# Diagnosing Forward Operator Error Using Optimal Transport

**Michael A. Puthawala, Cory D. Hauck, and Stanley J. Osher**

October 19, 2018

## 1 Abstract

We investigate overdetermined linear inverse problems for which the forward operator may not be given accurately. We introduce a new tool called the *structure*, based on the Wasserstein distance, and propose the use of this to diagnose and remedy forward operator error. Computing the structure turns out to use an easy calculation for a Euclidean homogeneous degree one distance, the Earth Mover's Distance, based on recently developed algorithms. The structure is proven to distinguish between noise and signals in the residual and gives a plan to help recover the true direct operator

Micheal A. Puthawala
Department of Mathematics University of California Los Angeles, CA 90095
E-mail: mputhawala@ucla.edu
This author's research was sponsored by Department of Energy grant DOE-SC0013838 and NSF (STROBE), NSFC 11671005.

Cory D. Hauck
Computational Mathematics Group, Computer Science and Mathematics Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831 USA
E-mail: hauckc@ornl.gov
This author's research was sponsored by the Office of Advanced Scientific Computing Research and performed at the Oak Ridge National Laboratory, which is managed by UT-Battelle, LLC under Contract No. DE-AC05-00OR22725.

Stanley J. Osher
Department of Mathematics University of California Los Angeles, CA 90095
E-mail: sjo@math.ucla.edu
This author's research was sponsored by Department of Energy grant DOE-SC0013838 and NSF (STROBE), NSFC 11671005.

9  in some interesting cases. We expect to use this technique not only to diagnose the
10 error, but also to correct it, which we do in some simple cases presented below.


## 2 Introduction

### 2.1 Motivation

From medical imaging [1] to petroleum engineering [24] to meteorology [4], inverse
problems are ubiquitous in science, engineering and mathematics. The goal of such
problems is to recover an unknown quantity $u$ given a known forward operator $L$
and measurement $b$ such that $L(u) = b$. In this work we consider the case where $L$ is
a linear operator and write $L(u) \equiv Lu$. While this choice facilitates a simple analysis
in some places, the computational techniques developed here can be extended to
consider non-linear operators.

A considerable amount of work has been dedicated to solving inverse problems
for a variety of forward operators, especially when $L$ is linear. Powerful techniques
have been developed that perform well in the presence of noise in $b$, singularities in
$L$ and various constraints on the solution $u$ [20].

Despite some great successes in the field of inverse problems, there are still math-
ematical challenges that are difficult to address. One of these, which is important in a
bevy of applications, is the calibration of forward operators. For example, computed
tomography (CT) machines are calibrated using known phantoms for which the de-
sired reconstruction is known exactly [30]; in synthetic aperture radar, reflectors
provide a known ground truth on which devices and reconstruction algorithms are
tuned [12]; and in some plasma imaging problems, the forward model has unknown
parameters, and the model itself is possibly incomplete [33].

Often the calibration problem can be formulated mathematically by considering
a family of forward operators $L_\theta$, parameterized by $\theta \in \Theta \subset \mathbb{R}^p$, with a unique $\hat{\theta}$ such
that $L_{\hat{\theta}}$ best represents the underlying physical system. In other words, there exists a
$\hat{\theta}$ such that $L = L_{\hat{\theta}}$ [29,33]. If $\hat{\theta}$ is estimated poorly, then an accurate approximation
of $u$ is often impossible, even with very sophisticated inverse procedures.

The problem of detecting forward operator error is similar to that of blind decon-
volution in image processing [5], where the task is to identify a blurring kernel and
recover an image from a given blurry signal. The application of the blurring operator
with the image can also be represented in the form $Lu = b$ where the action of $L$
gives the convolution with the blurring kernel. One important difference between the
calibration problem considered here and the problem of blind deconvolution is that
we will be considering overdetermined problems.


### 2.2 Prior Work

Methods for detecting and correcting for errors within the forward operator exist.
One approach is total least squares [16], which generalizes the standard least squares
method by allowing for error in $L$. This is expressed by the minimization problem

$$\min_{\mathbf{v},\mathbf{J}} \|\mathbf{L} - \mathbf{J}\|_F^2 + \|\mathbf{b} - \mathbf{J}\mathbf{v}\|_2^2 , \tag{1}$$

where $\mathbf{L}$ is the matrix representations of $L$, $\mathbf{b}$ is the vector representation of $b$, and $\|\cdot\|_F$ is the Frobenius norm.

This approach has the advantage of being relatively easy to analyze, robust under noise in the entries of $\mathbf{L}$ and solvable using standard linear algebra software. However, for calibration problems, the goal is not to remove entry-wise error in $\mathbf{L}_\theta$. Instead we seek a value of $\theta \approx \hat{\theta}$. Total least squares provides good reconstructions when $\mathbf{L}$ is a matrix whose entries are corrupted by noise. However it requires modification in order to be applied to the parametric calibration problem. In particular, adding the requirement $J = L_\theta$ for $\theta \in \Theta$ to Eq. 1 make the resulting minimization problem more difficult to solve, and so may require code beyond standard linear algebra software.

Another common approach for calibration is based on Bayesian techniques [19]. In this setting measured data (possibly noisy) is assumed to be the sum of model output and a discrepancy function, both of which are modeled as Gaussian processes. We do not go into details of the Bayesian approach in this paper but intend to make comparisons with the EMD approach in future work. However, it is worth noting that the results in this paper do not rely on a Gaussian noise model.

Our work is motivated in part by [7,8,34], where the authors use the quadratic Wasserstein metric to solve Full-Waveform Inversion (FWI) problems. In particular, it is demonstrated that the quadratic Wasserstein metric, as opposed to the $L_2$ norm, provides an effective measure of the misfit between given data and computed solution.

## 2.3 Our contribution

In this paper we introduce a new tool, called the structure, that is based on the Earth Mover's Distance (EMD) from optimal transport. We show that the structure is sensitive to modeling errors in $L$, but insensitive to noise in $b$. For simple functional forms of $L_\theta$, we demonstrate that the structure can successfully recover the correct parameter $\hat{\theta}$. The method can be implemented as a wrapper around existing inverse problem solvers and thus can be easily integrated into preexisting work flows for solving inverse problems with minimal modifications to existing code bases. Moreover, due to recent advancements in the calculation of the EMD [21,22], the additional cost is reasonable.

Our work extends that of [7,8,34] by considering different inverse problems, a more general noise model, and we use a different Wasserstein metric. See section 4.4 for more detail. We also show that new algorithms for computing the EMD can be combined with inverse problem solvers to diagnose forward operator error in general inverse problems.

## 3 Background

### 3.1 Inverse Problems

Let $\mathcal{U} \subset L^\infty(X)$ and $\mathcal{B} \subset L^\infty(Y)$ be function spaces defined over bounded rectangular domains $X \subset \mathbb{R}^{d_x}$ and $Y \subset \mathbb{R}^{d_y}$, respectively. We consider problems which come

from the discretization of the linear equation

$$\mathcal{L}f = g \tag{2}$$

where $f \in \mathcal{U}$, $g \in \mathcal{B}$, and $\mathcal{L} : \mathcal{U} \to \mathcal{B}$ is a bounded linear operator.

To discretize Eq. 2, we assume that for some $\Delta x > 0$ and $\Delta y > 0$, $X$ and $Y$ can be partitioned into hypercubes $K^x$ and $K^y$, respectively, of size $= \Delta x^{d_y}$ and $\Delta y^{d_y}$, respectively, such that $X = \cup_i \overline{K_i^x}$ and $Y = \cup_j \overline{K_j^y}$. We then let

$$\mathcal{U}_{\Delta x} = \{f_{\Delta x} \in \mathcal{U} : f_{\Delta x}|_{K_x} \text{ is constant for all } K_x \subset X\} \tag{3}$$

$$\mathcal{B}_{\Delta y} = \{g_{\Delta y} \in \mathcal{B} : g_{\Delta y}|_{K_y} \text{ is constant for all } K_y \subset Y\}. \tag{4}$$

The discrete version of Eq. 2 takes the form

$$Lu = b, \tag{5}$$

where $u \in \mathcal{U}_{\Delta x}$, $b \in \mathcal{B}_{\Delta y}$, and $L : \mathcal{U}_{\Delta x} \to \mathcal{B}_{\Delta y}$ is a bounded linear operator that approximates $\mathcal{L}$. The exact forms of $L$, $u$, and $b$ depend on the discretization. In the appendix, we present a discretization based on the assumption that $\mathcal{L}$ is generated by line integrals over paths $\mathcal{P}_y \subset X$ that are parameterized by elements $y \in Y$.

Solving Eq. 5 directly may not be practical if the condition number of $L$ is large, as noise in $b$ can be strongly amplified in the inversion process. A variational approach to address this difficulty is instead to solve

$$\tilde{u} = \tilde{L}^{-1}b \equiv \underset{v \in \mathcal{U}_{\Delta x}}{\operatorname{argmin}} \|Lv - b\|_2^2 + \Phi(v; \lambda), \tag{6}$$

where $\Phi : \mathcal{U}_{\Delta x} \to \mathbb{R}^+$ is a regularizing functional with parameter $\lambda \in \mathbb{R}^+$. If $\Phi = 0$, then Eq. 6 gives the least squares solution of Eq. 5. Nontrivial examples of $\Phi$ (which may require more regularity than $L^\infty(X)$) include

1. $\Phi(v; \lambda) = \lambda \|Cv\|_2^2$, where the linear operator $C$ approximates a differential operator (Generalized Tikhonov regularization);
2. $\Phi(v; \lambda) = \lambda \operatorname{TV}(v)$ (Total Variation regularization [27]);
3. $\Phi(v; \lambda) = \lambda \|Cv\|_1$, where $C$ is a transformation to a space in which $u$ is known to be sparse (Basis Pursuit in Compressed Sensing [14]);
4. a weighted sum of the coefficients in some basis of $U$ (such as a wavelet basis [23, 6] or singular vectors [18]).

These regularization methods are able to stably invert the operator $L$, at least approximately in the sense that $L\tilde{u} = L\tilde{L}^{-1}b \approx b$. Moreover, solutions of Eq. 6 are able to mitigate the effect of error within $b$; that is, even if $b$ is corrupted (e.g. by noise), $\tilde{u}$ will be a reasonable reconstruction. In contrast, a modest error in $L$ will likely result in a terrible reconstruction, regardless of the choice of $\Phi$. An example of this behavior is given in Fig. 1.

For the purposes of this paper, we assume that there exists a family $\{L_\theta\}_{\theta \in \Theta}$ of forward operators parameterized by $\theta \in \Theta$, and a unique $\hat{\theta} \in \Theta$ such that $L_{\hat{\theta}} = L$. Given a noisy measurement $b + \eta$, where $\eta$ is the noise, and a model parameter $\theta$, the approximate reconstruction of $u$, based on the regularization in Eq. 6 with operator $L_\theta$, is given by

$$\tilde{u}_{\theta, \eta} = \tilde{L}_\theta^{-1}(b + \eta). \tag{7}$$

(a) Ground truth, $u$.        (b) $u_\theta$ when $\theta = 2.3 = \hat{\theta}$     (c) $u_\theta$ when $\theta = 2.4 \not\approx \hat{\theta}$
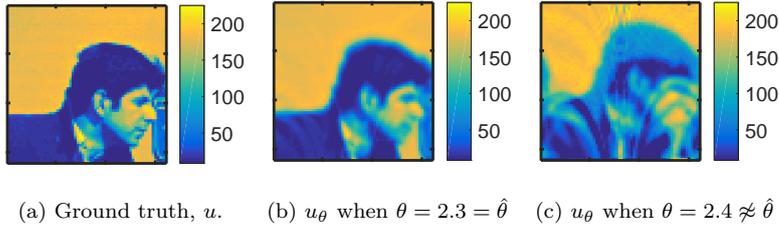
Fig. 1: Demonstration of the sensitivity in the reconstruction in Eq. 6 to errors in the forward operator. In this example $L = L_{\hat{\theta}}$ is the 'academic operator' from [29], $\theta$ is the parameter $R$ in [29, Table 1], and $\hat{\theta} = 2.3$. In this problem Tikhonov regularization was used to define the approximate inverse in Eq. 6.

where the tilde denotes the solution to a regularized problem of the form in Eq. 6 (where the choice of $\Phi$ is understood). This notation will be used throughout the remainder of the paper.

We define the residual as

$$r_{\theta,\eta} = (b + \eta) - L\tilde{u}_{\theta,\eta} = (I - L_\theta \tilde{L}_\theta^{-1})(b + \eta) \tag{8}$$

where $I$ is the identity operator. The residual is the main object that we study to determine when the parameter $\theta$ is poorly chosen.

3.2 Earth Mover's Distance

A key tool in our analysis of forward operator error is the Earth Mover's Distance. Below we summarize the presentation in [22].

**Definition 1 (Wasserstein Distance)** *Let $\Omega \subset \mathbb{R}^d$ be convex and compact, and let $c\colon \Omega \times \Omega \to [0, +\infty)$ be a distance. Given two non-negative distributions $\rho_1\colon \Omega \to \mathbb{R}^+, \rho_2\colon \Omega \to \mathbb{R}^+$ such that $\int_\Omega \rho_1 = \int_\Omega \rho_2$. For a given $p \in \mathbb{N}$ the $p$'th Wasserstein distance is*

$$W_p(\rho_1, \rho_2) = \left( \min_{\pi \geq 0} \int_{\Omega \times \Omega} c(x^{(1)}, x^{(2)})^p \pi(x^{(1)}, x^{(2)}) dx^{(1)} dx^{(2)} \right)^{1/p},$$

$$\text{subject to:} \quad \int_\Omega \pi(x^{(1)}, x^{(2)}) dx^{(2)} = \rho_1(x^{(1)}), \tag{9}$$

$$\int_\Omega \pi(x^{(1)}, x^{(2)}) dx^{(1)} = \rho_2(x^{(2)}).$$

The function $c$ is called the ground metric and each feasible function $\pi$ is referred to as a transport plan. In this work we set $c(x^{(1)}, x^{(2)}) = \left\| x^{(1)} - x^{(2)} \right\|_2$. The Earth Mover's Distance we define here is a special case of the Wasserstein distance where $p = 1$.

**Definition 2 (Earth Mover's Distance)** *Let $\Omega \subset \mathbb{R}^d$ be convex and compact, and let $c\colon \Omega \times \Omega \to [0, +\infty)$ be a distance. Given two non-negative distributions*

$\rho_1 \colon \Omega \to \mathbb{R}^+, \rho_2 \colon \Omega \to \mathbb{R}^+$ *such that* $\int_\Omega \rho_1 = \int_\Omega \rho_2$. *The Earth Mover's Distance (EMD) between* $\rho_1$ *and* $\rho_2$ *is*

$$\text{EMD}(\rho_1, \rho_2) = W_1(\rho_1, \rho_2). \tag{10}$$

The EMD can also be written in the equivalent form [10]

$$\begin{aligned}
\text{EMD}(\rho_1, \rho_2) = \min_m &\int_\Omega \|m(x)\|_2 \, dx, \\
\text{subject to:} \quad &\nabla \cdot m(x) + \rho_2(x) - \rho_1(x) = 0, \\
&m(x) \cdot n(x) = 0 \quad \forall x \in \partial\Omega,
\end{aligned} \tag{11}$$

where $n(x)$ is the normal vector at $x \in \partial\Omega$. This formulation is the basis for recently developed algorithms in [21,22].

## 4 Applying EMD to inverse problems

### 4.1 Residual and operator correctness

In a variational reconstruction procedure, the quality of the fit can be investigated by an analysis of $r_{\theta,\eta}$ and $\Phi(\tilde{u}_{\theta,\eta})$. Generally, the larger $\lambda$ the larger the first term and the smaller the second and vice-versa. Typically the value of $\lambda$ is chosen in an attempt to balance these contributions [17,18]. However if an incorrect forward operator is used, $r_{\theta,\eta}$ will have an additional contribution that does not depend on $\lambda$.

The characterization above can be made precise in the case of Tikhonov regularization by introducing a matrix notation and using Generalized Singular Value Decomposition [15, Chapter 8.7.3]. To this end, let $n = \dim(\mathcal{U}_{\Delta x})$ and $m = \dim(\mathcal{B}_{\Delta y})$, and expand $u$ and $b$ in terms of characteristic basis functions:

$$u(x) = \sum_{j=1}^n u_j \chi_{K_j^x}(x) \quad \text{and} \quad b(y) = \sum_{i=1}^m b_i \chi_{K_i^y}(y). \tag{12}$$

Then Eq. 5 becomes

$$\mathbf{L}\mathbf{u} = \mathbf{b}. \tag{13}$$

where $\mathbf{u} = (u_1, \dots, u_n)$, $\mathbf{b} = (b_1, \dots, b_m)$, and $\mathbf{L}$ has components

$$L_{i,j} = \frac{1}{\Delta y^{d_y}} \int_Y \chi_{K_i^y} L \chi_{K_j^x} dy. \tag{14}$$

**Definition 3 (GSVD)** *Let* $\mathbf{A} \in \mathbb{R}^{m \times n}$ *and* $\mathbf{B} \in \mathbb{R}^{o \times n}$ *be two matrices such that* $\text{null}(\mathbf{A}) \cap \text{null}(\mathbf{B}) = \emptyset$. *The Generalized Singular Value Decomposition (GSVD) of the matrix pair* $(\mathbf{A}, \mathbf{B})$ *is given by*

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{Z}^T \quad and \quad \mathbf{B} = \mathbf{V}\mathbf{\Gamma}\mathbf{Z}^T, \tag{15}$$

*where* $\mathbf{U} \in \mathbb{R}^{m \times n}$ *and* $\mathbf{V} \in \mathbb{R}^{o \times n}$ *are orthogonal;* $\mathbf{Z} \in \mathbb{R}^{n \times n}$ *is invertible; and*

$$\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_n) \in \mathbb{R}^{n \times n} \quad and \quad \mathbf{\Gamma} = \text{diag}(\gamma_1, \dots, \gamma_n) \in \mathbb{R}^{n \times n} \tag{16}$$

159 *are diagonal matrices such that*

$$1 \geq \sigma_1 \geq \cdots \geq \sigma_n \geq 0 \quad and \quad 0 \leq \gamma_1 \leq \cdots \leq \gamma_n \leq 1, \tag{17}$$

160 *with* $\mathbf{\Sigma}^2 + \mathbf{\Gamma}^2 = \mathbf{I}$.

161    Using the GSVD, we obtain the following:

162 **Proposition 1 (Residual with Tikhonov regularization)** *Suppose* $\mathbf{L}\mathbf{u} = \mathbf{b}$, *where*
163 $\mathbf{L} \in \mathbb{R}^{m \times n}$ *and* $m > n$. *Let* $\tilde{\mathbf{u}}_{\theta,\eta}$ *be defined by Eq. 7 with* $\Phi(\mathbf{v};\lambda) = \lambda \|\mathbf{C}\mathbf{v}\|_2^2$, *where*
164 $\mathbf{C} \in \mathbb{R}^{o \times n}$, *and a noise vector* $\boldsymbol{\eta} \in \mathbb{R}^m$ *whose elements are independent and spheri-*
165 *cally symmetric—that is,* $\boldsymbol{\eta}$ *and* $\mathbf{Q}\boldsymbol{\eta}$ *have the same probability distribution function*
166 *for any orthogonal matrix* $\mathbf{Q} \in \mathbb{R}^{m \times m}$. *Assume that* $\mathrm{null}(\mathbf{L}_\theta) \cap \mathrm{null}(\mathbf{C}) = \emptyset$ *so that*
167 *the GSVD*

$$\mathbf{L}_\theta = \mathbf{U}_\theta \mathbf{\Sigma}_\theta \mathbf{Z}_\theta^T \qquad \mathbf{C} = \mathbf{V}_\theta \mathbf{\Gamma}_\theta \mathbf{Z}_\theta^T \tag{18}$$

*for the matrix pair* $(\mathbf{L}_\theta, \mathbf{C})$ *is well-defined. Then the residual* $\mathbf{r}_{\theta,\eta}$ *associated to* $\tilde{\mathbf{u}}_{\theta,\eta}$
*satisfies the bound*

$$\|\mathbf{r}_{\theta,\eta}\|_2^2 \leq \left\|(\mathbf{I} - \mathbf{U}_\theta \mathbf{U}_\theta^T)\mathbf{b}\right\|_2^2 + \|(\mathbf{b} - \mathbf{L}_\theta \mathbf{u})\|_2^2$$
$$+ \frac{1}{4}\lambda \left\|\mathbf{Z}_\theta^T \mathbf{u}\right\|_2^2 + \frac{m - n + \mathrm{Tr}(\hat{\mathbf{D}}_{\theta,\lambda}^2))}{m} \mathbb{E}\left[\|\boldsymbol{\eta}\|_2^2\right]. \tag{19}$$

168 *The proof of Proposition 1 is in the appendix.*

169    This result shows how calibration error can induce $O(1)$ terms (with respect
170 to the regularization parameter $\lambda$) into the residual, the first two terms in Eq. 19.
171 The noise that is orthogonal to the image of $\mathbf{L}_\theta$ also induces $O(1)$ terms, even if
172 $\theta = \hat{\theta}$. Thus it is important to develop tools that can differentiate between these two
173 contributions. For completeness, one should also consider regularization with more
174 general forms of $\Phi$. Unfortunately in many situations, the operator $\tilde{\mathbf{L}}_\theta^{-1}$ is nonlinear,
175 and a rigorous analysis in this vein is much more difficult.

176 4.2 Introduction to the structure

177 We introduce a mathematical tool to detect contributions to $r_{\theta,\eta}$ that are due to
178 errors in the operator $L$, i.e., when $\theta \neq \hat{\theta}$, and is insensitive to noise in the residual.
179 This tool, which we call the structure, is a functional built using the Earth Mover's
180 Distance (EMD).

181 **Definition 4 (Structure)** *For any* $f \in L^1(\Omega)$, *the structure of* $f$ *is*

$$\mathrm{struc}\,[f] = \mathrm{EMD}(f^+, f^-), \tag{20}$$

182 *where*
$$f^+(x) = \max(f(x) - \mu, 0) \quad and \quad f^-(x) = \max(\mu - f(x), 0) \tag{21}$$
183 *and* $\mu = \frac{1}{\|\Omega\|} \int_\Omega f(x)dx$.

184 The following proposition is proven in the appendix.

**Proposition 2 (Basic Properties of Structure)** *The operator* struc [·] *satisfies the following properties:*

1. *it is a semi-norm on $L^1(\Omega)$;*
2. *for all $g \in L^1(\Omega)$ and $c \in \mathbb{R}$,*

$$\text{struc}\,[g] = \text{struc}\,[g + c]\,; \tag{22}$$

3. struc $[c] = 0$ *for any constant $c \in \mathbb{R}$;*
4. *if $\rho_1 \colon \Omega \to \mathbb{R}^+$, $\rho_2 \colon \Omega \to \mathbb{R}^+$ and $\int_\Omega \rho_1 = \int_\Omega \rho_2$,*

$$\text{struc}\,[\rho_2 - \rho_1] = \text{EMD}(\rho_1, \rho_2). \tag{23}$$

Using struc [·] is a good strategy for detecting operator error for several reasons:

- The struc [·] is small when applied to piecewise noise and large when applied to a (non-constant) smooth function. (Rigorous statements this effect are made in Section 4.3 below). Thus struc $[r_{\theta,\eta}]$ will be small when the forward operator is correct and large when it is not. Although the struc [·] of a constant is zero, any such contribution to the residual can be discerned by applying a standard norm to its spatial average.
- With recent algorithmic advances [21, 22], the underlying EMD calculation for computing struc [·] can be performed quickly. For example when $\mathbf{b} \in \mathbb{R}^{256} \times \mathbb{R}^{256}$, the structure calculation takes less than a second on consumer grade hardware.
- Because its evaluation does not affect the actual inverse procedure, the structure calculation can be incorporated into existing work flows without altering old code. Thus it can be quickly integrated into an existing toolbox for solving inverse problems.
- The struc $[r_{\theta,\eta}]$ calculation produces not only a number, but also outputs a transport plan (see Figs. 4b, 4d). For certain classes of forward operators this additional information can be leveraged to correct forward operators with minimal tuning. This idea will be explored in future work.

4.3 Theoretical Results

In this section we establish some theoretical results which support the use of the structure as a tool for diagnosing structural errors in the forward operator of an inverse problem. The proofs of Theorems 1–2 are given in Appendix. A.

**Theorem 1 (Characterization of noise by structure)** *Given non-negative integers integers $d$ and $\ell$, let $\Omega = [0,1)^d$ and let $\mathcal{O}_\ell = \left\{\omega_{\ell,1}, \ldots, \omega_{\ell,2^{\ell d}}\right\}$ partition $\Omega$ into $2^{\ell d}$ hypercubes of volume $2^{-\ell d}$. Define $h_\ell \colon \Omega \to \mathbb{R}$ by*

$$h_\ell(y) = \eta_{\ell,1}\chi_{\ell,1}(y) + \cdots + \eta_{\ell,2^{\ell d}}\chi_{\ell,2^{\ell d}}(y) \tag{24}$$

*where*

$$\chi_{\ell,i}(y) = \begin{cases} 1, & x \in \omega_{\ell,i}, \\ 0, & x \notin \omega_{\ell,i}, \end{cases} \tag{25}$$

and $\{\eta_{\ell,i}\}_{i=1}^{2^{\ell d}}$ is a set i.i.d. random variables with mean $\mu$ and variance $\sigma^2$ (See Fig. 2 for a visualization of $h_\ell$.) If $\epsilon_\ell = 2^{-\ell}$, then as $\ell \to \infty$, $\epsilon_\ell \to 0$ and

$$\mathbb{E}\left[\text{struc}\left[h_\ell\right]\right] \leq \sigma \begin{cases} -\epsilon_\ell \log \epsilon_\ell, & d = 2, \\ 2\sqrt{d}\epsilon_\ell, & d > 2, \end{cases} \tag{26}$$

where the expectation is with respect to the weights $\eta_{\ell,i}$.

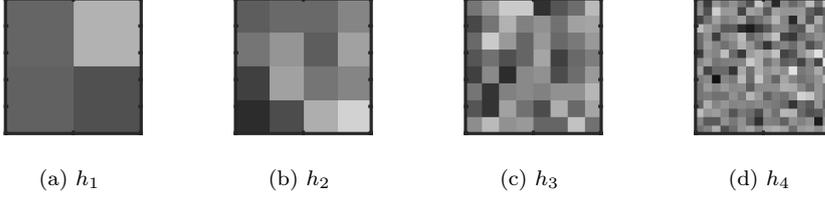| (a) $h_1$ | (b) $h_2$ | (c) $h_3$ | (d) $h_4$ |

Fig. 2: Example of $h_\ell$ when $d = 2$, $\mu = 0$, and $\sigma = 1$.

**Lemma 1 (L2 norm of Noise)** *Given the assumptions of Thm. 1, suppose further that $\mu = 0$. Then*

$$\mathbb{E}\left[\|h_\ell\|_2^2\right] = \sigma^2, \tag{27}$$

where the expectation is with respect to the weights $\eta_{\ell,i}$.

**Theorem 2 (Characterization of a smooth function by structure)** *Given the assumptions of Thm. 1, let $R_\ell \colon \mathcal{B} \to \mathcal{B}_{\epsilon_\ell}$. If*

$$R_\ell \phi(y) = \frac{1}{\omega_{\ell,i}} \int_{\omega_{\ell,i}} \phi(z)dz, \quad \forall y \in \omega_{\ell,i}. \tag{28}$$

*where $\phi \in C^1\left(\overline{Y}\right)$ then*

$$\left|\text{struc}\left[R_\ell \phi\right] - \text{struc}\left[\phi\right]\right| \leq C(|\nabla\phi|)\, d\epsilon_\ell^2, \tag{29}$$

*where the constant $C$ depends on the maximum of $\nabla\phi$ on $\overline{Y}$. In particular,*

$$\text{struc}\left[R_\ell \phi\right] \to \text{struc}\left[\phi\right] \quad \text{as } \ell \to +\infty. \tag{30}$$

4.4 Comparison with prior work

The work here is inspired, in part, by the study of seismic imaging inverse problems in [7,8,34]. There the authors measure the misfit between simulated and measured data using the Wasserstein distance squared $W_2^2(\rho_1, \rho_2) = (W_2(\rho_1, \rho_2))^2$. To handle the possibly negative distributions, the authors in [7,8,34] introduce the *misfit* function

$$\begin{aligned} d(f,g) = W_2^2 &\left(\frac{\max(f,0)}{\int \max(f,0)dx}, \frac{\max(g,0)}{\int \max(g,0)dx}\right) \\ &+ W_2^2\left(\frac{\max(-f,0)}{\int \max(-f,0)dx}, \frac{\max(-g,0)}{\int \max(-g,0)dx}\right) \end{aligned} \tag{31}$$

which plays a similar role to struc $[f - g]$ in this work. In [7, Section 2.6] the authors show that $d$ is insensitive to noise, with a scaling result that is similar to Thm. 1 up to a logarithmic factor. Specifically, if $f$ and $g$ are two non-negative functions such that $f - g$ has the form of $h_\ell$, defined in Eq. 24), with uniformly distributed noise, then

$$d(f, g) = O(\epsilon_\ell). \tag{32}$$

The approach taken in [7,8,34] differs from the approach in this paper in at least two key ways. First is the choice of $W_2^2$ rather than $W_1$. This has the following consequences:

- $W_2$ and $W_2^2$ have the property of *cyclic monotonicity* (see [9, Sec. 2.1] for a definition and proof), which can be used to show convexity of $d$ with respect to shifts, dilation and partial amplitude loss. In this work we make no such claims about the convexity of struc $[\cdot]$.
- As a semi-norm, the EMD (like all $W_p$ for $p \in [1, \infty)$) is a degree-one homogeneous functional and satisfies a triangle inequality (see [32, p. 94]. The functional $W_2^2$ has neither property. For example of the latter, let $f = 2\chi_{0,1/2}$, $h = 2\chi_{1/2,1}$ and $g = 2\chi_{1,3/2}$. Then $W_2^2(f, h) = \frac{1}{4}$, $W_2^2(h, g) = \frac{1}{4}$ but $W_2^2(f, g) = 1$, then

$$W_2^2(f, g) > W_2^2(f, h) + W_2^2(h, g). \tag{33}$$

- Redefining $d$ with $W_2$ instead of $W_2^2$ would recover a triangle inequality and degree-one homogeneity. However, the cost of such a modification would be to increase the sensitivity of $d$ to noise. Indeed, the scaling in Eq. 32 would change from $O(\epsilon_\ell)$ to $O(\epsilon_\ell^{1/2})$, which is significantly slower than the scaling in Thm. 1.
- Finally, $W_1$ is more directly analogous to the definition of work used throughout physics, distance times effort. Consider the case when

$$f(x) = \frac{1}{2}\chi_{[0,2]}(x) \quad g(x) = \frac{1}{2}\chi_{[1,3]}(x) \tag{34}$$

and the two transport plans

$$\pi_1(x_1, x_2) = \begin{cases} 1/2 \text{ if } x_2 = 1 + x_1 \text{ and } x_1 \in [0, 2] \\ 0 \text{ otherwise} \end{cases} \tag{35}$$

$$\pi_2(x_1, x_2) = \begin{cases} 1/2 \text{ if } x_2 = 2 + x_1 \text{ and } x_1 \in [0, 1] \\ 0 \text{ otherwise} \end{cases} \tag{36}$$

The cost of $\pi_1$ as measured by $W_2$ is twice that of $\pi_2$. Both plans cost the same as measured by $W_1$. In words $W_2$ 'prefers' to make many smaller movements as opposed to fewer larger movements, while $W_1$ is agnostic to such differences.

The second key difference between the approach in [7,8,34] and the approach taken here lies in the definition of $d$ and struc $[\cdot]$, both of which are used to address the fact that the Wasserstein metric is only defined for non-negative distributions with the same mass. It is worth noting that $d(f, g)$ and struc $[\cdot]$ could be defined using any Wassterstein metric. However, $d$ introduces several undesirable artifacts.

– The normalization in the definition means that

$$d(\lambda f, \lambda g) = d(f, g), \quad \forall \lambda \in \mathbb{R}^+. \tag{37}$$

In particular, unlike struc $[\cdot]$, it is not degree-one homogeneous.

– Special care is required in the case that $\max(f, 0) \equiv 0$ but $\max(g, 0) \not\equiv 0$. Indeed one of the reasons that the results in Eq. 32 require $f$ and $g$ to be positive and differ only by uniform noise is that small changes is the noise can alter the support of $\max(f, 0)$ and $\max(g, 0)$. The struc $[\cdot]$ has no such restrictions on the noise model.

– The struc $[\cdot]$ is continuous w.r.t. the $L_1(\Omega)$ norm provided that $\Omega$ is bounded (see Lemma 5). $d(f, g)$, however, is not. For example consider, the functions

$$f_\epsilon = \chi_{[\epsilon, 1-\epsilon]} - \epsilon \chi_{(1-\epsilon, 1]}, \quad g_\epsilon = -\epsilon \chi_{[0, \epsilon)} + \chi_{[\epsilon, 1-\epsilon]} - \epsilon \chi_{(1-\epsilon, 1]}. \tag{38}$$

Clearly $f_\epsilon - g_\epsilon \to 0$ in $L_1(\Omega)$ as $\epsilon \to 0$; however,

$$\lim_{\epsilon \to 0} d(f_\epsilon, g_\epsilon) \geq \lim_{\epsilon \to 0} \frac{1}{2}\left(1 + \frac{\epsilon}{4}\right)^2 = \frac{1}{2}. \tag{39}$$

This lack of continuity due to sign changes is one of the reasons for having restrictions on the noise model for $d(f, g)$.

– The kernel of struc $[\cdot]$ consists of constant functions, and so struc $[f - g] = 0 \iff f = g + c$ for some constant $c$. This $c$ is easily recovered by computing the difference between the averages if $f$ and $g$. On the other hand, the kernel of $d$ is

$$\mathrm{Ker}(d) = \left\{ \begin{array}{l} (f, g) \in L^1 \times L^1 : \max(f, 0) = \lambda_+ \max(g, 0) \text{ and} \\ \qquad\qquad \max(-f, 0) = \lambda_- \max(-g, 0) \quad \text{for } \lambda_+, \lambda_- \in \mathbb{R}^+ \end{array} \right\} \tag{40}$$

It is more difficult to account for such a kernel.

## 5 Numerical Results

In this section we present the results of several numerical experiments. We make two simplying assumptions. First, we let $X$ and $Y$ be two dimensional domains. This choice is motivated by ease of visualization as well as the availability of code to quickly compute the EMD in two dimensions. We, however, believe that our results generalize well to high dimensional problems. Second, we assume that $L_\theta$ is linear in $\theta$. This choice is for simplicity, but it also is a reasonable approximation for finding a local optimum. Indeed, if $L_\theta$ smoothly depends on $\theta$, then $L$ is locally linear:

$$L_{\hat{\theta} + \delta\theta} = L_{\hat{\theta}} + \nabla_\theta L(\hat{\theta}) \cdot \delta\theta + O(\delta\theta^2). \tag{41}$$

For each experiment, we provide with a known signal $u$ and a family of operators $\{L_\theta\}_{\theta \in \Theta}$. We then set $L = L_{\hat{\theta}}$ for some $\hat{\theta} \in \Theta$, generate a measurement $b = L_{\hat{\theta}} u$, and examine the behavior of struc $[r_{\theta, \eta}]$ as a function of $\theta$. The expectation is that

$$\hat{\theta} \approx \theta^* := \underset{\theta \in \Theta}{\mathrm{argmin}}\, \mathrm{struc}\, [r_{\theta, \eta}]. \tag{42}$$

| Parameter | Value | Parameter | Value | Ref. | Parameter | Value | Ref. |
|---|---|---|---|---|---|---|---|
| **Discretization**[1] | | **Inversion** | | | struc$[\cdot]$ | | |
| $\Delta x$ | 1/64 | $\Phi(\cdot, \lambda)$ | $\lambda\,\mathrm{TV}(u)$ | [27] | Max Iter | 8000 | [22] |
| $\Delta y$ | 1/100 | $\lambda$ | 10 | [27] | $\mathrm{EMD}_\mu$ | 7e-6 | [22] |
| | | $\mu$ | 100 | [14] | $\mathrm{EMD}_\tau$ | 3 | [22] |
| | | Bregman Iterations | 10 | [14] | | | |

Table 1: Numerical parameters for Experiments 1 - 3.

The first two experiments show that indeed $\theta^* \approx \hat{\theta}$ even with relatively high noise. The final experiment illustrates that the method performs better as the problem becomes more overdetermined. We report a figure of merit, the contrast, defined as:

$$\mathrm{cont}(F) = \frac{\max(F) - \min(F)}{\max(F) + \min(F)} \tag{43}$$

for any $F\colon \Theta \to \mathbb{R}^+$ that is not identically zero. The contrast measures the depth of a minimum, and the greater the contrast, the less the location of the minimum changes in the presence of additive noise in $F$. In all three experiments we compare the contrast of struc$[\cdot]$ with the discrete norms $\|\cdot\|_1$ and $\|\cdot\|_2$. For any $z \in \mathcal{B}_{\Delta y}$ these norms are given by,

$$\|z\|_1 = \Delta y^2 \sum_{i_1, i_2} |z_{i_1, i_2}| \quad \text{and} \quad \|z\|_2 = \Delta y \left( \sum_{i_1, i_2} z_{i_1, i_2}^2 \right)^{1/2} \tag{44}$$

We also generate plots of all three (semi-) norms as a function of the parameter $\theta$.

## 5.1 Implementation Details

The implementation of each of these experiments involves four basic steps: (i) the generation of the random forward operators $L_\theta$; (ii) generation of the signal $u$, measurement $b$ and noise $\eta$; (iii) calculation of $\tilde{u}_{\theta,\eta}$; and (iv) computation of the struc$[\cdot]$. The specific values of parameters needed to recreate our results are given in Table 1.

1. **Generation of the random forward operators.** Recall the definitions in Section 3.1. A forward operator $L_\theta$, even an academic one, but rather a the discretization of an operator $\mathcal{L}\colon \mathcal{U} \to \mathcal{B}$. In applications, $L_\theta$ models the action of some physical process which produces a measurement. For example in seismic imaging the forward operator is the propagation of a seismic wave [7], and in plasma imaging in tokamaks the forward operator couples the optics of the camera with the symmetries of the plasma [33].

   For our experiments, we presume that $\mathcal{L}$ is a Line Integral Operator (LIO). (See Appendix B for details.) If $f\colon X \to \mathbb{R}$ and $g\colon Y \to \mathbb{R}$, then for each $y \in Y$, $g(y)$ represents the integral of $f$ over some path $p(y)$. Some examples of common LIO are the Radon, Abel and Helical Abel transforms [29].

---

[1]  $\Delta x$ and $\Delta y$ both change for Experiment 3, however the other parameters are fixed.

291  2. **Generation of the signal, measurement and noise.** The underlying signal
292      $u \in \mathcal{U}_{\Delta x}$ is a series of concentric rings (see Fig. 3a). Then we apply $L_{\hat{\theta}}$ to $u$
293      to obtain a noiseless measurement $b \in \mathcal{B}_{\Delta y}$ (see Fig. 3b). The noisy signal (see
294      Fig. 3c) is generated by adding independent white noise $\eta$ with mean zero and
295      variance $\sigma$ to each element of $b$ so that

$$\text{SNR} = \frac{\|b\|_2}{\|\eta\|_2} \tag{45}$$

296      is at a specified level.



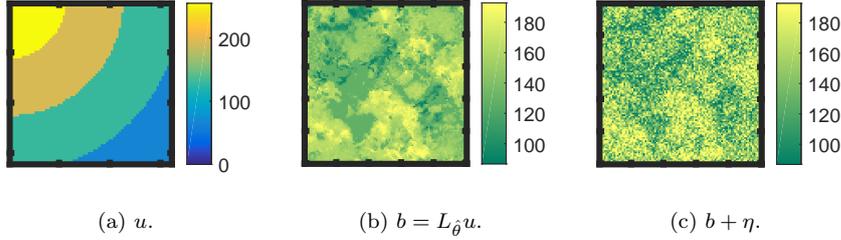(a) $u$.                              (b) $b = L_{\hat{\theta}}u$.                              (c) $b + \eta$.

Fig. 3: The signal $u$, measurement $b$, and noisy measure $b + \eta$ for Experiment 1.

3. **Computation of $\tilde{u}_{\theta,\eta}$.** Throughout these experiments, we use the inversion
   procedure of the form of Eq. 6 with $\Phi(v; \lambda) = \lambda \|\mathbf{C}v\|_1$ where $\mathbf{C}$ is a one-sided
   discrete approximation of the gradient operator:

$$(Cv)_{2i,j} = \frac{1}{dx}\left(v_{i,j} - v_{\ell-1,j}\right)$$
$$(Cv)_{2i+1,j} = \frac{1}{dy}\left(v_{i,j} - v_{i,j-1}\right) \tag{46}$$

297      where $v_{i,j}$ is the $i$'th x and $j$'th y component of the vector $\mathbf{v}$, and likewise
298      for $(\mathbf{C}v)_{i,j}$. This is TV regularization and has found wide success within image
299      processing, especially when the underlying signal to be recovered is piecewise
300      constant [14, 27].
301      To solve the resulting non-linear variational problem, we use the Split-Bregman
302      algorithm, specifically the Generalized Split-Bregman Algorithm (GSBA) of [14],
303      which requires specification of a step size parameter $\mu$ (called $\lambda$ in [14]). GSBA
304      requires the repeated solution of the linaer system $(\mathbf{L}^T\mathbf{L} + \lambda^2\mathbf{C}^T\mathbf{C})x = y$. The
305      matrix $(\mathbf{L}^T\mathbf{L} + \lambda^2\mathbf{C}^T\mathbf{C})$ is sparse and so we solve it using the L-BFGS [2, 35]
306      method (limited memory Broyden-Fletcher-Goldfarb-Shanno[3, 11, 13, 31]).
307  4. **Computation of the** struc $[\cdot]$**.** Computing struc $[\cdot]$ requires computing EMD.
308      The algorithm that we use is given in [21, 22, 28].

309  5.2 Experiment 1

310  This experiment is based on a normalized Eq. 41 where $p = 1$. Let $L_0$ and $L_1$ be two
311  operators generated as described in Appendix B. We define $\theta \in [0, 1]$ and

$$L_\theta = (1 - \theta)L_0 + \theta L_1. \tag{47}$$

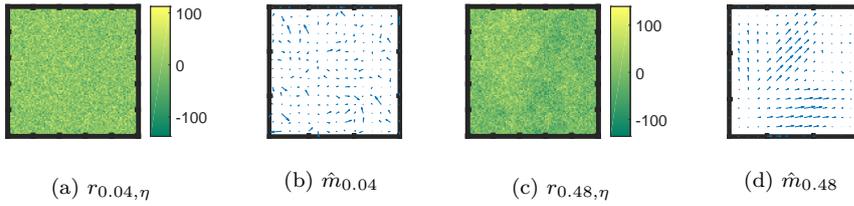(a) $r_{0.04,\eta}$    (b) $\hat{m}_{0.04}$    (c) $r_{0.48,\eta}$    (d) $\hat{m}_{0.48}$

Fig. 4: Results from Experiment 1. The residual and flow $\hat{m}_\theta$ that minimizes Eq. 11 for a given $\theta$. In Figs. 4b and 4d, the orientation of the arrows represents the direction $\hat{m}_\theta$, and the length of the arrows is proportional to the magnitude.

Fig. 4 is a plot of the residual for different values of $\theta$. In Fig. 4a, $\theta = 0.04$, and in Fig. 4c $\theta = 0.48$. Upon close inspection, one can see that from Fig. 4a that when $\theta$ is small the residual visually looks like white noise, whereas from Fig. 4c when $\theta$ is large the residual has underlying structure in addition to the noise. It is, however, difficult to see. Despite these two plots appearing similar they have very different structures, struc$[r_{0.04,\eta}] \approx 0.06$ and struc$[r_{0.48,\eta}] \approx 0.54$. The structure is also evident by looking at Figs. 4b, 4d, which are $m$ from Eq. 11. Note that when $\theta = 0.04$, $m$ is higgledy-piggledy, whereas when $\theta = 0.48$, $m$ appears much more orderly.

A plot of struc$[r_{\theta,\eta}]$ vs $\theta$ is given in Fig. 5. Clearly, struc$[r_{\theta,\eta}]$ is minimized when $\theta \approx 0$. Further, we note that struc$[r_{\theta,\eta}]$ is increasing as a function of $\theta$ when $\theta \in [0, 0.5]$, however then decreases. This is expected behavior around the minimum, however the problem is evidently not convex away from $\hat{\theta}$. This is important to keep in mind for future work.



(a) struc$\left[r_{\theta,\eta}\right]$ vs $\theta$.    (b) $\left\|r_{\theta,\eta}\right\|_1$ vs $\theta$.    (c) $\left\|r_{\theta,\eta}\right\|_2$ vs $\theta$.

Fig. 5: Results from Experiment 1. The value of $r_{\theta,\eta}$ as measured by struc$[\cdot]$, $\|\cdot\|_1$ and $\|\cdot\|_2$. In all examples the minimum occurs when $\theta = 0$ however the contrast is greatest for struc$[\cdot]$.

## 5.3 Experiment 2

Experiment 2 is also based on a normalized Eq. 41, however in this case $p = 2$ and $\hat{\theta} = \left(\frac{1}{2}, \frac{1}{2}\right)$. The true signal used in Experiment 2 is the same as in Experiment 1

(see Fig. 3a). This experiment studies the change in the contrast for $\operatorname{struc}\left[\cdot\right], \left\|\cdot\right\|_1$ and $\left\|\cdot\right\|_2$ as the SNR decreases. The results are summarized in Table 2.



(a) $\operatorname{struc}\left[r_{\theta,\eta}\right]$ vs $\theta$      (b) $\left\|r_{\theta,\eta}\right\|_1$ vs $\theta$      (c) $\left\|r_{\theta,\eta}\right\|_2$ vs $\theta$

Fig. 6: Results from Experiment 2. In these plots SNR $= 25$.



(a) $\operatorname{struc}\left[r_{\theta,\eta}\right]$ vs $\theta$      (b) $\left\|r_{\theta,\eta}\right\|_1$ vs $\theta$      (c) $\left\|r_{\theta,\eta}\right\|_2$ vs $\theta$

Fig. 7: Results from Experiment 2. In these plots SNR $= 5$.

| Contrast | $\operatorname{struc}\left[r_{\theta,\eta}\right]$ | $\left\|r_{\theta,\eta}\right\|_1$ | $\left\|r_{\theta,\eta}\right\|_2$ |
|---|---|---|---|
| SNR $= 25$ | 0.7547 | 0.3493 | 0.3544 |
| SNR $= 5$ | 0.5917 | 0.0398 | 0.0404 |

Table 2: Results from Experiment 2. The contrast for different choices of (semi)norms. Larger is better.

In all cases, the contrast of $\operatorname{struc}\left[\cdot\right]$ is greatest, and the contrast of $\operatorname{struc}\left[\cdot\right]$ relative to $\left\|\cdot\right\|_1$ of $\left\|\cdot\right\|_2$ increases as the problem becomes more noisy. This suggests that $\operatorname{struc}\left[\cdot\right]$ is a more robust choice of semi-norm for measuring the level of miscalibration of $L_\theta$, especially when noise levels are high.

## 5.4 Experiment 3

The final experiment examines the necessity of the overdetermined assumption of $L_\theta$. We repeat the setup of Experiment 2; however we fix the SNR $= 25$ and instead

adjust $\Delta y$ so that $L_\theta \colon \mathcal{U}_{\Delta x} \to \mathcal{B}_{\Delta y}$ becomes a square operator. We start with a fixed reference $\Delta y_0$, and consider

$$\mathcal{B}_{\Delta y_0} \cong \mathbb{R}^{100 \times 100} \quad \mathcal{B}_{4/3\Delta y_0} \cong \mathbb{R}^{75 \times 75} \quad \mathcal{B}_{2\Delta y_0} \cong \mathbb{R}^{50 \times 50} \quad \mathcal{B}_{4\Delta y_0} \cong \mathbb{R}^{25 \times 25}. \quad (48)$$

In all cases, $\mathcal{U}_{\Delta x} \cong \mathbb{R}^{25 \times 25}$ is fixed. Each of the $\mathcal{B}$ in Eq. 48 are plotted in Fig. 8. The values of

$$\theta^s = \underset{\theta \in \Theta}{\mathrm{argmin}}\, \mathrm{struc}\left[r_{\theta,\eta}\right] \quad \theta^1 = \underset{\theta \in \Theta}{\mathrm{argmin}}\, \|r_{\theta,\eta}\|_1 \quad \theta^2 = \underset{\theta \in \Theta}{\mathrm{argmin}}\, \|r_{\theta,\eta}\|_2 \quad (49)$$

as well as the contrast are recorded in Table 3. Finally, plots of $\mathrm{struc}\left[r_{\theta,\eta}\right]$, $\|r_{\theta,\eta}\|_1$, and $\|r_{\theta,\eta}\|_2$ vs $\theta$ as $\Delta y$ increases are Figs. 9, 10, 11, and 12.
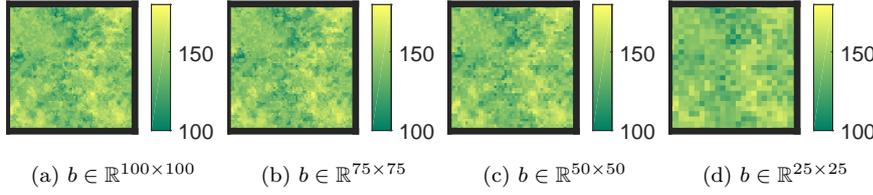


(a) $b \in \mathbb{R}^{100 \times 100}$      (b) $b \in \mathbb{R}^{75 \times 75}$      (c) $b \in \mathbb{R}^{50 \times 50}$      (d) $b \in \mathbb{R}^{25 \times 25}$

Fig. 8: Results from Experiment 3. Plot of $b$ for various choices of $\Delta y$ (see Eq. 48).



(a) $\mathrm{struc}\left[r_{\theta,\eta}\right]$ vs $\theta$      (b) $\left\|r_{\theta,\eta}\right\|_1$ vs $\theta$      (c) $\left\|r_{\theta,\eta}\right\|_2$ vs $\theta$
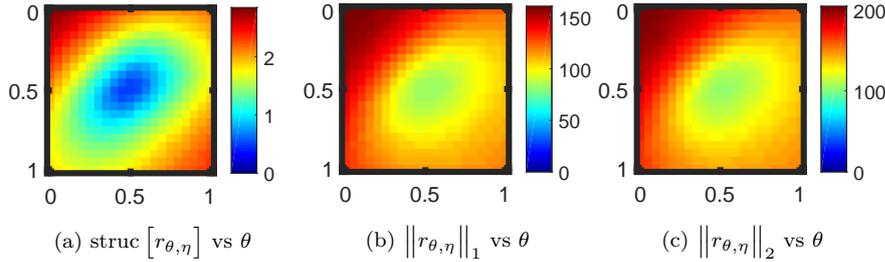
Fig. 9: Results from Experiment 3. In these plots $L \colon \mathbb{R}^{25 \times 25} \to \mathbb{R}^{100 \times 100}$. See Table 3 for the contrast, $\theta^s$, $\theta^1$, and $\theta^2$.

Throughout all trials of this experiment, $\theta^s$ was closer to $\hat\theta$ then $\theta^1$ or $\theta^2$. Additionally, the contrast is highest when the $\mathrm{struc}\left[\cdot\right]$ is used, except when $\mathcal{B}_{4\Delta y_0} \cong \mathbb{R}^{25 \times 25}$. These results show also that the degree to which the problem is overdetermined is indeed important. The more overdetermined the problem, the more nearly $\mathrm{struc}\left[r_{\theta,\eta}\right]$ is minimized at $\hat\theta$. Further, the more overdetermined the problem the greater the contrast of $\mathrm{struc}\left[\cdot\right]$ relative to $\|\cdot\|_1$ or $\|\cdot\|_2$. When $\mathcal{B}_{\Delta y_0} \cong \mathbb{R}^{100 \times 100}$, $\mathrm{cont}(\mathrm{struc}\left[r_{\theta,\eta}\right])$ is more than twice either $\mathrm{cont}(\|r_{\theta,\eta}\|_1)$ or $\mathrm{cont}(\|r_{\theta,\eta}\|_2)$. The ratio of $\mathrm{cont}(\mathrm{struc}\left[r_{\theta,\eta}\right])$ to either $\mathrm{cont}(\|r_{\theta,\eta}\|_1)$ or $\mathrm{cont}(\|r_{\theta,\eta}\|_2)$ decreases as $L$ becomes square, until finally $\mathcal{B}_{4\Delta y_0} \cong \mathbb{R}^{25 \times 25}$ and all three contrasts are similar. These results are consistent with Thms. 1 - 2, which together suggest that as $\Delta y$ decreases, the ability of $\mathrm{struc}\left[\cdot\right]$ to distinguish between noise and structure increases.

(a) struc $\left[r_{\theta,\eta}\right]$ vs $\theta$       (b) $\left\|r_{\theta,\eta}\right\|_1$ vs $\theta$       (c) $\left\|r_{\theta,\eta}\right\|_2$ vs $\theta$

Fig. 10: Results from Experiment 3. In these plots $L\colon \mathbb{R}^{25\times 25} \to \mathbb{R}^{75\times 75}$. See Table 3 for the contrast, $\theta^s$, $\theta^1$, and $\theta^2$.



(a) struc $\left[r_{\theta,\eta}\right]$ vs $\theta$       (b) $\left\|r_{\theta,\eta}\right\|_1$ vs $\theta$       (c) $\left\|r_{\theta,\eta}\right\|_2$ vs $\theta$

Fig. 11: Results from Experiment 3. In these plots $L\colon \mathbb{R}^{25\times 25} \to \mathbb{R}^{50\times 50}$. See Table 3 for the contrast, $\theta^s$, $\theta^1$, and $\theta^2$.



(a) struc $\left[r_{\theta,\eta}\right]$ vs $\theta$       (b) $\left\|r_{\theta,\eta}\right\|_1$ vs $\theta$       (c) $\left\|r_{\theta,\eta}\right\|_2$ vs $\theta$
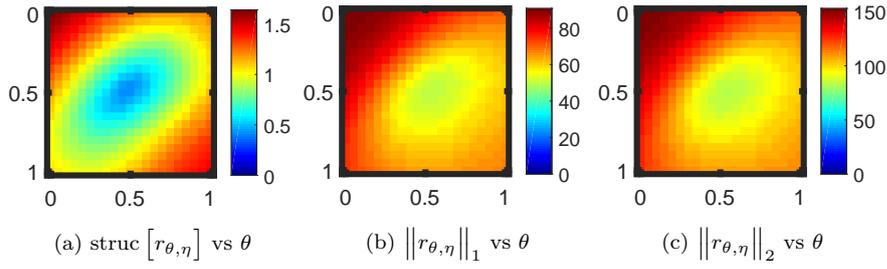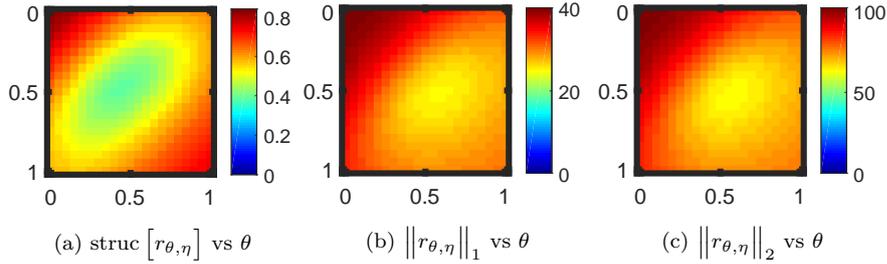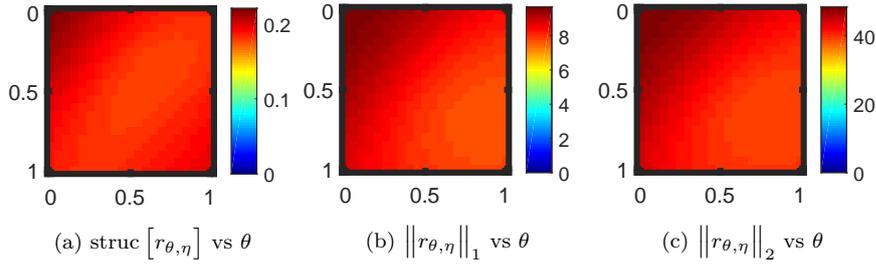
Fig. 12: Results from Experiment 3. In these plots $L\colon \mathbb{R}^{25\times 25} \to \mathbb{R}^{25\times 25}$. See Table 3 for the contrast, $\theta^s$, $\theta^1$, and $\theta^2$.

## 6 Conclusion

In this work we have developed a new functional called the structure, which is suitable for detecting forward operator error as it arises in inverse problems. The structure is defined by use of the Earth Mover's Distance (EMD), using a very rapid algorithm and a homogeneous degree one distance. The structure takes as input the residual from an existing inverse procedure, and can be computed quickly. We prove some apparently new results concerning the treatment of noise by EMD. Further, we consistent with these theoretical results we perform numerical experiments and show

|  | $\theta^s$ | $\theta^1$ | $\theta^2$ |
|---|---|---|---|
| $\mathbf{b} \in \mathbb{R}^{100 \times 100}$ | (0.55,0.55) | (0.55,0.55) | (0.55,0.55) |
| $\mathbf{b} \in \mathbb{R}^{75 \times 75}$ | (0.55,0.55) | (0.55,0.60) | (0.55,0.60) |
| $\mathbf{b} \in \mathbb{R}^{50 \times 50}$ | (0.55,0.50) | (0.55,0.65) | (0.60,0.60) |
| $\mathbf{b} \in \mathbb{R}^{25 \times 25}$ | (0.45,0.70) | (0.75,0.90) | (0.70,0.90) |
| Contrast | struc $\left[ r_{\theta,\eta} \right]$ | $\left\| r_{\theta,\eta} \right\|_1$ | $\left\| r_{\theta,\eta} \right\|_2$ |
| $\mathbf{b} \in \mathbb{R}^{100 \times 100}$ | 0.7044 | 0.3155 | 0.3215 |
| $\mathbf{b} \in \mathbb{R}^{75 \times 75}$ | 0.5877 | 0.2876 | 0.2931 |
| $\mathbf{b} \in \mathbb{R}^{50 \times 50}$ | 0.3677 | 0.2337 | 0.2376 |
| $\mathbf{b} \in \mathbb{R}^{25 \times 25}$ | 0.1116 | 0.1198 | 0.1125 |

Table 3: Results from Experiment 3. The above two tables record the location of the minimizer and contrast. Closer to $(0.5, 0.5)$ is better for $\theta$, and the larger the contrast the better.

that the structure is able to distinguish between error in the modeling of a forward operator, and noise in the signal of an inverse problem.

Our numerical results concern a model linear forward operator. On these problems the structure of the residual is indeed minimized when the correct forward operator is used and. The $L_1$ or $L_2$ norms of the residual are also minimized around the correct forward operator, the structure, however, is more localized and has better contrast around the minimum. Further, we observe that the degree to which the inverse problem is overdetermined is pivotal to the success of our procedure. The more over determined the problem, the more useful the structure. This is borne out by the analysis in the case of linear regularization, as well as the numerical results on more sophisticated problem.

In the future, we will extend our work to more sophisticated non-linear operators and promote our error detecting method into an error correcting method.

## A Proofs

*Proof (Proof of Proposition 1)* Given $\Phi(\mathbf{v}; \lambda) = \lambda \left\| \mathbf{Cv} \right\|_2^2$, the normal equations for Eq. 6 are

$$(\mathbf{L}_\theta^T \mathbf{L}_\theta + \lambda \mathbf{C}^T \mathbf{C}) \tilde{\mathbf{u}}_{\theta,\eta} = \mathbf{L}_\theta^T (\mathbf{b} + \boldsymbol{\eta}). \tag{50}$$

Therefore $\tilde{\mathbf{L}}_\theta^{-1} = (\mathbf{L}_\theta^T \mathbf{L}_\theta + \lambda \mathbf{C}^T \mathbf{C}^T)^{-1} \mathbf{L}_\theta^T$. Using the GSVD in Eq. Eq. 18, a direct calculation gives

$$\mathbf{L}_\theta \tilde{\mathbf{L}}_\theta^{-1} = \mathbf{U}_\theta \mathbf{D}_{\theta,\lambda} \mathbf{U}_\theta^T, \quad \text{where } \mathbf{D}_{\theta,\lambda} := \frac{\boldsymbol{\Sigma}_\theta^2}{\boldsymbol{\Sigma}_\theta^2 + \lambda \boldsymbol{\Gamma}_\theta^2} \in \mathbb{R}^{n \times n}. \tag{51}$$

Thus according to the definition of the residual in Eq. 8,

$$\mathbf{r}_{\theta,\eta} = (\mathbf{I} - \mathbf{L} \tilde{\mathbf{L}}^{-1})(\mathbf{b} + \boldsymbol{\eta}) = \mathbf{U}_\theta \hat{\mathbf{D}}_{\theta,\lambda} \mathbf{U}_\theta^T (\mathbf{b} + \boldsymbol{\eta}) + (\mathbf{I} - \mathbf{U}_\theta \mathbf{U}_\theta^T)(\mathbf{b} + \boldsymbol{\eta}) \tag{52}$$

where

$$\hat{\mathbf{D}}_{\theta,\lambda} := (\mathbf{I} - \mathbf{D}_{\theta,\lambda}) = \frac{\lambda \boldsymbol{\Gamma}_\theta^2}{\boldsymbol{\Sigma}_\theta^2 + \lambda \boldsymbol{\Gamma}_\theta^2} > 0. \tag{53}$$

We first bound two of the deterministic components of the residual. Using the GSVD,

$$\mathbf{U}_\theta \hat{\mathbf{D}}_{\theta,\lambda} \mathbf{U}_\theta^T \mathbf{b} = \mathbf{U}_\theta \hat{\mathbf{D}}_{\theta,\lambda} \mathbf{U}_\theta^T \mathbf{L}_\theta \mathbf{u} + \mathbf{U}_\theta \hat{\mathbf{D}}_{\theta,\lambda} \mathbf{U}_\theta^T (\mathbf{b} - \mathbf{L}_\theta \mathbf{u})$$
$$= \mathbf{U}_\theta \hat{\mathbf{D}}_{\theta,\lambda} \boldsymbol{\Sigma}_\theta \mathbf{Z}_\theta^T \mathbf{u} + \mathbf{U}_\theta \hat{\mathbf{D}}_{\theta,\lambda} \mathbf{U}_\theta^T (\mathbf{b} - \mathbf{L}_\theta \mathbf{u}). \tag{54}$$

375 Since $\left\|\hat{\mathbf{D}}_{\theta,\lambda}\right\|_2 \leq 1$ and $\mathbf{U}_\theta$ is orthogonal, it follows that

$$\left\|\mathbf{U}_\theta \hat{\mathbf{D}}_{\theta,\lambda} \mathbf{U}_\theta^T (\mathbf{b} - \mathbf{L}_\theta \mathbf{u})\right\|_2^2 \leq \|(\mathbf{b} - \mathbf{L}_\theta \mathbf{u})\|_2^2 \tag{55}$$

376 Furthermore, since

$$\hat{\mathbf{D}}_{\theta,\lambda} \boldsymbol{\Sigma}_\theta = \frac{\lambda \boldsymbol{\Gamma}_\theta^2 \boldsymbol{\Sigma}_\theta}{\boldsymbol{\Sigma}_\theta^2 + \lambda \boldsymbol{\Gamma}_\theta^2} \leq \frac{1}{2}\sqrt{\lambda}\boldsymbol{\Gamma}_\theta \leq \frac{1}{2}\sqrt{\lambda}\mathbf{I} \tag{56}$$

377 (where the inequalities between the diagonal matrices above are interpreted element-wise), it
378 follows that

$$\left\|\mathbf{U}_\theta \hat{\mathbf{D}}_{\theta,\lambda} \boldsymbol{\Sigma}_\theta \mathbf{Z}_\theta^T \mathbf{u}\right\|_2^2 \leq \left\|\hat{\mathbf{D}}_{\theta,\lambda} \boldsymbol{\Sigma}_\theta\right\|_2^2 \left\|\mathbf{Z}_\theta^T \mathbf{u}\right\|_2^2 \leq \frac{1}{4}\lambda \left\|\mathbf{Z}_\theta^T \mathbf{u}\right\|_2^2. \tag{57}$$

379 We next bound the noise component of the residual. Let $\mathbf{W}_\theta \in \mathbb{R}^{m \times (m-n)}$ be a matrix such
380 that $\mathbf{Q} := (\mathbf{U}_\theta | \mathbf{W}_\theta) \in \mathbb{R}^{m \times m}$ is orthogonal and set

$$\boldsymbol{\alpha} = \begin{pmatrix} \boldsymbol{\alpha}_\parallel \\ \boldsymbol{\alpha}_\perp \end{pmatrix} := \mathbf{Q}^T \boldsymbol{\eta} = \begin{pmatrix} \mathbf{U}_\theta^T \boldsymbol{\eta} \\ \mathbf{W}_\theta^T \boldsymbol{\eta} \end{pmatrix}. \tag{58}$$

381 Then

$$\left\|(\mathbf{I} - \mathbf{L}\tilde{\mathbf{L}}^{-1})\boldsymbol{\eta}\right\|_2^2 = \left\|\mathbf{U}_\theta \hat{\mathbf{D}}_{\theta,\lambda} \mathbf{U}_\theta^T \boldsymbol{\eta} + (\mathbf{I} - \mathbf{U}_\theta \mathbf{U}_\theta^T)\boldsymbol{\eta}\right\|_2^2 = \left\|\mathbf{U}_\theta \hat{\mathbf{D}}_{\theta,\lambda} \boldsymbol{\alpha}_\parallel\right\|_2^2 + \|\mathbf{W}_\theta \boldsymbol{\alpha}_\perp\|_2^2, \tag{59}$$

where the last equality uses the fact that the columns of $\mathbf{U}_\theta$ and $\mathbf{W}_\theta$ are orthogonal and $\mathbf{I} - \mathbf{U}_\theta \mathbf{U}_\theta^T = \mathbf{W}_\theta \mathbf{W}_\theta^T$. Due to the spherical symmetry assumption on $\boldsymbol{\eta}$, $\boldsymbol{\alpha}_\parallel$ and $\boldsymbol{\alpha}_\perp$ are spherically symmetric random variables of dimension $n$ and $m - n$, respectively, with components that are independent. Therefore

$$\mathbb{E}\left[\left\|\mathbf{U}_\theta \hat{\mathbf{D}}_{\theta,\lambda} \boldsymbol{\alpha}_\parallel\right\|_2^2\right] = \mathbb{E}\left[\left\|\hat{\mathbf{D}}_{\theta,\lambda} \boldsymbol{\alpha}_\parallel\right\|_2^2\right]$$
$$= \sum_{i=1}^n \left(\frac{\lambda \gamma_i^2}{\sigma_i^2 + \lambda \gamma_i^2}\right)^2 \mathbb{E}\left[\boldsymbol{\eta}_i^2\right] = \frac{1}{m}\operatorname{Tr}(\hat{\mathbf{D}}_{\theta,\lambda}^2)\mathbb{E}\left[\|\boldsymbol{\eta}\|_2^2\right] \tag{60}$$

382 and

$$\mathbb{E}\left[\mathbf{W}_\theta \|\boldsymbol{\alpha}_\perp\|_2^2\right] = \mathbb{E}\left[\|\boldsymbol{\alpha}_\perp\|_2^2\right] = \frac{m-n}{m}\mathbb{E}\left[\|\boldsymbol{\eta}\|_2^2\right]. \tag{61}$$

383 This completes the proof.

384 *Proof (Proof of Proposition 2)* It is convenient to write Eq. 11 in the abstract form

$$\operatorname{EMD}(\rho_1, \rho_2) = \min_{m \in C(\rho_1, \rho_2)} \mathcal{T}(m). \tag{62}$$

385 In addition, for any $f \in L^1(\Omega)$, let $m_f$ be a minimizer of $\mathcal{T}(f^+, f-)$ over $C(f^+, f-)$ so that
386 $\operatorname{struc}[f] = \mathcal{T}(m_f)$.
387 1. We check absolute homogeneity, positivity, and the triangle inequality.
    (a) To check absolute homogeneity, let $\lambda \in \mathbb{R}$ be a nonzero scalar. By linearity, $m \in C(|\lambda|f, |\lambda|g)$ if and only if $|\lambda|^{-1}m \in C(f, g)$. Therefore

$$\operatorname{EMD}(|\lambda|f, |\lambda|g) = \min_{m \in C(|\lambda|f, |\lambda|g)} \mathcal{T}(m)$$
$$= \min_{m \in C(f,g)} \mathcal{T}(|\lambda|m) = |\lambda| \min_{m \in C(f,g)} \mathcal{T}(m) = |\lambda|\operatorname{EMD}(f, g), \tag{63}$$

388 If $\lambda > 0$, Eq. 63 implies that

$$\operatorname{struc}[\lambda f] = \operatorname{EMD}(\lambda f^+, \lambda f^-) = |\lambda|\operatorname{EMD}(f^+, f^-) = |\lambda|\operatorname{struc}[f] \tag{64}$$

If $\lambda < 0$, then $(\lambda f)^\pm = |\lambda|f^\mp$. Again Eq. 63 implies that

$$\operatorname{struc}[\lambda f] = \operatorname{EMD}((\lambda f)^+, (\lambda f)^-) = \operatorname{EMD}(|\lambda|f^-, |\lambda|f^+)$$
$$= |\lambda|\operatorname{EMD}(f^-, f^+) = |\lambda|\operatorname{EMD}(f^+, f^-) = |\lambda|\operatorname{struc}[f]. \tag{65}$$

389 Finally, if $\lambda = 0$, then the fact that $\operatorname{struc}[\lambda f] = \lambda \operatorname{struc}[f] = 0$ is trivial.

(b) Positivity follows immediately from the positivity of EMD.

(c) The triangle inequality follows from the fact that

$$(f+g)^+ - (f+g)^- = (f^+ - f^-) + (g^+ - g^-) \tag{66}$$

for all $f, g \in L^1(\Omega)$. Thus if $m_f \in C(f^+, f^-)$ and $m_g \in C(g^+, g^-)$, then $m_f + m_g \in C\left((f+g)^+, (f+g)^-\right)$. Along with the triangle inequality for $\mathcal{T}$, this implies that

$$\text{struc}[f+g] \equiv \mathcal{T}(m_{f+g}) \leq \mathcal{T}(m_f + m_g) \leq \mathcal{T}(m_f) + \mathcal{T}(m_g) \equiv \text{struc}[f] + \text{struc}[g]. \tag{67}$$

2. Because $\frac{1}{\|\Omega\|} \int_\Omega (g+c)dx = \frac{1}{\|\Omega\|} \int_\Omega gdx + c$, we have that $g^+ = (g+c)^+$, and $g^- = (g+c)^-$. Therefore

$$\text{struc}[g+c] = \text{EMD}\left((g+c)^+, (g+c)^-\right) = \text{EMD}(g^+, g^-) = \text{struc}[g]. \tag{68}$$

3. Let $g = 0$ in Eq. 68 above. Then

$$\text{struc}[c] = \text{struc}[0] = 0, \quad \forall c \in \mathbb{R}. \tag{69}$$

4. Because the constraint in Eq. 11 involves only the difference of $\rho_1$ and $\rho_2$, it follows that $\text{EMD}(\rho_1, \rho_2) = \text{EMD}(\rho_1 + f, \rho_2 + f)$ for any non-negative $f \in L^1(\Omega)$. Moreover, because $\rho_2$ and $\rho_1$ have the same mass, the average of $\rho_2 - \rho_1$ is zero. Hence,

$$\begin{aligned}
\text{struc}[\rho_2 - \rho_1] &= \text{EMD}(\max(\rho_2 - \rho_1, 0), \max(\rho_1 - \rho_2, 0)) \\
&= \text{EMD}(\max(\rho_2 - \rho_1, 0) + \min(\rho_1, \rho_2), \max(\rho_1 - \rho_2, 0) + \min(\rho_1, \rho_2))
\end{aligned} \tag{70}$$

Since $\forall x, y \in \mathbb{R}, \max(x - y, 0) + \min(x, y) = x$, it follows from Eq. 70 that

$$\text{struc}[\rho_2 - \rho_1] = \text{EMD}(\rho_2, \rho_1) = \text{EMD}(\rho_1, \rho_2) \tag{71}$$

Before proving Thm. 1-3, we will first prove two useful lemmas, which will be used extensively.

**Lemma 2 (EMD triangle inequality)** *Let $\Omega \subset \mathbb{R}^n$ be a bounded set and $f, g, h \in L^\infty(\Omega)$ and $\int_\Omega f dx = \int_\Omega h dx = \int_\Omega g dx$. Then*

$$\text{EMD}(f, g) \leq \text{EMD}(f, h) + \text{EMD}(h, g). \tag{72}$$

*Proof* Recall from Prop. 2 that $\text{struc}[f - g] = \text{EMD}(f, g)$, then by the triangle inequality of $\text{struc}[\cdot]$,

$$\text{EMD}(f, g) = \text{struc}[f - g] \leq \text{struc}[f - h] + \text{struc}[h - g] = \text{EMD}(f, h) + \text{EMD}(h, g) \tag{73}$$

**Lemma 3 (struc$[\cdot]$ and EMD of the mean)** *$\Omega \subset \mathbb{R}^n$ be a bounded set and $f \in L^\infty(\Omega)$ and $\mu = \frac{1}{|\Omega|} \int_\Omega f dx$. Then*

$$\text{struc}[f] = \text{EMD}(f, \mu). \tag{74}$$

*Proof* Recall from Prop. 2 that $\text{EMD}(f, g) = \text{EMD}(f + h, g + h)$, therefore

$$\text{struc}[f] = \text{EMD}(f^+, f^-) = \text{EMD}(f^+ + (\mu - f^-), f^- + (\mu - f^-)) = \text{EMD}(f, \mu). \tag{75}$$

**Lemma 4 (EMD Subadditivity)** *If $\text{EMD}(f_1, g_1)$ and $\text{EMD}(f_2, g_2)$ are well defined, then so too is $\text{EMD}(f_1 + f_2, g_1 + g_2)$, and*

$$\text{EMD}(f_1 + f_2, g_1 + g_2) \leq \text{EMD}(f_1, g_1) + \text{EMD}(f_2, g_2). \tag{76}$$

*Proof* We use the Eq. 10 of the EMD. Let $\pi_1$ and $\pi_2$ satisfy the constraint of Eq. 9 for $\text{EMD}(f_1, g_1)$ and $\text{EMD}(f_2, g_2)$ resp. Then clearly

$$\int_\Omega (\pi_1 + \pi_2) dx^{(2)} = f_1 + f_2$$

$$\int_\Omega (\pi_1 + \pi_2) dx^{(1)} = g_1 + g_2$$

$$\pi_1 + \pi_2 \geq 0, \tag{77}$$

and so by the minimality of the EMD,

$$\begin{aligned}
\text{EMD}(f_1, g_1) + \text{EMD}(f_2, g_2) &= \int_{\Omega \times \Omega} c\pi_1 dx^{(1)} dx^{(2)} + \int_{\Omega \times \Omega} c\pi_2 dx^{(1)} dx^{(2)} \\
&= \int_{\Omega \times \Omega} c(\pi_1 + \pi_2) dx^{(1)} dx^{(2)} \\
&\geq \min_{\pi \geq 0} \int_{\Omega \times \Omega} c\pi dx^{(1)} dx^{(2)} \\
&= \text{EMD}(f_1 + f_2, g_1 + g_2)
\end{aligned} \tag{78}$$

where $\pi$ is subject to the constraints of Eq. 9 where $\rho_1 = f_1 + f_2$ and $\rho_2 = g_1 + g_2$.

**Lemma 5** (EMD **is bounded by the** $L_1$ **norm**) *Let $\Omega$ be a bounded set, and $l \geq \left\| x^{(1)} - x^{(2)} \right\|_2$ for all $x^{(1)}, x^{(2)} \in \Omega$. If $f, g : \Omega \to \mathbb{R}^+$ then*

$$\text{EMD}(f, g) \leq \frac{l}{2} \left\| f - g \right\|_{L^1(\Omega)}. \tag{79}$$

*Proof* Let $\gamma = \int_\Omega (f - g)^+ dx$ and $x^c$ be such that $\|x^c - x\|_2 \leq l/2 \; \forall x \in \Omega$ then

$$\begin{aligned}
\text{EMD}(f, g) = \text{struc}\,[f - g] &\leq \text{EMD}((f - g)^+, \gamma \delta_{x^c}) + \text{EMD}(\gamma \delta_{x^c}, (f - g)^-) \\
&\leq \frac{l}{2} \left\| (f - g)^+ \right\|_{L^1(\Omega)} + \frac{l}{2} \left\| (f - g)^- \right\|_{L^1(\Omega)} = \frac{l}{2} \left\| f - g \right\|_{L^1(\Omega)}
\end{aligned} \tag{80}$$

The last two lines could use a few details between them.

**Lemma 6** (**Expectation bound by the standard deviation**) *Let $\eta$ be a scalar random variable with zero mean such that $\text{Var}[\eta]$ is finite. Then $\mathbb{E}\left[|\eta|\right] \leq \sqrt{\text{Var}[\eta]}$.*

*Proof* Let $\psi$ be the probability distribution for $\eta$. By the Cauchy-Schwarz inequality,

$$\mathbb{E}\left[|\eta|\right] \equiv \int_{-\infty}^\infty |x| \psi(x) dx \leq \left( \int_{-\infty}^\infty x^2 \psi(x) dx \right)^{\frac{1}{2}} \left( \int_{-\infty}^\infty \psi(x) dx \right)^{\frac{1}{2}} = \left( \text{Var}[\eta] \right)^{1/2}. \tag{81}$$

We now proceed to the proof of Theorem 2, but first it is helpful to give a brief summary. To bound the EMD from above, we give a candidate transport plan that is based on the multigrid strategy depicted in Fig. 13 for the case $d = 2$. In this case, the strategy is to divide the domain into square windows with two square panels per side, as shown in Figure 13. The mass in each window is then redistributed in such a way that the new distribution is constant on each window. Each window then becomes a panel in a window that is a factor a factor of two larger in each dimension, and the process is repeated until the distribution on the entire square is constant. For $d > 2$, the plan is the same, except that each window is a hypercube $2^d$ panels. The cost of the complete transport plan can be bounded by the sum of the costs of the transport plan for each step. These costs are computed in the proof below and their sum leads to the bound in Theorem 1.
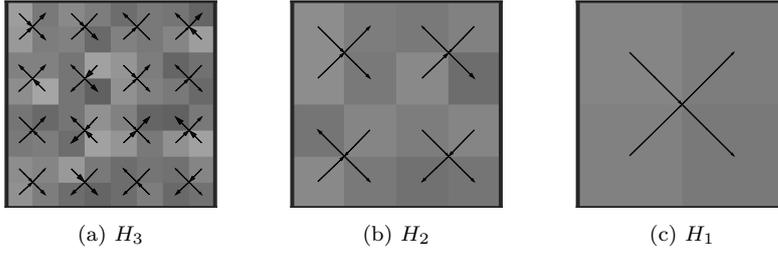
(a) $H_3$          (b) $H_2$          (c) $H_1$

Fig. 13: The multigrid idea of Theorem 1 when $\ell = 3$. At each step, a transport plan is computed in each 2x2 window. Then the same problem is solved at the next coarser scale. In the above figures, the arrow tip area is proportional to the mass transported at each substep. The function $H_i$ is defined in Eq. 91.

*Proof (Proof of Theorem 2)*

Since $\text{struc}\,[h_\ell] = \text{struc}\,[h_\ell - \bar{\mu}]$ we can assume, without loss of generality, that $\bar{\mu} = 0$. Consider the case $\ell = 1$, which will be used for the general setting later. We construct a two-step plan that first moves all of the mass in $h_1^+$ to the point $y^c = (1/2, \ldots, 1/2)$ at the center of the domain and then moves the mass from $y^c$ to $h_1^-$.[2]

Let $\gamma = \int_\Omega h_1^+ dy = \int_\Omega h_1^- dy$, $\mu_0 = \int_\Omega h_1 dy$, and $\gamma_{1,k} = |\eta_{1,k} - \mu_0||\omega_{1,k}|$. Then $\text{EMD}(h_1^+, \gamma\delta_{y^c}) = \text{EMD}(\gamma\delta_{y^c}, h_1^-)$ and

$$\text{struc}\,[h_1] \equiv \text{EMD}(h_1^+, h_1^-) \leq \text{EMD}(h_1^+, \gamma\delta_{y^c}) + \text{EMD}(\gamma\delta_{y^c}, h_1^-)$$

$$= \sum_{k=1}^{2^d} \text{EMD}\left(|\eta_{1,k} - \mu_0|\chi_{1,k}, \gamma_{1,k}\delta_{y^c}\right). \tag{82}$$

Thus we turn our attention to computing the terms in the sum above. First,

$$\text{EMD}(|\eta_{1,k} - \mu_0|\chi_{1,k}, \gamma_{1,k}\delta_{y^c}) = |\eta_{1,k} - \mu_0|\,\text{EMD}(\chi_{1,k}, |\omega_{1,k}|\,\delta_{y^c}). \tag{83}$$

There is only one one admissible transport plan (see from Eq. 10) between $\chi_{1,k}$ and $|\omega_{1,k}|\delta_{y^c}$; it simply moves the mass around each point of $\omega_{1,k}$ to $y^c$:

$$\pi\left(x^{(1)}, x^{(2)}\right) = \chi_{1,k}(x^{(1)}) \times \delta_{y^c}(x^{(2)}) \tag{84}$$

If we consider the more general case where $\omega_{1,k}$ has side length $l$, then upon a change of coordinates,

$$\text{EMD}(\chi_{1,k}, |\omega_{1,k}|\delta_{y^c}) = \int_\Omega \int_\Omega \left\|x^{(1)} - x^{(2)}\right\|_2 \chi_{1,k}(x^{(1)}) \times \delta_{y^c}(x^{(2)})dx^{(1)}dx^{(2)}$$

$$= \int_{\omega_{1,k}} \int_\Omega \left\|x^{(1)} - x^{(2)}\right\|_2 \delta_{y^c}(x^{(2)})dx^{(1)}dx^{(2)}$$

$$= \int_{\omega_{1,k}} \left\|x^{(1)} - y^c\right\|_2 dx^{(1)} = \int_{[0,l]^d} \left\|x^{(1)}\right\|_2 dx^{(1)}$$

$$\leq \sqrt{d} \int_{[0,l]^d} \left\|x^{(1)}\right\|_\infty dx^{(1)} \leq \sqrt{d}\frac{l^{d+1}}{2}$$

$$\tag{85}$$

---

[2] While the definition of the EMD in Eq. 10 is still well-defined for delta function, the formula in Eq. 11 is not. Thus while we use Eq. 11 for numerical calculations, we often rely on Eq. 10 for theoretical bounds.

431   Substituting Eq. 83 and Eq. 85 into Eq. 82 gives

$$\text{struc}\,[h_1] \leq \sum_{i=1}^{2^d} |\eta_{1,k} - \mu_0| \frac{\sqrt{d}l^{d+1}}{2} = \frac{\sqrt{d}}{2^{d+2}} \sum_{k=1}^{2^d} |\eta_{1,k} - \mu_0|, \tag{86}$$

432   where we have used the fact that when $\ell = 1, l = 2^{-1}$. A standard calculation shows that

$$\text{Var}(|\eta_{1,k} - \mu|) \leq \text{Var}(|\eta_{1,k}|), \quad i = 1, \dots, 2^d. \tag{87}$$

433   Further, w.l.o.g. $\mathbb{E}\left[\eta_{1,k}\right] = 0$ and Lemma 6 give:

$$\mathbb{E}\left[|\eta_{1,k} - \mu_0|\right] \leq \sigma \tag{88}$$

434   with Eq. 86 and get

$$\mathbb{E}\left[\text{struc}\,[h_1]\right] \leq \frac{\sqrt{d}2^d}{2^{(d+2)}} \sum_{k=1}^{2^d} \mathbb{E}\left[|n_{1,k} - \mu_0|\right] \leq \frac{\sqrt{d}2^d}{2^{(d+2)}}\sigma = \frac{\sqrt{d}}{4}\sigma. \tag{89}$$

Now we consider the case when $\ell > 1$. Define the functions

$$H_\ell(y) = h_\ell(y) = \sum_{k=1}^{2^{\ell d}} \eta_{\ell,k}\chi_{\ell,k}(y) \tag{90}$$

$$H_i(y) = \sum_{k=1}^{2^{id}} \mu_{i,k}\chi_{i,k}(y), \text{ where } \mu_{i,k} = \frac{1}{|\omega_{i,k}|} \int_{\omega_{i,k}} H_{i+1}(y)dy, \quad i = 0, 1, \dots, \ell-1. \tag{91}$$

Instances of $H_i$ are shown in Fig. 13. The function $h_\ell$ can be written as the telescoping sum

$$h_\ell = H_\ell = (H_\ell - H_{\ell-1}) + (H_{\ell-1} - H_{\ell-2}) + \cdots + (H_2 - H_1) + (H_1 - H_0) + H_0. \tag{92}$$

435   Moreover, because $H_i = \sum_{k=1}^{2^{d(i-1)}} H_i\chi_{i-1,k}$, it follows that

$$H_i - H_{i-1} = \sum_{k=1}^{2^{d(i-1)}} s_{i-1,k}, \quad \text{where } s_{i-1,k}(y) = \left(H_i(y) - \mu_{i-1,k}\right)\chi_{i-1,k}(y). \tag{93}$$

We apply struc $[\cdot]$ to Eq. 92, using Eq. 93, the triangle inequality, and the fact that struc $[H_0] = 0$ (because it is a constant). The result is

$$\text{struc}\,[h_\ell] \leq \sum_{i=1}^{\ell} \text{struc}\,[H_i - H_{i-1}] \leq \sum_{i=1}^{\ell} \sum_{k=1}^{2^{d(i-1)}} \text{struc}\,\left[s_{i-1,k}\right]. \tag{94}$$

436   To evaluate struc $\left[s_{i-1,k}\right]$, we repeat the argument used to generate Eq. 86. This gives

$$\text{struc}\,\left[s_{i-1,k}\right] \equiv \text{EMD}(s_{i-1,k}^+, s_{i-1,k}^-) \leq \frac{\sqrt{d}l^{d+1}}{2} \sum_{k':\omega_{i,k'}\subset\omega_{i-1,k}} |\mu_{i,k'} - \mu_{i-1,k}|. \tag{95}$$

437   By construction,

$$\mu_{i-1,k} = 2^{-d} \sum_{k':\omega_{i,k'}\subset\omega_{i-1,k}} \mu_{i,k}. \tag{96}$$

It follows that the random variable $(\mu_{i+1,k'} - \mu_{i,k})$ that appears in Eq. 95 has zero mean. Thus Lemma 6 applies and

$$\mathbb{E}\left[|\mu_{i,k'} - \mu_{i-1,k}|\right] \leq \left(\mathrm{Var}[|\mu_{i,k'} - \mu_{i-1,k}|]\right)^{\frac{1}{2}} \leq \left(\mathrm{Var}[|\mu_{i,k'}|]\right)^{\frac{1}{2}} := \sigma_i, \tag{97}$$

where the last two inequalities above follows from standard probability theory. Also, because of Eq. 96, another standard probablity result gives

$$\sigma_i = 2^{-\frac{d}{2}} \sigma_{i+1} = \cdots = 2^{-\frac{d}{2}(\ell-i)} \sigma_\ell, \quad i = 1, \ldots, \ell. \tag{98}$$

We now take the expectation of Eq. 95, using the fact that $\omega_{i,k'}$ has side length $l = 2^{-i}$, along with the triangle and Eq. 98,. The result is

$$\mathbb{E}\left[\mathrm{struc}\left[s_{i-1,k}\right]\right] \leq \sqrt{d} 2^{-i(d+1)-1} \sum_{k':\omega_{i,k'} \subset \omega_{i-1,k}} 2^{-\frac{d}{2}(\ell-i)} \sigma_\ell = \sqrt{d} 2^{-\frac{id}{2}-i+d-\frac{d\ell}{2}-1} \sigma_\ell \tag{99}$$

Substituting this bound into Eq. 94 gives

$$\mathbb{E}\left[\mathrm{struc}\left[h_\ell\right]\right] \leq \sum_{i=1}^{\ell} \sum_{k=1}^{2^{d(i-1)}} \sqrt{d} 2^{-\frac{id}{2}-i+d-\frac{d\ell}{2}-1} \sigma_\ell = \frac{\sqrt{d} \sigma_\ell}{2^{1+\frac{\ell d}{2}}} \sum_{i=1}^{\ell} \left(2^{\frac{d}{2}-1}\right)^i \tag{100}$$

If $d = 2$, then $2^{\frac{d}{2}-1} = 1$ and Eq. 100 becomes

$$\mathbb{E}\left[\mathrm{struc}\left[h_\ell\right]\right] = \mathbb{E}\left[\mathrm{struc}\left[H_\ell\right]\right] \leq \frac{2\sigma_\ell}{2^{1+i}} \ell = \frac{\sigma_\ell \ell}{2^\ell}. \tag{101}$$

If $d \geq 3$, then $2^{\frac{d}{2}-1}/(2^{\frac{d}{2}-1} - 1) \leq 4$, so the geometric sum in Eq. 100 is

$$\sum_{i=1}^{\ell} \left(2^{\frac{d}{2}-1}\right)^i = \frac{2^{\left(\frac{d}{2}-1\right)(\ell+1)} - 2^{\frac{d}{2}-1}}{2^{\frac{d}{2}-1} - 1} \leq \frac{2^{\frac{d}{2}-1} 2^{\left(\frac{d}{2}-1\right)\ell}}{2^{\frac{d}{2}-1} - 1} \leq 2^{\frac{\ell d}{2}-\ell+2}. \tag{102}$$

Thus for $d \geq 3$,

$$\mathbb{E}\left[\mathrm{struc}\left[h_\ell\right]\right] \leq \sqrt{d} \sigma_\ell \frac{2^{\frac{\ell\sqrt{d}}{2}-\ell+2}}{2^{1+\frac{\ell\sqrt{d}}{2}}} = \sqrt{d} \sigma_\ell 2^{-\ell+1} \tag{103}$$

Finally, setting $\epsilon = 2^{-\ell}$ gives

$$\mathbb{E}\left[\mathrm{struc}\left[h_\ell\right]\right] \leq \sigma \begin{cases} -\epsilon_\ell \log(\epsilon_\ell) & \text{when } d = 2 \\ 2\sqrt{d} \epsilon_\ell & \text{when } d > 2 \end{cases} \tag{104}$$

This completes the proof.

*Proof (Proof of Lemma 1)* The proof follows directly from the definition of $h_\ell$ in the statement of Thm. 1:

$$\mathbb{E}\left[\|h_\ell\|_2^2\right] = \mathbb{E}\left[\int_{[0,1)^d} (h_\ell(y))^2 \, dy\right] = \sum_{k=1}^{2^{\ell d}} \mathbb{E}\left[\eta_{\ell,k}^2\right] 2^{-\ell d} = 2^{-\ell d} \sum_{k=1}^{2^{\ell d}} \sigma^2 = \sigma^2. \tag{105}$$

*Proof (Proof of Theorem 2)*

Without loss of generality, assume that $\phi$ is positive a.e. (If not, simply replace $\phi$ by $\phi - \mathrm{ess\,inf}\,\phi$ and use Eq. 68.) By construction, $\phi$ and $R_\ell\phi$ have the same average over $Y$, which we denote by $\mu$. Thus by Lemmas 2 and 3,

$$\mathrm{struc}\left[R_\ell\phi\right] = \mathrm{EMD}(R_\ell\phi, \mu) \leq \mathrm{EMD}(R_\ell\phi, \phi) + \mathrm{EMD}(\phi, \mu) = \mathrm{EMD}(R_\ell\phi, \phi) + \mathrm{struc}\left[\phi\right]. \tag{106}$$

Hence

$$\text{struc}\,[R_\ell\phi] - \text{struc}\,[\phi] \le \text{EMD}(R_\ell\phi, \phi). \tag{107}$$

One the other hand, switching the roles of $R_\ell\phi$ and $\phi$ Eq. 106 gives

$$\text{struc}\,[\phi] - \text{struc}\,[R_\ell\phi] \le \text{EMD}(R_\ell\phi, \phi) \tag{108}$$

Together Eq. 107 and Eq. 107 imply the bound

$$|\text{struc}\,[R_\ell\phi] - \text{struc}\,[R_\ell\phi]\,| \le \text{EMD}(R_\ell\phi, \phi). \tag{109}$$

We now bound $\text{EMD}(R_\ell\phi, \phi)$. For any $\ell, i$ $\int_{\omega_{\ell,i}} R_\ell\phi dy = \int_{\omega_{\ell,i}} \phi dy$. Thus by Lemma 4,

$$\text{EMD}(R_\ell\phi, \phi) \le \sum_{i=1}^{2^{\ell d}} \text{EMD}(R_\ell\phi\chi_{\ell,i}, \phi\chi_{\ell,i}) \tag{110}$$

and further by Lemma 5, for $i = 1, \ldots, 2^{\ell d}$

$$\text{EMD}(R_\ell\phi\chi_{\ell,i}, \phi\chi_{\ell,i}) \le \|R_\ell\phi - \phi\|_{L^1(\omega_{\ell,i})}\, d^{1/2} 2^{-\ell} \tag{111}$$

Now we bound $\|R_\ell\phi - \phi\|_{L^1(\omega_{\ell,i})}$. Since $\phi \in C^1\left(\overline{Y}\right)$, it follows that, for $y \in \omega_{\ell,i}$

$$|R_\ell\phi(y) - \phi(y)| = \frac{1}{|\omega_{\ell,i}|} \left| \int_{\omega_{\ell,i}} (\phi(y') - \phi(y)) dy' \right|$$

$$\le \sup_{y\in\omega_{\ell,i}} |\nabla\phi(y)| \sup_{y\in\omega_{\ell,i}} |y' - y| \le d^{1/2} 2^{-\ell} \sup_{y\in\omega_{\ell,i}} |\nabla\phi(y)| \tag{112}$$

Therefore

$$\|R_\ell\phi - \phi\|_{L^1(\omega_{\ell,i})} \le |\omega_{\ell,i}| d^{1/2} 2^{-\ell} \sup_{y\in\omega_{\ell,i}} |\nabla\phi(y)| = d^{1/2} 2^{-(d+1)\ell} \sup_{y\in\omega_{\ell,i}} |\nabla\phi(y)|. \tag{113}$$

Combining Eq. 109, Eq. 111, and Eq. 113 yields

$$|\text{struc}\,[R_\ell\phi] - \text{struc}\,[\phi]\,| \le \sum_{i=1}^{2^{\ell d}} d 2^{-(d+2)\ell} \sup_{y\in\omega_{\ell,i}} |\nabla\phi(y)| \le d 2^{-2\ell} \sup_{y\in Y} |\nabla\phi(y)| \equiv C(|\nabla\phi|) d\epsilon_\ell^2, \tag{114}$$

where $C(|\nabla\phi|) = \sup_{y\in Y} |\nabla\phi(y)|$ and $\epsilon_\ell = 2^{-\ell}$. This completes the proof.

## B Line Integral Operators

Recall from Section 3 the spaces $\mathcal{U}$ and $\mathcal{B}$ of functions defined on domains $X$ and $Y$, respectively. An operator $\mathcal{L} \colon \mathcal{U} \to \mathcal{B}$ is a line integral operators (LIO), if $\forall f \in \mathcal{U}$,

$$(\mathcal{L}f)(y) = \int_{P_y} f(x) dx = \int_0^1 f(\hat{x}(t; y))\hat{x}'(t; y) dt, \tag{115}$$

where for each $y \in Y$, $P_y = \{\hat{x}(t; y) : t \in (0, 1)\} \subset X$, and $\hat{x}(t; y)$ is continuous in $t$ and $y$. In particular, if $f$ is a continuous on $X$, then $\mathcal{L}f$ is continuous on $Y$. Figs. 14b and 14a illustrate a LIO in two dimensions. The recipe we used to generate examples of $\hat{x}$ is given below.

To discretize $\mathcal{L}$, we generate a path $P_y$ for each hypercube $\omega \subset Y$. Line integrals along these paths are approximated via quadrature. For all LIOs, we use same the quadratures, and $X$, and $Y$.

To construct the LIO for Experiments 1 - 3, we do the following.

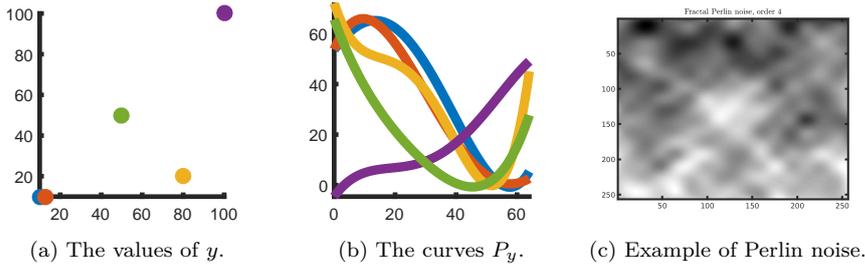(a) The values of $y$.　　　(b) The curves $P_y$.　　　(c) Example of Perlin noise.

Fig. 14: An example of a LIO. Points on the right are used to generate curves on the left of the same color. Coefficients for the parameterization in Eq. 117 of $P_y$ come from Perlin noise.

1. **Construction of numerical grids.** In all of our computational examples, the domains $X$ and $Y$ are unit squares in $\mathbb{R}^2$. We discretize these domains with $N^x$ and $N^y$ points, respectively, on each side and define grid points

$$x_{i,j} = (i\Delta x, j\Delta x), \quad 0 \le i, j \le N^x - 1, \tag{116a}$$
$$y_{k,l} = (k\Delta y, l\Delta y), \quad 0 \le k, l \le N^y - 1, \tag{116b}$$

where $\Delta x = 1/N^x$ and $\Delta y = 1/N^y$. We then generate values $u_{i,j}$ by sampling a prescribed function at the points $x_{i,j}$. An illustrative example is given in Fig. 3a, where piecewise smooth rings have been sampled on a $64 \times 64$ grid.

2. **Generation of smooth paths.** To form $\hat{x}$, we first sample coefficients $\alpha_{p,r}$ for $p = 0, \ldots, 4$ and $r = 1, 2$ from Perlin noise [25, 26] of order four. In Fig. 14c, a realization of one such coefficient as a function of $y$ is shown on a $256 \times 256$ grid. Given these coefficients, we let $\bar{x} = (x^{(1)}, x^{(2)})$ be polynomials in $t$:

$$\bar{x}^{(r)}(t; y_{k,l}) = \sum_{p=0}^{4} \frac{\alpha_{p,r}(y_{k,l})}{p!} t^p, \quad r = 1, 2, \tag{117}$$

and then let $\hat{x}$ be the following normalization of $\bar{x}$:

$$\hat{x}^{(r)}(t; y_{k,l}) = \frac{\bar{x}^{(r)}(t; y_{k,l}) - \min_s \bar{x}^{(r)}(s; y_{k,l})}{\max_s \bar{x}^{(r)}(s; y_{k,l}) - \min_t \bar{x}^{(r)}(s; y_{k,l})}, \quad r = 1, 2. \tag{118}$$

3. To generate the components of $\mathbf{L}$, we first compute

$$I_{k,l} = \left\{ (i,j) | \exists s \in [0,1] \text{ s.t.} (i,j) = \underset{(i,j)}{\operatorname{argmin}} \left\| x_{i,j} - \hat{x}(s; y_{k,l}) \right\| \right\}. \tag{119}$$

and then set the values of $\mathbf{L}$ directly by

$$L_{(k,l),(i,j)} = \begin{cases} \frac{1}{|I_{k,l}|} & \text{if } (i,j) \in I_{k,l}. \\ 0 & \text{else} \end{cases} \tag{120}$$

## References

1. Simon R Arridge. Optical tomography in medical imaging. Inverse problems, 15(2):R41, 1999.
2. Stephen Becker. Lbfgsb (l-bfgs-b) mex wrapper, 2012–2015.
3. Charles George Broyden. The convergence of a class of double-rank minimization algorithms 1. general considerations. IMA Journal of Applied Mathematics, 6(1):76–90, 1970.

4. Moustafa T Chahine. Inverse problems in radiative transfer: Determination of atmospheric parameters. Journal of the Atmospheric Sciences, 27(6):960–967, 1970.

5. Tony F Chan and Jianhong Jackie Shen. Image processing and analysis: variational, PDE, wavelet, and stochastic methods, volume 94. Siam, 2005.

6. Ingrid Daubechies. Orthonormal bases of compactly supported wavelets. Communications on pure and applied mathematics, 41(7):909–996, 1988.

7. Bjorn Engquist and Brittany D Froese. Application of the wasserstein metric to seismic signals. arXiv preprint arXiv:1311.4581, 2013.

8. Bjorn Engquist, Brittany D Froese, and Yunan Yang. Optimal transport for seismic full waveform inversion. arXiv preprint arXiv:1602.01540, 2016.

9. Lawrence C Evans. Partial differential equations and monge-kantorovich mass transfer. Current developments in mathematics, 1997(1):65–126, 1997.

10. Lawrence C Evans and Wilfrid Gangbo. Differential equations methods for the Monge-Kantorovich mass transfer problem, volume 653. American Mathematical Soc., 1999.

11. Roger Fletcher. A new approach to variable metric algorithms. The computer journal, 13(3):317–322, 1970.

12. Anthony Freeman. Sar calibration: An overview. IEEE Transactions on Geoscience and Remote Sensing, 30(6):1107–1121, 1992.

13. Donald Goldfarb. A family of variable-metric methods derived by variational means. Mathematics of computation, 24(109):23–26, 1970.

14. Tom Goldstein and Stanley Osher. The split bregman method for l1-regularized problems. SIAM journal on imaging sciences, 2(2):323–343, 2009.

15. Gene H Golub. Matrix computations. Johns Hopkins University Press, 1996.

16. Gene H Golub, Per Christian Hansen, and Dianne P O'Leary. Tikhonov regularization and total least squares. SIAM Journal on Matrix Analysis and Applications, 21(1):185–194, 1999.

17. Per Christian Hansen. Analysis of discrete ill-posed problems by means of the l-curve. SIAM review, 34(4):561–580, 1992.

18. Per Christian Hansen and Dianne Prost O'Leary. The use of the l-curve in the regularization of discrete ill-posed problems. SIAM Journal on Scientific Computing, 14(6):1487–1503, 1993.

19. Marc C Kennedy and Anthony O'Hagan. Bayesian calibration of computer models. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 63(3):425–464, 2001.

20. Andreas Kirsch. An introduction to the mathematical theory of inverse problems, volume 120. Springer Science & Business Media, 2011.

21. Wuchen Li, Stanley Osher, and Wilfrid Gangbo. A fast algorithm for earth mover's distance based on optimal transport and l1 type regularization. arXiv preprint arXiv:1609.07092, 2016.

22. Wuchen Li, Ernest K Ryu, Stanlet Osher, Wotao Yin, and Wolfred Gangbo. A parallel method for earth mover's distance. Journal of Scientific Computing, page 75(1), 2018.

23. Stephane G Mallat. Multiresolution approximations and wavelet orthonormal bases of $l^2(r)$. Transactions of the American mathematical society, 315(1):69–87, 1989.

24. Dean S Oliver, Albert C Reynolds, and Ning Liu. Inverse theory for petroleum reservoir characterization and history matching. Cambridge University Press, 2008.

25. Ken Perlin. An image synthesizer. ACM Siggraph Computer Graphics, 19(3):287–296, 1985.

26. Ken Perlin. Improving noise. In ACM Transactions on Graphics (TOG), volume 21, pages 681–682. ACM, 2002.

27. Leonid I Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal algorithms. Physica D: Nonlinear Phenomena, 60(1-4):259–268, 1992.

28. Ernest Ryu, Yongxin Chen, Wuchen Li, and Stanley Osher. Vector and matrix optimal mass transport: Theory, algorithm, and applications. arXiv, 2017.

29. Kai Schneider, Romain Nguyen van Yen, Nicolas Fedorczak, Frederic Brochard, Gerard Bonhomme, Marie Farge, and Pascale Monier-Garbet. Tomographic reconstruction of tokamak plasma light emission using wavelet-vaguelette decomposition. In APS Meeting Abstracts, 2012.

30. Uwe Schneider, Eros Pedroni, and Antony Lomax. The calibration of ct hounsfield units for radiotherapy treatment planning. Physics in Medicine & Biology, 41(1):111, 1996.

31. David F Shanno. Conditioning of quasi-newton methods for function minimization. Mathematics of computation, 24(111):647–656, 1970.

32. Cédric Villani. Optimal transport: old and new, volume 338. Springer Science & Business Media, 2008.
33. Andreas Wingen, MW Shafer, Ezekial A Unterberg, Judith C Hill, and Donald L Hillis. Regularization of soft-x-ray imaging in the diii-d tokamak. Journal of Computational Physics, 289:83–95, 2015.
34. Yunan Yang, Björn Engquist, Junzhe Sun, and Brittany F Hamfeldt. Application of optimal transport and the quadratic wasserstein metric to full-waveform inversion. Geophysics, 83(1):R43–R62, 2018.
35. Ciyou Zhu, Richard H Byrd, Peihuang Lu, and Jorge Nocedal. Lbfgs-b: Fortran subroutines for large-scale bound constrained optimization. Report NAM-11, EECS Department, Northwestern University, 1994.