

1 We thank the reviewers for their positive comments and constructive suggestions. We first address the key common
2 concerns and then move to individual reviewers. The paper will be updated accordingly in the camera-ready version.

3 **Comparison with prior works (R1, R2 and R4):** Besides the novel theoretical contribution, our work distinguished
4 itself from prior works on Bayesian sparse neural network by imposing a spike-and-slab prior with the Dirac spike
5 function. Hence automatically, all posterior samples are from exact sparse DNN models. In contrast, Blundell et al.
6 (ICML 2015) and Deng et al. (NIPS 2019) considered spike-and-slab priors with a Gaussian and Laplacian spike
7 respectively, while Ghosh et al. (ICML 2018) chose to use the popular horseshoe shrinkage prior. These existing works
8 actually yield posteriors over the dense DNN model space despite applying sparsity induced priors. In order to derive
9 explicit sparse inference results, user has to additionally determine certain pruning rules on the posterior, which is
10 similar to variational dropout (Molchanov et al. ICML 2017). Thus in our mind, they still belong to the category of
11 Bayesian pruning methods. However, we do agree with the reviewers that additional experiments could be made for
12 better comparison. For example, we conduct Horseshoe BNN under Simulation setting I: it achieves test RMSEs of
13 1.02 ± 0.01 (A) and 1.24 ± 0.03 (B), with neuron-wise sparsity ratio of 0.79 ± 0.19 (A) and 0.81 ± 0.09 (B). This
14 performance is worse than our method. Note that more experiments will be added in the final version.

15 **Extend current result to more complicated network models (R1, R2 and R3):** Our developed theory should work
16 as long as it is a regular DNN (i.e., networks only involve fully-connected layers). Unfortunately we are unable to
17 perform super large-scale experiments due to limited computing resources at hand. Extending current results to more
18 complicated networks (convolutional layer, residual network, etc.) is not trivial. Conceptually, it requires design of
19 structured sparsity (e.g., group sparsity in Neklyudov et al. NIPS 2017) to serve the purpose of faster prediction.
20 Theoretically, it requires deeper understanding of the expressive ability (i.e. approximation error) and capacity (i.e.,
21 packing or covering number) of the network model space. By intuition, we conjecture that the generalized theory should
22 still hold, but it will be a future work to provide rigorous theoretical support. To illustrate the practical value of our
23 method for complex tasks, as a preliminary experiment, we apply a 2-Conv-2-FC network on Fashion-MNIST. The
24 testing accuracy is 90.07% with 60% sparsity (connection-wise), where the baseline by dense model is 90.65%.

25 **R2:** Please refer to common concerns for major questions, and below is our response to minor questions:

26 - To induce structured group sparsity, say in CNN modeling, we can let all weights in the same filter share the same
27 binary indicator, such that the whole filter is turned on/off simultaneously. The theoretical correctness of this approach
28 is left as our future work.

29 - The reviewer is correct. Molchanov et al. (ICML 2017) is a Bayesian pruning method, and the word "frequentist" on
30 Line 112/113 is a typo that should be removed.

31 - The five datasets are chosen for numerical study since they have fairly larger sample size.

32 - Throughout the paper, we always report the standard deviation. The "standard error" mentioned in the supplementary
33 material is a typo.

34 **R3:** Our results can be extended to different variation distribution families as long as they contain a Dirac component
35 at zero (although the technical details might differ), and the choice of normal slab distribution in this paper is merely for
36 technical simplicity. Without such a Dirac component, the variational posterior can no longer automatically induce
37 sparse inferences. Therefore, those cases are not under our consideration.

38 The authors don't possess cutting-edge infrastructure, thus reporting the CPU/GPU time of our implementation is
39 not very meaningful. Relatively, under same computational environment, our method takes roughly twice more time
40 than the frequentist counterparts, which is majorly due to the reason that there are more parameters to optimize in our
41 algorithm. Developing more computational efficient algorithm could be a future direction.

42 **R4:** The proposed choice of λ does encourage posterior sparsity, but it doesn't play a dominating role on how sparse
43 the posterior would be. Instead, given the value of λ , it is the likelihood of the data that helps to adaptively choose the
44 data-dependent optimal posterior sparsity level. λ plays as a regularization, and we provide theoretical suggestion such
45 that it doesn't mask the importance of the data likelihood. In contrast, pruning methods such as AGP and LOT require
46 user-input knowledge to explicitly determine the sparsity level of the result. In our implementation of AGP and LOT,
47 we tried grid search for the sparsity levels ranging from 95% to 5%, and the level that yields the best testing accuracy is
48 chosen and reported.

49 For minor mistakes,

50 - Thanks for pointing out the typo and the outdated references. In Line 18, the sum should be iterated over $\gamma \in \Gamma^T$.

51 - Condition 4.4 is not required by Lemma 4.1 since we prove for cases either $\lim n(r_n^* + \xi_n^*) = \infty$ or $\lim n(r_n^* + \xi_n^*) \neq \infty$.

52 - For the question raised, as has been described in Line 214/215, the posterior sparsity is measured by $\sum_{i=1}^T \phi_i / T$ and
53 an individual connection (of the VB point estimator) is considered inactive if $\phi_i < 0.5$.