

High resolution neural connectivity from incomplete tracing data using nonnegative spline regression

Authored by:

Kameron D. Harris
Stefan Mihalas
Eric Shea-Brown

Abstract

Whole-brain neural connectivity data are now available from viral tracing experiments, which reveal the connections between a source injection site and elsewhere in the brain. These hold the promise of revealing spatial patterns of connectivity throughout the mammalian brain. To achieve this goal, we seek to fit a weighted, nonnegative adjacency matrix among 100 μ m brain voxels using viral tracer data. Despite a multi-year experimental effort, injections provide incomplete coverage, and the number of voxels in our data is orders of magnitude larger than the number of injections, making the problem severely underdetermined. Furthermore, projection data are missing within the injection site because local connections there are not separable from the injection signal. We use a novel machine-learning algorithm to meet these challenges and develop a spatially explicit, voxel-scale connectivity map of the mouse visual system. Our method combines three features: a matrix completion loss for missing data, a smoothing spline penalty to regularize the problem, and (optionally) a low rank factorization. We demonstrate the consistency of our estimator using synthetic data and then apply it to newly available Allen Mouse Brain Connectivity Atlas data for the visual system. Our algorithm is significantly more predictive than current state of the art approaches which assume regions to be homogeneous. We demonstrate the efficacy of a low rank version on visual cortex data and discuss the possibility of extending this to a whole-brain connectivity matrix at the voxel scale.

1 Paper Body

Although the study of neural connectivity is over a century old, starting with pioneering neuroscientists who identified the importance of networks for determining brain function, most knowledge of anatomical neural network structure

is limited to either detailed description of small subsystems [2, 9, 14, 26] or to averaged connectivity between larger regions [7, 21]. We focus our attention on spatial, structural connectivity at the mesoscale: a coarser scale than that of single neurons or cortical columns but finer than whole brain regions. Thanks to the development of new tracing techniques, image processing algorithms, and high-throughput methods, data at this resolution are now accessible in animals such as the fly [12, 19] and mouse [15, 18]. We present a novel regression technique tailored to the challenges of learning spatially refined mesoscale connectivity from neural tracing experiments. We have designed this technique with neural data in mind and will use this 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

+
+
+
VISp
+

Figure 1: A, We seek to fit a matrix W which reproduces neural tracing experiments. Each column of W represents the expected signal in target voxels given an injection of one unit into a single source voxel. B, In the work of Oh et al. [18], a regionally homogeneous connectivity matrix was fit using a predefined regional parcellation to constrain the problem. We propose that smoothness of W is a better prior. C, The mouse’s visual field can be represented in azimuth/altitude coordinates. This representation is maintained in the retinotopy, a smoothly varying map replicated in many visual areas (e.g. [8]). D, Assuming locations in VISp (the primary visual area) project most strongly to positions which represent the same retinotopic coordinates in a secondary visual area, then we expect the mapping between upstream and downstream visual areas to be smooth.

language to describe our method, but it is a general technique to assimilate spatial network data or infer smooth kernels of integral equations. Obtaining a spatially-resolved mesoscale connectome will reveal detailed features of connectivity, for example unlocking cell-type specific connectivity and microcircuit organization throughout the brain [13]. In mesoscale anterograde tracing experiments, a tracer virus is first injected into the brain. This infects neurons primarily at their cell bodies and dendrites and causes them to express a fluorescent protein in their cytoplasm, including in their axons. Neurons originating in the source injection site are then imaged to reveal their axonal projections throughout the brain. Combining many experiments with different sources then reveals the pathways that connect those sources throughout the brain. This requires combining data across multiple animals, which appears justified at the mesoscale [18]. We assume there exists some underlying nonnegative, weighted adjacency matrix $W \geq 0$ that is common across animals. Each experiment can be thought of as an injection x , and its projections y , so that $y = Wx$ as in Fig. 1A. Uncovering the unknown W from multiple experiments (x_i, y_i) for $i = 1, \dots, n_{inj}$ is then a multivariate regression problem: Each x_i is an image of the brain which represents the strength of the signal within the injection site.

Likewise, every y_i is an image of the strength of signal elsewhere, which arises due to the axonal projections of neurons with cell bodies in the injection site. The unknown matrix W is a linear operator which takes images of the brain (injections) and returns images of the brain (projections). In a previous paper, Oh et al. [18] were able to obtain a 213×213 regional weight matrix using 469 experiments with mice (Fig. 1B). They used nonnegative least squares to find the unknown regional weights in an overdetermined regression problem. Our aim is to obtain a much higher-resolution connectivity map on the scale of voxels, and this introduces many more challenges. First, the number of voxels in the brain is much larger than the number of injection experiments we can expect to perform; for mouse with $100 \mu\text{m}$ voxels this is $O(10^5)$ versus $O(10^3)$ [15, 18]. Also, the injections that are performed will inevitably leave gaps in their coverage of the brain. Thus specifying W is underdetermined. Second, there is no way to separately image the injections and projections. In order to construct them, experimenters image the brain once by serial tomography and fluorescence

microscopy. The injection sites can be annotated by finding infected cell bodies, but there is no way to disambiguate fluorescence from the cell bodies and dendrites from that of local injections. Projection strength is thus unknown within the injection sites and the neighborhood occupied by dendrites. Third, fitting full-brain voxel-wise connectivity is challenging since the number of elements in W is the square of the number of voxels in the brain. Thus we need compressed representations of W as well as efficient algorithms to perform inference. The paper proceeds as follows. In Section 2, we describe our assumption that the mesoscale connectivity W is smoothly-varying in space, as could be expected from the presence of topographic maps across much of cortex. Later, we show that using this assumption as a prior yields connectivity maps with improved cross-validation performance. In Section 3, we present an inference algorithm designed to tackle the difficulties of underdetermination, missing data, and size of the unknown W . To deal with the gaps and ill-conditioning, we use smoothness as a regularization on W . We take an agnostic approach, similar to matrix completion [5], to the missing projection data and use a regression loss function that ignores residuals within the injection site. Finally, we present a low rank version of the estimator that will allow us to scale to large matrices. In Section 4, we test our method on synthetic data and show that it performs well for sparse data that is consistent with the regression priors. This provides evidence that it is a consistent estimator. We demonstrate the necessity of both the matrix completion and smoothing terms for good reconstruction. In Section 5, we then apply the spline-smoothing method to recently available Allen Institute for Brain Science (Allen Institute) connectivity data from mouse visual cortex [15, 18]. We find that our method is able to outperform current spatially uniform regional models, with significantly reduced cross-validation errors. We also find that a low rank version is able to achieve approximately 23% compression of the original data, with the optimal solution very close to the full rank optimum. Our method is a superior predictor to the existing regional model for visual system data, and the success of the

low rank version suggests that this approach will be able to reveal whole-brain structural connectivity at unprecedented scale. All of our supplemental material and data processing and optimization code is available for download from: <https://github.com/kharris/high-res-connectivity-nips-2016>.

2

Spatial smoothness of mesoscale connectivity

The visual cortex is a collection of relatively large cortical areas in the posterior part of the mammalian brain. Visual stimuli sensed in the retina are relayed through the thalamus into primary visual cortex (VISp), which projects to higher visual areas. We know this partly due to tracing projections between these areas, but also because neurons in the early visual areas respond to visual stimuli in a localized region of the visual field called their receptive fields [11]. An interesting and important feature of visual cortex is the presence of topographic maps of the visual field called the retinotopy [6, 8, 10, 20, 25]. Each eye sees a 2-D image of the world, where two coordinates, such as azimuth and altitude, define a point in the visual field (Fig. 1C). Retinotopy refers to the fact that cells are organized in cortical space by the position of their receptive fields; nearby cells have similar receptive field positions. Furthermore, these retinotopic maps reoccur in multiple visual areas, albeit with varying orientation and magnification. Retinotopy in other areas downstream from VISp, which do not receive many projections directly from thalamus, are likely a function of projections from VISp. It is reasonable to assume that areas which code for similar visual locations are most strongly connected. Then, because retinotopy is smoothly varying in cortical space and similar retinotopic coordinates are the most strongly connected between visual areas, the connections between those areas should be smooth in cortical space (Fig. 1C and D). Retinotopy is a specific example of topography, which extends to other sensory systems such as auditory and somatosensory cortex [22]. For this reason, connectivity may be spatially smooth throughout the brain, at least at the mesoscale. This idea can be evaluated via the methods we introduce below: if a smooth model is more predictive of held-out data than another model, then this supports the assumption. 3

3

Nonnegative spline regression with incomplete tracing data

We consider the problem of fitting an adjacency operator $W : T \rightarrow S \rightarrow R^+$ to data arising from n_{inj} injections into a source space S which projects to a target space T . Here S and T are compact subsets of the brain, itself a compact subset of R^3 . In this mathematical setting, S and T could be arbitrary sets, but typically $S = T$ for the ipsilateral data we present here.¹ The source S and target T are discretized into n_x and n_y cubic voxels, respectively. The discretization of W is then an adjacency $n \times n$ matrix $W \in R^{n_y \times n_x}$. Mathematically, we define the tracing data as a set of pairs $x_i \in R^{n_x}$ and $y_i \in R^{n_y}$, the source and target tracer signals at each voxel for experiments $i = 1, \dots, n_{inj}$. We would like to fit a linear model, a matrix W such that $y_i \approx W x_i$. We assume an observation model $y_i = W x_i + \epsilon_i$ iid

with $\epsilon_i \sim N(0, \Sigma)$ multivariate Gaussian random variables with zero mean

and covariance matrix $\Sigma \in \mathbb{R}^{n_y \times n_y}$. The true data are not entirely linear, due to saturation effects of the fluorescence signal, but the linear model provides a tractable way of ‘credit assignment’ of individual source voxels’ contributions to the target signal [18]. Finally, we assume that the target projections are unknown within the injection site. In other words, we only know y_j outside the support of x_j , which we denote $\text{supp } x_j$, and we wish to only evaluate error for the observable voxels. Let $\mathbf{I} \in \mathbb{R}^{n_y \times n_{inj}}$, where the j th column $\mathbf{I}_j = \mathbf{1}_{\text{supp } x_j}$, the indicator of the complement of the support. We define the orthogonal projector $\mathbf{P} : \mathbb{R}^{n_y \times n_{inj}} \rightarrow \mathbb{R}^{n_y \times n_{inj}}$ as $\mathbf{P}(A) = A - \mathbf{I} \mathbf{I}^T A$, the entrywise product of A and \mathbf{I} . This operator zeros elements of A which correspond to the voxels within each experiment’s injection site. The operator \mathbf{P} is similar to what is used in matrix completion [5], here in the context of regression rather than recovery. These assumptions lead to a loss function which is the familiar ‘2-loss applied to the projected residuals: $\frac{1}{2} \|\mathbf{P}(\mathbf{W}\mathbf{X} - \mathbf{Y})\|_F^2$ (1) $\sum_{j=1}^{n_{inj}}$

where $\mathbf{Y} = y_1, \dots, y_{n_{inj}}$ and $\mathbf{X} = x_1, \dots, x_{n_{inj}}$ are data matrices. Here $\|\cdot\|_F$ is the Frobenius norm, i.e. the ‘2-norm of the matrix as a vector: $\|\mathbf{A}\|_F = \|\text{vec}(\mathbf{A})\|_2$, where $\text{vec}(\mathbf{A})$ takes a matrix and converts it to a vector by stacking consecutive columns. We next construct a regularization penalty. The matrix \mathbf{W} represents the spatial discretization of a two-point kernel W . An important assumption for W is that it is spatially smooth. Function space norms of the derivatives of W , viewed as a real-valued function on $T \times S$, are a natural way to measure the roughness of this function. For this study, we chose the squared L2-norm of the Laplacian $\int \|\nabla W\|^2 dx$, $T \times S$

which is called the thin plate spline bending energy [24]. In the discrete setting, this becomes the squared ‘2-norm of a discrete Laplacian applied to \mathbf{W} :

$\frac{1}{2} \|\mathbf{L} \text{vec}(\mathbf{W})\|_2^2 = \mathbf{L}_y \mathbf{W} + \mathbf{W}^T \mathbf{L}_x \mathbf{F}$. (2) The operator $\mathbf{L} : \mathbb{R}^{n_x \times n_y} \rightarrow \mathbb{R}^{n_x \times n_y}$ is the discrete Laplacian operator or second finite difference matrix on $T \times S$. The equality in Eqn. (2) results from the fact that the Laplacian on the product space $T \times S$ can be decomposed as $\mathbf{L} = \mathbf{L}_x \otimes \mathbf{I}_y + \mathbf{I}_x \otimes \mathbf{L}_y$ [17]. Using the well-known Kronecker product identity for linear matrix equations

$\mathbf{B}^T \otimes \mathbf{A} \text{vec}(\mathbf{X}) = \text{vec}(\mathbf{Y}) \iff \mathbf{A}\mathbf{X}\mathbf{B} = \mathbf{Y}$ (3) gives the result in Eqn. (2) [23], which allows us to efficiently evaluate the Laplacian action. As for boundary conditions, we do not want to impose any particular values at the boundary, so we choose the finite difference matrix corresponding to a homogeneous Neumann (zero derivative) boundary condition. $\frac{1}{2} \|\cdot\|_{\text{Ipsilateral}}$ refers to connections within the same cerebral hemisphere. For contralateral (opposite hemisphere) connectivity, S and T are disjoint subsets of the brain corresponding to the two hemispheres. $\frac{1}{2}$ It is straightforward to avoid smoothing across region boundaries by imposing Neumann boundary conditions at the boundaries; this is an option in our code available online.

4

Combining the loss and penalty terms, Eqn. (1) and (2), gives a convex optimization problem for inferring the connectivity:

n_{inj}

$$\mathbf{L}_y \mathbf{W} + \mathbf{W}^T \mathbf{L}_x \mathbf{F} \cdot (\mathbf{P1}) \mathbf{W} = \arg \min_{\mathbf{W}} \frac{1}{2} \|\mathbf{P}(\mathbf{W}\mathbf{X} - \mathbf{Y})\|_F^2 + \frac{\lambda}{2} \|\mathbf{L} \text{vec}(\mathbf{W})\|_2^2$$

In the final form, we absorb the noise variance σ^2 into the regularization hyperparameter λ and rescale the penalty so that it has the same dependence on the problem size n_x , n_y , and n_{inj} as the loss. We solve the optimization (P1) using the L-BFGS-B projected quasi-Newton method, implemented in C++ [3, 4]. The gradient is efficiently computed using matrix algebra. Note that (P1) is a type of nonnegative least squares problem, since we can use Eqn. (3) to convert it into $\min_{\mathbf{w}} \|\mathbf{A}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$, with

where $\mathbf{A} = \text{diag}(\text{vec}(\mathbf{W}))^T \mathbf{X}^T \mathbf{I}_{n_{inj}} \mathbf{y}$, $\mathbf{y} = \text{diag}(\text{vec}(\mathbf{W}))^T \text{vec}(\mathbf{Y})$, and $\mathbf{w} = \text{vec}(\mathbf{W})$. Furthermore, without the nonnegativity constraint the estimator is linear and has an explicit solution. However, the design matrix \mathbf{A} will have dimension $(n_{inj} \times (n_y n_x))$, with $O(n_y^3 n_{inj})$ entries if $n_x = O(n_y)$. The dimensionality of the problem prevents us from working directly in the tensor product space. And since the model is a structured matrix regression problem [1], the usual representer theorems [24], which reduce the dimensionality of the estimator to effectively the number of data points, do not immediately apply. However, we hope to elucidate the connection to reproducing kernel Hilbert spaces in future work.

3.1 Low rank version

The largest object in our problem is the unknown connectivity \mathbf{W} , since in the underconstrained setting $n_{inj} n_x, n_y$. In order to improve the scaling of our problem with the number of voxels, we reformulate it with a compressed version of \mathbf{W} :

$$\min_{\mathbf{U}, \mathbf{V}} \|\mathbf{U} \mathbf{V}^T - \mathbf{W}\|_F^2 + \lambda \|\mathbf{U}\|_F^2 + \lambda \|\mathbf{V}\|_F^2 \quad (\text{P2}) \quad (\mathbf{U}^T, \mathbf{V}^T) = \arg \min_{\mathbf{U}^T, \mathbf{V}^T} \|\mathbf{U} \mathbf{V}^T - \mathbf{W}\|_F^2 + \lambda \|\mathbf{U}\|_F^2 + \lambda \|\mathbf{V}\|_F^2$$

Here, $\mathbf{U} \in \mathbb{R}^{n_y \times r}$ and $\mathbf{V} \in \mathbb{R}^{n_x \times r}$ for some fixed rank r , so that the optimal connectivity $\mathbf{W} = \mathbf{U} \mathbf{V}^T$ is given in low rank, factored form. Note that we use nonnegative factors rather than constrain $\mathbf{U}, \mathbf{V} \geq 0$, since this is a nonlinear constraint. This has the advantage of automatically computing a nonnegative matrix factorization (NMF) of \mathbf{W} . The NMF is of separate scientific interest, to be pursued in future work, since it decomposes the connectivity into a relatively small number of projection patterns, which has interpretations as a clustering of the connectivity itself. In going from the full rank problem (P1) to the low rank version (P2), we lose convexity. So the usual optimization methods are not guaranteed to find a global optimum, and the clustering just mentioned is not unique. However, we have also reduced the size of the unknowns to the potentially much smaller matrices \mathbf{U} and \mathbf{V} , if $r \ll n_y, n_x$. If $n_x = O(n_y)$, we have only $O(n_y r)$ unknowns instead of $O(n_y^2)$. Evaluating the penalty term still requires computation of $n_y n_x$ terms, but this can be performed without storing them in memory. We use a simple projected gradient method with Nesterov acceleration in Matlab to find a local optimum for (P2) [3], and will present and compare these results to the solution of (P1) below. As before, computing the gradients is efficient using matrix algebra. This method has been used before for NMF [16].

4

Test problem

We next apply our algorithms to a test problem consisting of a one-dimensional
 $\varphi_{\text{brain},?}$ where the source and target space $S = T = [0, 1]$. The true connectivity kernel corresponds to a Gaussian profile about the diagonal plus a bump: $(\frac{1}{2} - |x - y|)^2 \exp(-(x - 0.8)^2) + (y - 0.1)^2$
 $W_{\text{true}}(x, y) = \exp(-x^2) + 0.9 \exp(-(y - 0.4 - 0.2)^2)$

W_{full} without φ

W_{true} 1

W_{full} with $\varphi=0$

W_{rank} 20

1 0.9

0.9 0.8

0.8

1

0.9 6 0.8 5

0.7 0.7

0.7

0.6

0.6

0.5

0.5

0.4

0.4

0.3

0.3

0.9

0.8

0.8

0.8

0.8

0.7

0.7

0.7

0.6

4

0.6 0.5

0.6 0.5

0.5

0.4

0.4

3

0.5 0.4

0.4

0.3

0.2

0.9

0.6

0.5
0.2
0.9
0.7
0.6
1
0.9
0.3
2
0.2
0.2
0.1
0.4 0.3
0.3
0.3
0.2
0.1
0.2
0.1
0
0.1
0
0.2
1 0.1
0.1
0
0 0
0.2
0.4
0.6
0.8
1
0 ?0.1
0 0
0.2
0.4
0.6
0.8
?0.1
1
0 0
0.2
0.4
0.6
0.8
0.1

0.1
1
0 0
0.2
0.4
0.6
0.8
1

Figure 2: Comparison of the true (far left) and inferred connectivity from 5 injections. Unless noted, $\eta = 100$. Second from left, we show the what happens when we solve (P1) without the matrix completion term P^η . The holes in the projection data cause patchy and incorrect output. Note the colorbar range is 6% that in the other cases. Second from right is the result with P^η but without regularization, solving (P1) for $\eta = 0$. There, the solution does not interpolate between injections. Far right is a rank $r = 20$ result using (P2), which captures the diagonal band and off-diagonal bump that make up W^{true} . In this case, the low rank result has less relative error (9.6%) than the full rank result (11.1%, not shown).

See the left panel of Fig. 2. The input and output spaces were discretized using $n_x = n_y = 200$ points. Injections are delivered at random locations within S , with a width of $0.12 + 0.1$ where $\eta \sim \text{Uniform}(0, 1)$. The values of x are set to 1 within the injection region and 0 elsewhere, y is set to 0 within the injection region, and we take noise level $\eta = 0.1$. The matrices $L_x = L_y$ are the 5-point finite difference Laplacians for the rectangular lattice. Example output of (P1) and (P2) is given for 5 injections in Fig. 2. Unless stated otherwise, $\eta = 100$. The injections, depicted as black bars in the bottom of each sub-figure, do not cover the whole space S but do provide good coverage of the bump, otherwise there is no information about that feature. We depict the result of the full rank algorithm (P1) without the matrix completion term P^η , the result including P^η but without smoothing ($\eta = 0$), and the result of (P2) with rank $r = 20$. The full rank solution is not shown, but is similar to the low rank one. Figure 2 shows the necessity of each term within the algorithm. Leaving out the matrix completion P^η leads to dramatically biased output since the algorithm uses incorrect values $y_{\text{supp}}(x) = 0$. If we include P^η but neglect the smoothing term by setting $\eta = 0$, we also get incorrect output: without smoothing, the algorithm cannot fill in the injection site holes nor can it interpolate between injections. However, the low rank result accurately approximates the true connectivity W^{true} , including the diagonal profile and bump, achieving 9.6% relative error measured as $\|W - W^{\text{true}}\|_F / \|W^{\text{true}}\|_F$. The full rank version is similar, but in fact has slightly higher 11.1% relative error.

5

Finding a voxel-scale connectivity map for mouse cortex

We next apply our method to the latest data from the Allen Institute Mouse Brain Connectivity Atlas, obtained with the API at <http://connectivity.brain-map.org>. Briefly, in each experiment mice were injected with adeno-associated virus expressing a fluorescent protein. The virus infects neurons in the injection

site, causing them to produce the protein, which is transported throughout the axonal and dendritic processes. The mouse brains for each experiment were then sliced, imaged, and aligned onto the common coordinates in the Allen Reference Atlas version 3 [15, 18]. These coordinates divide the brain volume into $100 \mu\text{m} \times 100 \mu\text{m} \times 100 \mu\text{m}$ voxels, with approximately 5×10^5 voxels in the whole brain. The fluorescent pixels in each aligned image were segmented from the background, and we use the fraction of segmented versus total pixels in a voxel to build the vectors x and y . Since cortical dendrites project locally, the signal outside the injection site is mostly axonal, and so the method reveals anterograde axonal projections from the injection site. From this dataset, we selected 28 experiments which have 95% of their injection volumes contained within the visual cortex (atlas regions VISal, VISam, VISl, VISp, VISpl, VISpm, VISli, VISpor, VISrl, and VISa) and injection volume less than 0.7 mm^3 . For this study, we present only the results for ipsilateral connectivity, where $S = T$ and $n_x = n_y = 7497$. To compute the smoothing penalty, we used the 7-point finite-difference Laplacian on the cubic voxel lattice. 6

Model Voxel MSErel Regional MSErel Regional 107% (70%) 48% (6.8%) Voxel 33% (10%) 16% (2.3%) Table 1: Model performance on Allen Institute Mouse Brain Connectivity Atlas data. Cross-validation errors of the voxel model (P1) and regionally homogeneous models are shown, with training errors in parentheses. The errors are computed in both voxel space and regional space, using the relative mean squared error MSErel, Eqn. (4). In either space, the voxel model shows reduced training and cross-validation errors relative to the regional model.

In order to evaluate the performance of the estimator, we employ nested cross-validation with 5 inner and outer folds. The full rank estimator (P1) was fit for $\lambda = 103, 104, \dots, 1012$ on the training data. Using the validation data, we then selected the λ_{opt} that minimized the mean square error relative to the average squared norm of the prediction $W X$ and truth Y , evaluating errors outside the injection sites: $2kP^? (W X - Y)^2_{k2F}$. (4) $\text{MSErel} = kP^? (W X)^2_{k2F} + kP^? (Y)^2_{k2F}$ This choice of normalization prevents experiments with small kY from dominating the error. This error metric as well as the ℓ_2 -loss adopted in Eqn. (P1) both more heavily weight the experiments with larger signal. After selection of λ_{opt} , the model was refit to the combined training and validation data. In our dataset, $\lambda_{\text{opt}} = 105$ was selected for all outer folds. The final errors were computed with the test datasets in each outer fold. For comparison, we also fit a regional model within the cross-validation framework, using nonnegative least squares. To do this, similar to the study by Oh et al. [18], we constrained the connectivity $W_{kl} = W_{Ri} R_j$ to be constant for all voxels k in region R_i and l in region R_j . The results are shown in Table 1. Errors were computed according to both voxels and regions. For the latter, we integrated the residual over voxels within the regions before computing the error. The voxel model is more predictive of held-out data than the regional model, reducing the voxel and regional MSErel by 69% and 67%, respectively. The regional model is designed for inter-region connectivity. To allow an easier comparison with the voxel model, we here include within

region connections. We find that the regional model is a poor predictor of voxel scale projections, with over 100% relative voxel error, but it performs okay at the regional scale. The training errors, which reflect goodness of fit, were also reduced significantly with the voxel model. We conclude that the more flexible voxel model is a better estimator for these Allen Institute data, since it improves both the fits to training data as well as cross-validation skill. The inferred visual connectivity also exhibits a number of features that we expect. There are strong local projections (similar to the diagonal in the test problem, Fig. 2) along with spatially organized projections to higher visual areas. See Fig. 3, which shows example projections from source voxels within VISp. These are just two of 7497 voxels in the full matrix, and we depict only a 2-D projection of 3-D images. The connectivity exhibits strong local projections, which must be filled in by the smoothing since within the injection sites the projection data are unknown; it is surprising how well the algorithm does at capturing short-range connectivity that is translationally invariant. There are also long-range bumps in the higher visual areas, medial and lateral, which move with the source voxel. This is a result of retinotopic maps between VISp and downstream areas. The supplementary material presents a view of this high-dimensional matrix in movie form, allowing one to see the varying projections as the seed voxel moves. We encourage the reader to view the supplemental movies, where movement of bumps in downstream regions hints at the underlying retinotopy: <https://github.com/kharris/high-res-connectivity-nips-2016>. 5.1

Low rank inference successfully approximates full rank solution for visual system

We next use these visual system data, for which the full rank solution was computed, to test whether the low rank approximation can be applied. This is an important stepping stone to an eventual inference of spatial connectivity for the full brain. First, we note that the singular value spectrum of the fitted W_{full} (now using all 28 injections and 5×10^5) is heavily skewed: 95% of the energy can be captured with 21 of 7497 components, and 99% with 67 components. However, this does not directly imply that a nonnegative factorization will

7

Figure 3: Inferred connectivity using all 28 selected injections from visual system data. Left, Projections from a source voxel (blue) located in VISp to all other voxels in the visual areas. The view is integrated over the superior-inferior axis. The connectivity shows strong local connections and weaker connections to higher areas, in particular VISam, VISal, and VISl. Movies of the inferred connectivity (full, low rank, and the low rank residual) for varying source voxel are available in the supplementary material. Center, For a source 800 μ m voxels away, the pattern of anterograde projections is similar, but the distal projection centers are shifted, as expected from retinotopy. Right, The residuals between the full rank and rank 160 result from solving (P2), for the same source voxel as in the center. The residuals are an order of magnitude less than typical features of the connectivity. To test this, we fit a low rank decomposition directly to all 28 visual injection data using (P2) with rank $r = 160$ and $\epsilon = 10^{-5}$. The output of the optimization procedure yields U and

$\|V\|_F$, and we find that the low rank output is very similar to the full result W full fit to the same data (see also Fig. 3, which visualizes the residuals): $\|kU - V\|_F / \|W\|_F = 13\%$. $\|kW - W\|_F / \|W\|_F$ This close approximation is despite the fact that the low rank solution achieves a roughly 23% compression of the 7497 \times 7497 matrix.

Assuming similar compressibility for the whole brain, where the number of voxels is 5×10^5 , would mean a rank of approximately 104. This is still a problem in $O(10^9)$ unknowns, but these bring the memory requirements of storing one matrix iterate in double precision from approximately 1.9 TB to 75 GB, which is within reach of commonly available large memory machines.

6

Conclusions

We have developed and implemented a new inference algorithm that uses modern machine learning ideas: matrix completion loss, a smoothing penalty, and low rank factorization to assimilate sparse connectivity data into complete, spatially explicit connectivity maps. We have shown that this method can be applied to the latest Allen Institute data from multiple visual cortical areas, and that it significantly improves cross-validated predictions over the current state of the art and unveils spatial patterning of connectivity. Finally, we show that a low rank version of the algorithm produces very similar results on these data while compressing the connectivity map, potentially opening the door to the inference of whole brain connectivity from viral tracer data at the voxel scale.

Acknowledgements We acknowledge the support of the UW NIH Training Grant in Big Data for Neuroscience and Genetics (KDH), Boeing Scholarship (KDH), NSF Grant DMS-1122106 and 1514743 (ESB & KDH), and a Simons Fellowship in Mathematics (ESB). We thank Liam Paninski for helpful insights at the outset of this project. We wish to thank the Allen Institute founders, Paul G. Allen and Jody Allen, for their vision, encouragement, and support. This work was facilitated through the use of advanced computational, storage, and networking infrastructure provided by the Hyak supercomputer system at the University of Washington.

2 References

[1] Argyriou, A., Micchelli, C. A., and Pontil, M. (2009). When Is There a Representer Theorem? Vector Versus Matrix Regularizers. *J. Mach. Learn. Res.*, 10:2507–2529.

8

[2] Bock, D. D., Lee, W.-C. A., Kerlin, A. M., Andermann, M. L., Hood, G., Wetzel, A. W., Yurgenson, S., Soucy, E. R., Kim, H. S., and Reid, R. C. (2011). Network anatomy and in vivo physiology of visual cortical neurons. *Nature*, 471(7337):177–182. [3] Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press, New York, NY, USA. [4] Byrd, R., Lu, P., Nocedal, J., and Zhu, C. (1995). A Limited Memory Algorithm for Bound Constrained Optimization. *SIAM Journal on Scientific Computing*,

16(5):1190?1208. [5] Candes, E. and Tao, T. (2010). The Power of Convex Relaxation: Near-Optimal Matrix Completion. *IEEE Transactions on Information Theory*, 56(5):2053?2080. [6] Chaplin, T. A., Yu, H.-H., and Rosa, M. G. (2013). Representation of the visual field in the primary visual area of the marmoset monkey: Magnification factors, point-image size, and proportionality to retinal ganglion cell density. *Journal of Comparative Neurology*, 521(5):1001?1019. [7] Felleman, D. J. and Van Essen, D. C. (1991). Distributed Hierarchical Processing in the Primate. *Cerebral Cortex*, 1(1):1?47. [8] Garrett, M. E., Nauhaus, I., Marshel, J. H., and Callaway, E. M. (2014). Topography and Areal Organization of Mouse Visual Cortex. *Journal of Neuroscience*, 34(37):12587?12600. [9] Glickfeld, L. L., Andermann, M. L., Bonin, V., and Reid, R. C. (2013). Cortico-cortical projections in mouse visual cortex are functionally target specific. *Nature Neuroscience*, 16(2). [10] Goodman, C. S. and Shatz, C. J. (1993). Developmental mechanisms that generate precise patterns of neuronal connectivity. *Cell*, 72, Supplement:77?98. [11] Hubel, D. H. and Wiesel, T. N. (1962). Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of Physiology*, 160(1):106?154.2. [12] Jenett, A., Rubin, G. M., Ngo, T.-T. B., Shepherd, D., Murphy, C., Dionne, H., Pfeiffer, B. D., Cavallaro, A., Hall, D., Jeter, J., Iyer, N., Fetter, D., Hausenfluck, J. H., Peng, H., Trautman, E. T., Svirska, R. R., Myers, E. W., Iwinski, Z. R., Aso, Y., DePasquale, G. M., Enos, A., Hulamm, P., Lam, S. C. B., Li, H.-H., Lavery, T. R., Long, F., Qu, L., Murphy, S. D., Rokicki, K., Safford, T., Shaw, K., Simpson, J. H., Sowell, A., Tae, S., Yu, Y., and Zugates, C. T. (2012). A GAL4-Driver Line Resource for *Drosophila* Neurobiology. *Cell Reports*, 2(4):991?1001. [13] Jonas, E. and Kording, K. (2015). Automatic discovery of cell types and microcircuitry from neural connectomes. *eLife*, 4:e04250. [14] Kleinfeld, D., Bharioke, A., Blinder, P., Bock, D. D., Briggman, K. L., Chklovskii, D. B., Denk, W., Helmstaedter, M., Kaufhold, J. P., Lee, W.-C. A., Meyer, H. S., Micheva, K. D., Oberlaender, M., Prohaska, S., Reid, R. C., Smith, S. J., Takemura, S., Tsai, P. S., and Sakmann, B. (2011). Large-Scale Automated Histology in the Pursuit of Connectomes. *The Journal of Neuroscience*, 31(45):16125?16138. [15] Kuan, L., Li, Y., Lau, C., Feng, D., Bernard, A., Sunkin, S. M., Zeng, H., Dang, C., Hawrylycz, M., and Ng, L. (2015). Neuroinformatics of the Allen Mouse Brain Connectivity Atlas. *Methods*, 73:4?17. [16] Lin, C.-J. (2007). Projected Gradient Methods for Nonnegative Matrix Factorization. *Neural Computation*, 19(10):2756?2779. [17] Lynch, R. E., Rice, J. R., and Thomas, D. H. (1964). Tensor product analysis of partial difference equations. *Bulletin of the American Mathematical Society*, 70(3):378?384. [18] Oh, S. W., Harris, J. A., Ng, L., Winslow, B., Cain, N., Mihalas, S., Wang, Q., Lau, C., Kuan, L., Henry, A. M., Mortrud, M. T., Ouellette, B., Nguyen, T. N., Sorensen, S. A., Slaughterbeck, C. R., Wakeman, W., Li, Y., Feng, D., Ho, A., Nicholas, E., Hirokawa, K. E., Bohn, P., Joines, K. M., Peng, H., Hawrylycz, M. J., Phillips, J. W., Hohmann, J. G., Wohnoutka, P., Gerfen, C. R., Koch, C., Bernard, A., Dang, C., Jones, A. R., and Zeng, H. (2014). A mesoscale connectome of the mouse brain. *Nature*, 508(7495):207?214. [19] Peng, H., Tang, J., Xiao, H., Bria, A., Zhou, J., Butler, V., Zhou, Z., Gonzalez-Bellido, P. T., Oh, S. W., Chen, J., Mitra, A., Tsien, R.

W., Zeng, H., Ascoli, G. A., Iannello, G., Hawrylycz, M., Myers, E., and Long, F. (2014). Virtual finger boosts three-dimensional imaging and microsurgery as well as terabyte volume image visualization and analysis. *Nature Communications*, 5. [20] Rosa, M. G. and Tweedale, R. (2005). Brain maps, great and small: Lessons from comparative studies of primate visual cortical organization. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 360(1456):665-691. [21] Sporns, O. (2010). *Networks of the Brain*. The MIT Press, 1st edition. [22] Udin, S. B. and Fawcett, J. W. (1988). Formation of Topographic Maps. *Annual Review of Neuroscience*, 11(1):289-327. [23] Van Loan, C. F. (2000). The ubiquitous Kronecker product. *Journal of Computational and Applied Mathematics*, 123(1-2):85-100. [24] Wahba, G. (1990). *Spline Models for Observational Data*. SIAM. [25] Wang, Q. and Burkhalter, A. (2007). Area map of mouse visual cortex. *The Journal of Comparative Neurology*, 502(3):339-357. [26] White, J. G., Southgate, E., Thomson, J. N., and Brenner, S. (1986). The Structure of the Nervous System of the Nematode *Caenorhabditis elegans*. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, 314(1165):1-340.