: $\mathcal{L}_{all} = \mathcal{L}_{task} + \lambda_{load} \mathcal{L}_{load} + \lambda_{lasso} \mathcal{L}_{lasso} + \lambda_{expert} \mathcal{L}_{expert}$
 - 6: $W^e = W^e - \alpha \frac{\partial}{\partial W^e} \mathcal{L}_{all}(x, y; W^e, W^g)$
 - 7: $W^g = W^g - \alpha \frac{\partial}{\partial W^g} \mathcal{L}_{all}(x, y; W^e, W^g)$
 - 8: **if** $\mathcal{L}_{task} < t$ **then**
 - 9: **for all** $W_{(c,k)}^e \in W^e$ **do**
 - 10: $W_{(c,k)}^e = 0$, if $\|W_{(c,k)}^e\|_2 < \gamma$
-

144 **Mitosis training.** Memory might become a bottleneck during training if we initialize all experts
145 with full softmax. Therefore, we design one training scheme, called mitosis training, to reduce
146 memory requirement. The method is to initialize with a smaller model (fewer number of experts)
147 and then gradually breed to a bigger one after noisy cloning shown in Fig. 3. For each cloning, the
148 sparsity is inherited so that less memory is required. For example, in one of our experiments, we only
149 need 3.25x memory with 64 experts compared to a full softmax implementation.

150 **The final training algorithm.** Our final training objective, \mathcal{L}_{all} , consists of a combination of the
151 related contributions discussed above. We describe our training procedure in Algorithm 1.