
Bayesian methods for neural identification in *C. elegans* based on continuous relaxations for permutations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The nematode *C. elegans* is a unique model organism for neuroscientists as its
2 connectome, or neural wiring diagram, has been known for at least three decades.
3 Despite this knowledge, an understanding of the functional significance of these
4 synaptic connections has remained elusive. Now several groups can routinely
5 image the activity of a large fraction of neurons in the head of the worm, providing
6 a unique opportunity to probe this organism. We propose a hierarchical Bayesian
7 framework that combines strong prior information with data from many experi-
8 ments to estimate posteriors over the functional connectivity weights. However,
9 these attempts are stifled by a major obstacle: in many cases it is not clear exactly
10 which neurons are being imaged, so to combine information across experiments
11 one must solve a matching, or permutation inference, problem.

12 In this work we introduce new variational methods that enable the joint inference
13 of connectivity weights and neural identity. Working with actual permutations
14 would involve evaluating and differentiating an intractable partition function. As an
15 alternative, we build upon recent continuous relaxation techniques [Jang et al., 2016,
16 Maddison et al., 2016], extending them from the original case of the probability
17 simplex, to the Birkhoff polytope, the convex hull of permutation matrices. We test
18 our method with simulated data from the true connectome and known covariates
19 (neural position) and show that our approach outperforms many alternatives in
20 identifying neurons.

21 1 Introduction

22 The nematode *C. elegans* plays a special role as a model organism in neuroscience, since its neural net-
23 work is stereotyped from animal to animal and its complete neural wiring diagram is known [Varshney
24 et al., 2011]. Modern calcium imaging technology enables simultaneous measurements of hundreds
25 of these neurons simultaneously [Kato et al., 2015, Nguyen et al., 2016]. Thus, the time seems right to
26 employ modern statistical methods to summarize what can be learned about the functional connectome
27 in this system, and to suggest new experiments to constrain our uncertainty further.

28 To this end, Bayesian inference stands out as the most suited methodological framework, as it allow
29 us to represent hierarchical probabilistic structures, and to integrate our strong priors on the system
30 components (e.g. sparsity patterns in the connectome, and a priori knowledge of approximate neural
31 positions). In the most general setup, we would be interested on posterior inference of a generic
32 dynamical system that dictates the distribution of next neural state, given history, system’s input and
33 behavior. Learning and inference in dynamical systems with MCMC methods is rather standard, even
34 in cases with complicated latent structures De Freitas et al. [2001], Paninski et al. [2010]. Further,
35 methods to account for the hierarchical aspect, i.e., incorporating information from many worms

are also widely available [Gelman et al., 2014]. However, we note a fundamental technical hurdle complicates our efforts of integrating across-specimens information: in practice, associating recorded traces to neuron names is a painstaking process; experimenters consider the location of the neuron along with its pattern of activity to perform this matching, but the process is laborious and the results prone to error. In the lack of this association, it is impossible to represent recordings canonically, where indexes point to actual neuron names, common to all worms. This technical problem, then, prevent us from automatically applying hierarchical methods.

In this work we present a method for overcoming this hurdle, by incorporating inference over *canonicalizing* permutations. For the sake of simplicity, we focus on the most elementary non-trivial dynamical system. Specifically, given the connectome [Varshney et al., 2011] encoded as a $N = 282$ dimensional (number of somatic neurons) adjacency matrix $A \in \{0, 1\}^{N,N}$ (Figure 1A) and J worms, we consider the following hierarchical model with shared linear dynamics (represented by a weight matrix W) to represent the recorded traces $Y_t^{(j)} \in \mathcal{R}^N$ during T timesteps:

$$Y_t^{(j)} = P^{(j)} (W \odot A) Y_{t-1}^{(j)} P^{(j)\top} + \varepsilon_t^{(j)}, \quad (1)$$

$$\varepsilon_t^{(j)} \sim \mathcal{N}(0, I_N), \quad W \sim \mathcal{N}(0, \sigma_w^2 I_N), \quad P^{(j)} \sim \text{Uniform}(\mathcal{P}_N^{(j)}).$$

The operation $W \odot A$ represents the component-wise product with the adjacency matrix, inducing sparseness in the linear system, and reducing dimensionality. For each worm, we represent by the permutation matrix $P^{(j)}$ the matching between recorded indexes and a canonical arbitrary order. These permutations are confined to the a subset $\mathcal{P}_N^{(j)}$, that contains the admissible permutations, based on covariate information. Specifically, we use neural position along the worm’s body to constrain the possible neural identities for a given recorded neuron. We use the known positions of each neuron [Lints et al., 2005], approximating the worm as a one-dimensional object with neurons locations distributed as in Fig. 1B. Then, given reported positions of the neurons (Figure 1B) we can conceive a binary confusion matrix $\mathcal{C}^{(j)}$ so that $\mathcal{C}_{mn}^{(j)} = 1$ (observed) neuron m is close enough to (canonical) neuron n ; i.e., if their distance is smaller than a tolerance η (Figure 1C). Also, absolute certainty of neural identity can be encoded in this matrix, by imposing $\mathcal{C}_{mn}^{(j)} = 1$ in for only true index n .

In this setup, then, we are concerned with joint posterior inference of $p(\{W, P^{(j)}\})$. Although this problem may be also addressed with MCMC, and in practice poor mixing is observed, motivating the use or alternative tools. Here, we cast this problem as an instance of the *variational inference* (VI) framework [Blei et al., 2017] and develop new tools the applicability of this framework to the case where the latent variables are permutation, a case that substantially deviates from the standard practice. In section 2 we detail our VI formulation and summarize our developed methods. Finally, in section 3 we show our findings; notably, the supremacy of our method over the naive MCMC sampler.

2 New methods for variational inference of latent permutations

Consider a latent variable model determined by a prior over the latent $z \sim p(z)$ and a likelihood $p(y|z)$ for the observed data y . In the VI framework, instead of accessing the perhaps intractable posterior $p(z|y)$ one aims to find the distribution $q(z; \nu)$ among a certain variational family, parameterized by $\nu \in \mathcal{V}$, such that it minimizes its discrepancy with $p(z|y)$. Typically, one considers the KL divergence:

$$\nu^* = \arg \min_{\nu \in \mathcal{V}} KL(p(z|y) \| q(z; \nu)). \quad (2)$$

In turn, one can show that the above problem is equivalent to the maximization of the *evidence lower bound* (ELBO):

$$\nu^* = \arg \max_{\nu \in \mathcal{V}} ELBO(q(z; \nu)) \equiv E_{q(z; \nu)}(\log p(y|z)) - KL(q(z; \nu) \| p(z)). \quad (3)$$

To maximize equation (3) one usually appeals to stochastic optimization methods [Kushner and Yin, 1987]: specifically, all the expectations involved in (3) are approximated by Monte Carlo samples, and gradient descent iterations are then performed to this approximation. One critical component is the choice of the Monte Carlo approximation. Perhaps the most common choice is through the so called *score function estimator*, which bases upon the identity $h(\nu) \nabla_\nu \log h(\nu) = \nabla_\nu h(\nu)$. Unfortunately,

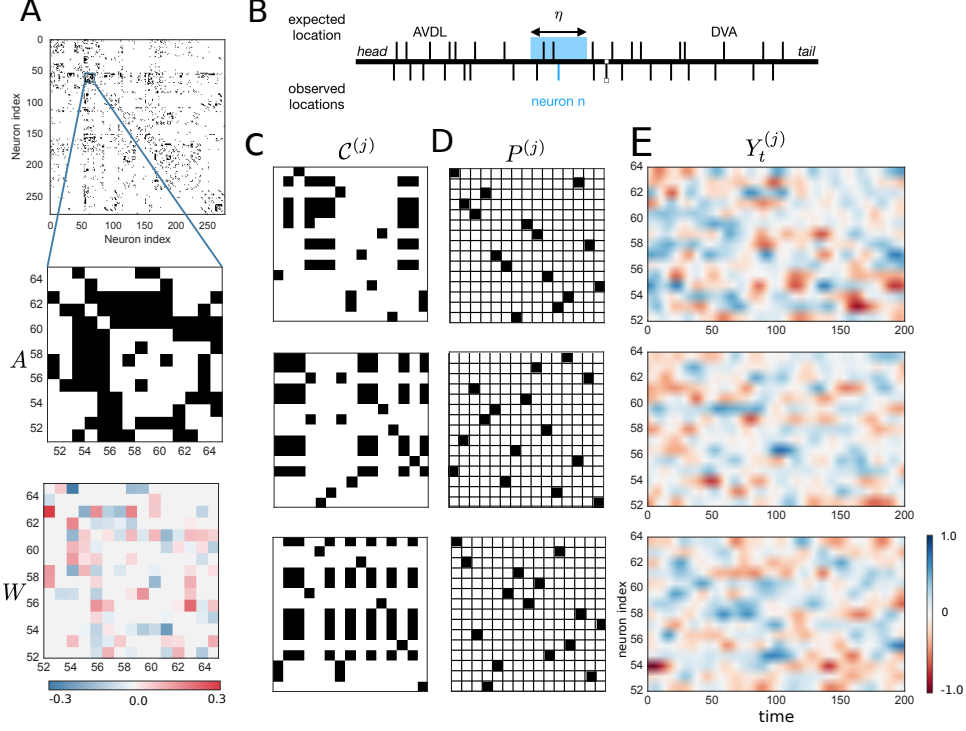


Figure 1: Hierarchical model outline. **A** (top) actual connectome A , from [Varshney et al., 2011], (center) zoom-in to 14 neurons, (bottom) sampled matrix W consistent with A . **B** reference linear location information [White et al., 1986, Lints et al., 2005] can be incorporated in our model to constrain possibilities. **C** in detail, location constraints are represented through a confusion matrix $\mathcal{C}^{(j)}$ that encodes which assignments are possible for each neuron. **D** permutations $P^{(j)}$ are chosen in consistency with constraints. **E** dynamical system observations $Y_t^{(j)}$ are sampled from $P^{(j)}$ -permuted copies of W .

82 this estimator, also referred to as REINFORCE [Williams, 1992], cannot be applied to permutations,
83 since it involves the evaluation and differentiation of a likelihood which is intractable for any non-
84 trivial distribution over permutations (computing the partition function involves a summation over $N!$
85 terms).

86 An appealing alternative comes from the re-parameterization trick Kingma and Welling [2013], which
87 leads to a new gradient estimator if one can re-parameterize z as a differentiable function of a noise
88 distribution and the parameters; i.e., if for certain f and $\xi \sim p(\xi)$ one has $z = f(\xi, \nu)$. In the
89 case of discrete random variables a re-parameterization always exists and it is given by the *Gumbel*
90 *trick* [Papandreou and Yuille, 2011, Balog et al., 2017], which states that one can sample from
91 any discrete distribution by perturbing each potential with Gumbel i.i.d noise, and then finding the
92 configuration with the maximum value. Unfortunately, the underlying f to this re-parameterization is
93 the non-differentiable arg max operator, precluding the use of gradient descent methods.

94 Recent work by [Jang et al., 2016, Maddison et al., 2016] proposed a solution to this problem, by
95 replacing the arg max by a temperature (τ) dependent softmax approximation, which in the limit
96 converges to the original arg max. By combining the Gumbel trick with the softmax approximation,
97 they conceived the *Concrete* or *Gumbel-Softmax* distribution, and obtain explicit distribution formulae.
98 Then, they showed one can learn on a discrete latent variable model using the re-parameterization trick
99 and gradient descent, by replacing the original ELBO with the surrogate arising by this continuous
100 relaxation, as long as τ is chosen in a reasonable range: not too high as it would lead to a degenerate
101 distribution in the simplex; but also not too low, to avoid too high variances of the gradients.

102 We developed three methods for extending the above to permutations. We name them *stick-breaking*,
103 *rounding* and *Gumbel-Sinkhorn* methods. We refer the reader to sections 3.1 and 3.2 of Linderman
104 et al. [2017] and section 4 of Anonymous [2018] for details, respectively. Here we briefly summarize
105 them: in all of them the primary geometric object is the Birkhoff polytope, the convex hull of

permutation matrices, and analog to the probability simplex in this case. For the stick-breaking construction, we generalize to this polytope the one that exists in the simplex [Linderman et al., 2015], surmounting a new complication; of being able to consistently “break the stick” while satisfying both the row and column constraints that characterize a doubly stochastic matrix. For the rounding construction, we start by a noise distribution and force it to be close to permutation matrices by pulling them towards the extreme-points of the Birkhoff polytope. Finally, for the Gumbel-Sinkhorn method we notice that the so-called *Sinkhorn operator*, or infinite and successive row and column normalization of a matrix, is a natural extension of the softmax operator. With this, we are able to conceive the Gumbel-Sinkhorn distribution, which approximates the sampling of a relevant discrete distribution. Importantly, while stick-breaking and rounding yield explicit densities, Gumbel-Sinkhorn does not. However, there are ways to circumvent this difficulty, and overall we observe the latter performs the best.

3 Results

We compared against three methods: (i) naive variational inference, where we do not enforce the constraint that $X^{(j)}$ be a permutation and instead treat each row of $X^{(j)}$ as a Dirichlet distributed vector; (ii) MCMC, where we alternate between sampling from the conditionals of W (Gaussian) and $X^{(j)}$, from which one can sample by proposing local swaps, as described in Diaconis [2009], and (iii) maximum a posteriori estimation (MAP). Our MAP algorithm alternates between the optimizing estimate of W given $\{X^{(m)}, Y^{(m)}\}$ using linear regression and finding the optimal $X^{(j)}$. The second step requires solving a quadratic assignment problem (QAP) in $X^{(j)}$; that is, it can be expressed as $\text{Tr}(AXBX^T)$ for matrices A, B . We used the QAP solver proposed by Vogelstein et al. [2015].

We found that our method outperforms each baseline. Specifically, we show that our method outperforms alternatives when there are many possible candidates (Table 1) and when only a small proportion of neurons are known with certitude (Table 2). Altogether, these results indicate our method enables a more efficient use of information than its alternatives. This is consistent with other results showing faster convergence of variational inference over MCMC [Blei et al., 2017], especially with simple Metropolis-Hastings proposals. We conjecture that MCMC could eventually obtain similar if not better results, if current local proposals—swapping pairs of labels—were replaced by more involved ones.

	10		30		45		60	
	1 worm	4 worms	1 Worm	4 worms	1 worm	4 worms	1 worms	4 worms
NAIVE VI	.34	.32	.16	.16	.13	.12	.11	.12
MAP	.34	.32	.17	.17	.14	.13	.13	.12
MCMC	.34	.65	.18	.28	.14	.17	.13	.15
VI	.79	.94	.4	.69	.25	.51	.21	.44

Table 1: Accuracy in the C.elegans neural identification problem, for varying mean number of candidate neurons (10, 30, 45, 60) and number of worms.

	40.%		30.%		20.%		10.%	
	$\eta = 0.1$	$\eta = 0.2$	$\eta = 0.1$	$\eta = 0.2$	$\eta = 0.1$	$\eta = 0.2$	$\eta = 0.1$	$\eta = 0.2$
Naive VI .43	.41	.33	.31	.23	.22	.12	.1	
MAP	.42	.41	.33	.32	.23	.22	.12	.11
MCMC	.85	.80	.52	.46	.3	.26	.15	.12
VI	.97	.96	.92	.84	.74	.58	.44	.23

Table 2: Accuracy in inferring true neural identity for different of proportion of known neurons, and two values of η .

136 References

- 137 Anonymous. Learning latent permutations with gumbel-sinkhorn networks. *International Conference*
138 *on Learning Representations*, 2018.
- 139 M. Balog, N. Tripuraneni, Z. Ghahramani, and A. Weller. Lost relatives of the gumbel trick. *arXiv*
140 *preprint arXiv:1706.04161*, 2017.
- 141 D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians.
142 *Journal of the American Statistical Association*, 2017.
- 143 N. De Freitas, C. Andrieu, P. Højén-Sørensen, M. Niranjana, and A. Gee. Sequential monte carlo
144 methods for neural networks. In *Sequential Monte Carlo methods in practice*, pages 359–379.
145 Springer, 2001.
- 146 P. Diaconis. The Markov chain Monte Carlo revolution. *Bulletin of the American Mathematical*
147 *Society*, 46(2):179–205, 2009.
- 148 A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. *Bayesian data*
149 *analysis*, volume 2. CRC press Boca Raton, FL, 2014.
- 150 E. Jang, S. Gu, and B. Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint*
151 *arXiv:1611.01144*, 2016.
- 152 S. Kato, H. Kaplan, T. Schrödel, S. Skora, T. Lindsay, E. Yemini, S. Lockery, and M. Zimmer. Global
153 brain dynamics embed the motor command sequence of *Caenorhabditis elegans*. *Cell*, 163(3):656 –
154 669, 2015. ISSN 0092-8674.
- 155 D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*,
156 2013.
- 157 H. Kushner and G. Yin. Stochastic approximation algorithms for parallel and distributed processing.
158 *Stochastics: An International Journal of Probability and Stochastic Processes*, 22(3-4):219–250,
159 1987.
- 160 S. Linderman, M. Johnson, and R. P. Adams. Dependent multinomial models made easy: Stick-
161 breaking with the pólya-gamma augmentation. In *Advances in Neural Information Processing*
162 *Systems*, pages 3456–3464, 2015.
- 163 S. W. Linderman, G. E. Mena, H. Cooper, L. Paninski, and J. P. Cunningham. Reparameterizing the
164 Birkhoff polytope for variational permutation inference. *arXiv preprint arXiv:1710.09508*, 2017.
- 165 R. Lints, Z. F. Altun, H. Weng, T. Stephney, G. Stephney, M. Volaski, and D. H. Hall. WormAtlas
166 Update. 2005.
- 167 C. J. Maddison, A. Mnih, and Y. W. Teh. The concrete distribution: A continuous relaxation of
168 discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- 169 J. P. Nguyen, F. B. Shipley, A. N. Linder, G. S. Plummer, M. Liu, S. U. Setru, J. W. Shaevitz, and A. M.
170 Leifer. Whole-brain calcium imaging with cellular resolution in freely behaving *Caenorhabditis*
171 *elegans*. *Proceedings of the National Academy of Sciences*, 113(8):E1074–E1081, 2016.
- 172 L. Paninski, Y. Ahmadian, D. G. Ferreira, S. Koyama, K. R. Rad, M. Vidne, J. Vogelstein, and W. Wu.
173 A new look at state-space models for neural data. *Journal of computational neuroscience*, 29(1-2):
174 107–126, 2010.
- 175 G. Papandreou and A. L. Yuille. Perturb-and-map random fields: Using discrete optimization to learn
176 and sample from energy models. In *Computer Vision (ICCV), 2011 IEEE International Conference*
177 *on*, pages 193–200. IEEE, 2011.
- 178 L. R. Varshney, B. L. Chen, E. Paniagua, D. H. Hall, and D. B. Chklovskii. Structural properties of
179 the *Caenorhabditis elegans* neuronal network. *PLoS Computational Biology*, 7(2):e1001066, 2011.
- 180 J. T. Vogelstein, J. M. Conroy, V. Lyzinski, L. J. Podrazik, S. G. Kratzer, E. T. Harley, D. E. Fishkind,
181 R. J. Vogelstein, and C. E. Priebe. Fast approximate quadratic programming for graph matching.
182 *PLOS one*, 10(4):e0121002, 2015.

- 183 J. G. White, E. Southgate, J. N. Thomson, and S. Brenner. The structure of the nervous system of
184 the nematode *Caenorhabditis elegans*: the mind of a worm. *Phil. Trans. R. Soc. Lond.*, 314:1–340,
185 1986.
- 186 R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
187 learning. *Machine Learning*, 8(3–4):229–256, 1992.