
On permutation symmetries in Bayesian neural network posteriors: Appendix

Anonymous Author(s)

Affiliation

Address

email

1 A Additional details on the alignment method

2 In the main paper, to align the distributions with respect to permutation matrices we argue to use the
3 Wasserstein distance rather than the Kullback-Leibler (KL) divergence. Indeed, by considering the KL
4 divergence $\text{KL}[P_{\#}q_1 \parallel q_0]$ between Gaussians we have

$$\text{KL}[P_{\#}q_1 \parallel q_0] = \log \det \text{diag}(s_0) - \log \det \text{diag}(Ps_1) + \text{Tr}(\text{diag}(Ps_1s_0^{-1})) + \quad (1)$$

$$(\mathbf{m}_0 - P\mathbf{m}_1)^\top \text{diag}(s_0^{-1})(\mathbf{m}_0 - P\mathbf{m}_1) \quad (2)$$

5 It's easy to verify that the first three terms do not depend on P , leading to just a distance between
6 means and disregarding any covariance information. In the figure below, we visualize the difference
7 between doing linear assignment problem (LAP) with the KL cost and LAP with the Wasserstein cost.

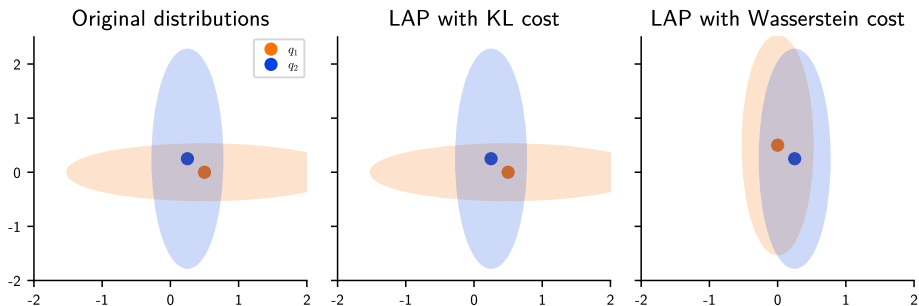


Figure A1: **Alignment using different objectives.** Given two distributions symmetrical w.r.t. the $y = x$ plane, using the KL cost LAP results in the identity permutation (which fails to recover the symmetry), while the Wasserstein cost better aligns the two distributions

8 B Experimental setup

9 Tables 2 and 3 show details on the multilayer perceptrons (MLPs) and convolutional neural networks
10 (CNNs) base architectures used in our experimental campaign, while Table 1 reports the hyper-
11 parameters used in the experiments. Note that differently from Entezari et al. [5] and Ainsworth et al.
12 [2], we don't use data augmentation. A possible protocol for handling data augmentation in Bayesian
13 neural networks (BNNs) is presented by Osawa et al. [14] and involves carefully tuning the likelihood
14 temperature to correctly counting the number of data points.

Table 1: Hyperparameters used for the experiments

Dataset	CIFAR10	MNIST
Model	ResNet20	MLP
Data Aug.	False	False
Batch size	500	500
Temperature	1.0	1.0
Test samples	128	128
Train samples	1	1
VI std. init	0.01	0.01
Base features	16	512
Prior var	0.01	0.0025
Learning rate	0.000001	0.000001
Train epochs	1000	1000

Table 2: MLP

Layer	Dimensions
Linear-ReLU	$512 \times D_{in}$
Linear-ReLU	512×512
Linear-ReLU	512×512
Linear-Softmax	$D_{out} \times 512$

Table 3: ResNet20

Layer	Dimensions
Conv2D	$16 \times 3 \times 3 \times D_{in}$
Residual Block	$\begin{matrix} 3 \times 3, 16 \\ 3 \times 3, 16 \end{matrix} \times 3$
Residual Block	$\begin{matrix} 3 \times 3, 32 \\ 3 \times 3, 32 \end{matrix} \times 3$
Residual Block	$\begin{matrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{matrix} \times 3$
AvgPool	8×8
Linear-Softmax	$D_{out} \times 64$

15 B.1 Computing platform

16 The experiments have been performed using JAX [4] and run on two AWS p4d.24xlarge instances
 17 with 8 NVIDIA A100 GPUs. Experiments were conducted using in the eu-west-1 region, which has a
 18 carbon efficiency of 0.62 kgCO₂eq/kWh. A cumulative of 6500 hours of computation was performed
 19 on GPUs and it includes interactive sessions as well as small experiments with very low GPU usage,
 20 providing a pessimistic estimation of the true utilization. Total emissions are estimated to be 1007.5
 21 kgCO₂eq of which 100 percents were directly offset by AWS.

22 C Additional results

23 We present timings obtained by profiling the time needed to solve the sum of bilinear assignment
 24 problems (SOLAP) with the Wasserstein cost, as well as the time for the deterministic case [2]. In
 25 Fig. A2 we show the results for MLP and ResNet20 architectures, varying the model width. It
 26 is evident that, in the majority of cases, the algorithm completes within a minute. Moreover, as
 27 anticipated, in case of variational inference (VI) solving our distribution alignment problem for wide
 28 neural networks is more computationally demanding compared to merely matching weights from
 29 stochastic gradient descent (SGD) solutions.

30 Finally, we also test our setup on the CIFAR100 dataset [11]. Surprisingly, we were not able to
 31 replicate the same level of performance as in the other cases. In Fig. A3, we see that, despite
 32 converging well, we fall short to find zero-barrier solutions. Similarly to the comments of Ainsworth
 33 et al. [2], we also stress that the failure to align distributions does not rule out the existence of a proper
 34 permutation map that the algorithm couldn't find. Nonetheless, this raises a number of questions:
 35 the Bayesian posterior is the product of two ingredients, the prior and the likelihood, conditioned to
 36 observing a dataset.

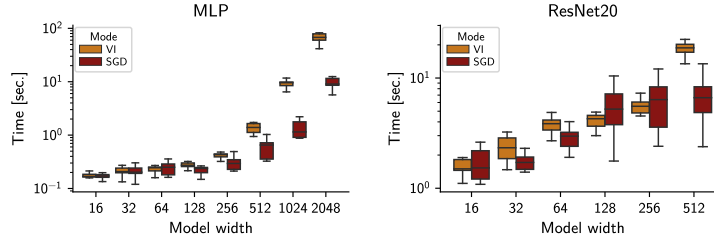


Figure A2: **Timings.** Profile of the algorithms to align the VI solutions and to match weights from SGD solutions.

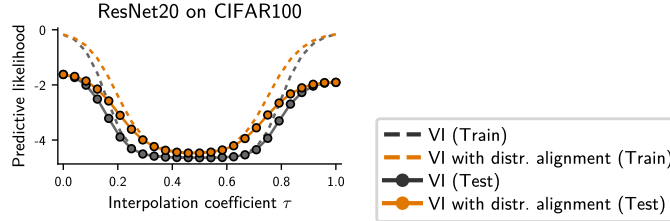


Figure A3: **Alignment failure.** The method proposed fails to recover zero-barrier solutions for CIFAR100.

37 D A primer on variational inference for Bayesian neural networks

38 VI is a classic tool to tackle intractable Bayesian inference [9, 3]. VI casts the inference problem into
 39 an optimization-based procedure to compute a tractable approximation of the true posterior. Assume
 40 a generic parametric model f parameterized by some unknown parameters θ (i.e. $f(\cdot, \theta)$) and a
 41 collection of data $\mathbf{y} \in \mathbb{R}^N$ corresponding to some input points $\mathbf{X} = \{\mathbf{x}_i \mid \mathbf{x}_i \in \mathbb{R}^{D_{\text{in}}}\}_{i=1, \dots, N}$. In
 42 our setting, we have a probabilistic model $p(\mathbf{y} \mid f(\mathbf{X}; \theta))$ with parameters θ , a prior distributions on
 43 them $p(\theta)$ and a set of observations $\{\mathbf{X}, \mathbf{y}\}$. In a nutshell, the general recipe of VI consists of (i)
 44 introducing a set \mathcal{Q} of distributions; (ii) defining a tractable objective that “measure” the distance
 45 between any arbitrary distribution $q(\theta) \in \mathcal{Q}$ and the true posterior $p(\theta \mid \mathbf{y})$; and finally (iii) providing
 46 a programmatic way to find the distribution $\tilde{q}(\theta)$ that minimizes such distance. In practice, $q(\theta)$
 47 has some free parameters ν (also known as *variational parameters*), which are optimized such that
 48 the approximating distribution $q(\theta; \nu)$ is as closer as possible to the true posterior $p(\theta \mid \mathbf{y})$. We can
 49 derive the variational objective starting from the definition of the KL,

$$\begin{aligned} \text{KL}[q(\theta; \nu) \parallel p(\theta \mid \mathbf{y})] &= \mathbb{E}_{q(\theta; \nu)} [\log q(\theta; \nu) - \log p(\theta \mid \mathbf{y})] = & (3) \\ &= \mathbb{E}_{q(\theta; \nu)} [\log q(\theta; \nu) - \log p(\mathbf{y} \mid \theta) - \log p(\theta)] + \log p(\mathbf{y}) \end{aligned}$$

50 Rearranging we have that

$$\log p(\mathbf{y}) - \text{KL}[q(\theta; \nu) \parallel p(\theta \mid \mathbf{y})] = \mathbb{E}_{q(\theta; \nu)} [\log q(\theta; \nu) - \log p(\mathbf{y} \mid \theta) - \log p(\theta)] \quad (4)$$

51 The r.h.s. of the equation defines our variational objective, also known as evidence lower bound
 52 (ELBO), that can be arranged as follows,

$$\mathcal{L}_{\text{ELBO}}(\nu) = \underbrace{\mathbb{E}_{q(\theta; \nu)} \log p(\mathbf{y} \mid \theta)}_{\text{Model fitting term}} - \underbrace{\text{KL}[q(\theta; \nu) \parallel p(\theta)]}_{\text{Regularization term}}. \quad (5)$$

53 This formulation highlights the property of this objective, which is made of two components: the first
 54 one is the expected log-likelihood under the approximate posterior q and measures how the model fits
 55 the data. The second term, on the other hand, has the regularization effect of penalizing posteriors
 56 that are far from the prior as measured by the KL. Before diving into the challenges of optimization
 57 of the ELBO, we shall spend a brief moment discussing the form of the approximating distribution q .
 58 One of the simplest and easier choice is the mean field approximation [7], where each variable θ_i
 59 is taken to be independent with respect to the remaining θ_{-i} . Effectively, this imposes a factorization
 60 of the posterior,

$$q(\theta; \nu) = \prod_{i=1}^K q(\theta_i; \nu_i) \quad (6)$$

61 where ν_i is the set of variational parameters for the parameter θ_i . On top of this approximation, $q(\theta_i)$
 62 is often chosen to be Gaussian,

$$q(\theta_i) = \mathcal{N}(\mu_i, \sigma_i^2) \quad (7)$$

63 Now, the collection of all means and variances $\{\mu_i, \sigma_i^2\}_{i=1}^K$ defines the set of variational parameters
 64 to optimize.

65 For BNNs the analytic evaluation of the ELBO (and its gradients) is always untractable due the non-
 66 linear nature of the expectation of the log-likelihood under the variational distribution. Nonetheless,
 67 this can be easily estimated via Monte Carlo integration [13], by sampling N_{MC} times from q_ν ,

$$\mathbb{E}_{q(\theta; \nu)} \log p(\mathbf{y} | \theta) \approx \frac{1}{N_{\text{MC}}} \sum_{j=1}^{N_{\text{MC}}} \log p(\mathbf{y} | \tilde{\theta}_j), \quad \text{with } \tilde{\theta}_j \sim q(\theta; \nu) \quad (8)$$

68 In practice, this is as simple as re-sampling the weights and the biases for all the layers N_{MC} times
 69 and computing the output for each new sample.

70 We now have a tractable objective that needs to be optimized with respect to the variational parameters
 71 ν . Very often the KL term is known, making its differentiation trivial. On the other hand the expectation
 72 of the likelihood is not available, making the computation of its gradients more challenging. This
 73 problem can be solved using the so-called *reparameterization trick* [19, 10]. The reparameterization
 74 trick aims at constructing θ as an invertible function \mathcal{T} of the variational parameters ν and of another
 75 random variable ε , so that $\theta = \mathcal{T}(\varepsilon; \nu)$. Generally, a \mathcal{T} that suits this constraint might not exists;
 76 Ruiz et al. [18] discuss how to build “weakly” dependent transformation \mathcal{T} for distributions like
 77 Gamma, Beta and Log-normal. For discrete distributions, instead, one could use a continuous
 78 relaxation, like the Concrete [12]. ε is chosen such that its marginal $p(\varepsilon)$ does not depend on the
 79 variational parameters. With this parameterization, \mathcal{T} separates the deterministic components of q
 80 from the stochastic ones, making the computation of its gradient straightforward. For a Gaussian
 81 distribution with mean μ and variance σ^2 , \mathcal{T} corresponds to as simple scale-location transformation
 82 of an isotropic Gaussian noise,

$$\theta \sim \mathcal{N}(\mu, \sigma^2) \iff \theta = \mu + \sigma\varepsilon \quad \text{with } \varepsilon \sim \mathcal{N}(0, 1). \quad (9)$$

83 This simple transformation ensures that $p(\varepsilon) = \mathcal{N}(0, 1)$ does not depends on the variational param-
 84 eters $\nu = \{\mu, \sigma^2\}$. The gradients of the ELBO can be therefore computed as

$$\nabla_\nu \mathcal{L}_{\text{ELBO}} = \mathbb{E}_{p(\varepsilon)} [\nabla_\theta \log p(\mathbf{y} | \theta) |_{\theta=\mathcal{T}(\varepsilon; \nu)} \nabla_\nu \mathcal{T}(\varepsilon; \nu)] - \nabla_\nu \text{KL} [q(\theta; \nu) || p(\theta)] . \quad (10)$$

85 The gradient $\nabla_\theta \log p(\mathbf{y} | \theta)$ depends on the model and it can be derived with automatic differentiation
 86 tools [1, 15], while $\nabla_\nu \mathcal{T}(\varepsilon; \nu)$ doesn’t have any stochastic components and therefore can be known
 87 deterministically. Note that the reparameterization trick can be also used when the KL is not
 88 analytically available. In that case, we would end up with,

$$\nabla_\nu \mathcal{L}_{\text{ELBO}} = \mathbb{E}_{p(\varepsilon)} [\nabla_\theta \log p(\mathbf{y} | \theta) + \log q(\theta; \nu) - \log p(\theta)]_{\theta=\mathcal{T}(\varepsilon; \nu)} \nabla_\nu \mathcal{T}(\varepsilon; \nu) \quad (11)$$

89 Roeder et al. [17] argue that when we believe that $q(\theta; \nu) \approx p(\mathbf{y} | \theta)$, Eq. (11) should be preferred
 90 over Eq. (10) even if computing analytically the KL is possible. Note that this case is very unlikely for
 91 BNN posteriors, and that the additional randomness introduced by the Monte Carlo estimation of the
 92 KL could be harmful.

93 In case of large datasets and complex models, the formulation summarized in Eq. (10) can be
 94 computationally challenging, due to the evaluation of the likelihood and its gradients N_{MC} times.
 95 Assuming factorization of the likelihood,

$$p(\mathbf{y} | \theta) = p(\mathbf{y} | f(\mathbf{X}; \theta)) = \prod_{i=1}^N p(y_i | f(\mathbf{x}_i; \theta)) \quad (12)$$

96 this quantity can be approximated using mini-batching [6, 8]. Recalling \mathbf{y} as the set of labels of
 97 our dataset with N examples, by taking $\mathcal{B} \subset \mathbf{y}$ as a random subset of \mathbf{y} , the likelihood term can be
 98 estimated in an unbiased way as

$$\log p_\theta(\mathbf{y} | \theta) \approx \frac{N}{M} \sum_{y_i \sim \mathcal{B}} \log p(y_i | \theta). \quad (13)$$

99 where M is the number of points in the minibatch. At the cost of increase “randomness”, we can use
 100 Eq. (10) to compute the gradients of the ELBO with the minibatch formulation in Eq. (13). Stochastic
 101 optimization, e.g. any version of SGD, will converge to a local optimum provided with a decreasing
 102 learning rate and sufficient gradient updates [16].

103 **References**

- 104 [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean,
105 M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser,
106 M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens,
107 B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden,
108 M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-Scale Machine Learning on
109 Heterogeneous Systems, 2015. Software available from tensorflow.org.
- 110 [2] S. Ainsworth, J. Hayase, and S. Srinivasa. Git Re-Basin: Merging Models modulo Permutation Symmetries.
111 In *The Eleventh International Conference on Learning Representations*, 2023.
- 112 [3] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational Inference: A Review for Statisticians. 112
113 (518):859–877, 2017.
- 114 [4] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke,
115 J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. JAX: Composable Transformations of Python+NumPy
116 Programs, 2018.
- 117 [5] R. Entezari, H. Sedghi, O. Saukh, and B. Neyshabur. The Role of Permutation Invariance in Linear Mode
118 Connectivity of Neural Networks. In *International Conference on Learning Representations*, 2022.
- 119 [6] A. Graves. Practical Variational Inference for Neural Networks. In J. Shawe-Taylor, R. S. Zemel, P. L.
120 Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*,
121 pages 2348–2356. Curran Associates, Inc., 2011.
- 122 [7] G. E. Hinton and D. van Camp. Keeping the Neural Networks Simple by Minimizing the Description
123 Length of the Weights. In *Proceedings of the Sixth Annual Conference on Computational Learning Theory*,
124 1993. ISBN 0897916115.
- 125 [8] M. D. Hoffman, D. M. Blei, C. Wang, and J. W. Paisley. Stochastic Variational Inference. 14(1):1303–1347,
126 2013.
- 127 [9] M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An Introduction to Variational Methods for
128 Graphical Models. 37(2):183–233, 1999-11-01.
- 129 [10] D. P. Kingma and M. Welling. Auto-Encoding Variational Bayes. In *International Conference on Learning
130 Representations*, 2014.
- 131 [11] A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- 132 [12] C. J. Maddison, A. Mnih, and Y. W. Teh. The Concrete Distribution: A Continuous Relaxation of Discrete
133 Random Variables. In *Proceedings of the 5th International Conference on Learning Representations*.
134 OpenReview.net, 2017.
- 135 [13] N. Metropolis and S. Ulam. The monte carlo method. 44(247):335–341, 1949. PMID: 18139350.
- 136 [14] K. Osawa, S. Swaroop, M. E. E. Khan, A. Jain, R. Eschenhagen, R. E. Turner, and R. Yokota. Practical Deep
137 Learning with Bayesian Principles. In *Advances in Neural Information Processing Systems*, volume 32,
138 pages 4287–4299. Curran Associates, Inc., 2019.
- 139 [15] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimselshein,
140 L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner,
141 L. Fang, J. Bai, and S. Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library.
142 In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett, editors, *Advances
143 in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- 144 [16] H. Robbins and S. Monro. A Stochastic Approximation Method. 22(3):400–407, 1951.
- 145 [17] G. Roeder, Y. Wu, and D. Duvenaud. Sticking the Landing: Simple, Lower-Variance Gradient Estimators
146 for Variational Inference. In I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N.
147 Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages
148 6925–6934, 2017.
- 149 [18] F. R. Ruiz, M. Titsias, and D. Blei. The Generalized Reparameterization Gradient. In D. Lee, M. Sugiyama,
150 U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,
151 volume 29. Curran Associates, Inc., 2016.
- 152 [19] T. Salimans and D. A. Knowles. Fixed-Form Variational Posterior Approximation through Stochastic
153 Linear Regression. 8(4):837–882, 2013.