



# On the representation and learning of monotone triangular transport maps

Ricardo Baptista, Youssef Marzouk, Olivier Zahm

## ► To cite this version:

Ricardo Baptista, Youssef Marzouk, Olivier Zahm. On the representation and learning of monotone triangular transport maps. Foundations of Computational Mathematics, 2023, 10.1007/s10208-023-09630-x . hal-03060198v1

**HAL Id: hal-03060198**

**<https://inria.hal.science/hal-03060198v1>**

Submitted on 30 Jan 2023 (v1), last revised 9 Jan 2024 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution| 4.0 International License

# On the representation and learning of monotone triangular transport maps

Ricardo Baptista<sup>1</sup>, Youssef Marzouk<sup>1</sup>, and Olivier Zahm<sup>2</sup>

<sup>1</sup>*Massachusetts Institute of Technology  
Cambridge, MA 02139 USA  
[rsb@mit.edu](mailto:rsb@mit.edu), [ymarz@mit.edu](mailto:ymarz@mit.edu),*

<sup>2</sup>*Université Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK  
38000, Grenoble, France  
[olivier.zahm@inria.fr](mailto:olivier.zahm@inria.fr)*

**Abstract:** Transportation of measure provides a versatile approach for modeling complex probability distributions, with applications in density estimation, Bayesian inference, generative modeling, and beyond. Monotone triangular transport maps—approximations of the Knothe–Rosenblatt (KR) rearrangement—are a canonical choice for these tasks. Yet the representation and parameterization of such maps have a significant impact on their generality and expressiveness, and on properties of the optimization problem that arises in learning a map from data (e.g., via maximum likelihood estimation). We present a general framework for representing monotone triangular maps via invertible transformations of smooth functions. We establish conditions on the transformation such that the associated infinite-dimensional minimization problem has no spurious local minima, i.e., all local minima are global minima; and we show for target distributions satisfying certain tail conditions that the unique global minimizer corresponds to the KR map. Given a sample from the target, we then propose an adaptive algorithm that estimates a sparse semi-parametric approximation of the underlying KR map. We demonstrate how this framework can be applied to joint and conditional density estimation, likelihood-free inference, and structure learning of directed graphical models, with stable generalization performance across a range of sample sizes.

**Keywords and phrases:** Transportation of measure, Knothe–Rosenblatt rearrangement, normalizing flows, monotone functions, infinite-dimensional optimization, adaptive approximation, multivariate polynomials, wavelets, density estimation.

**MSC2020 subject classifications:** 65C20, 49Q22, 62G07.

## Contents

1	Introduction . . . . .	1
2	Triangular transport for density estimation and simulation . . . . .	4
3	Representing and learning continuous monotone functions . . . . .	6
4	Adaptive parameterization of transport maps . . . . .	13
5	Numerical experiments . . . . .	17
6	Conclusions . . . . .	24
	Acknowledgments . . . . .	26
A	Proofs and theoretical details . . . . .	26
B	Multi-index refinement for the wavelet basis . . . . .	35
C	Architecture details of alternative methods . . . . .	36
	References . . . . .	36

## 1. Introduction

Many sampling, estimation, and inference algorithms seek to characterize a somehow intractable or complex probability distribution  $\mu$  on  $\mathbb{R}^d$ . Transportation of measure provides a useful and versatile approach to this problem. The underlying idea is to construct a coupling of  $\mu$  with a tractable “reference” distribution  $\nu$  on  $\mathbb{R}^d$ —for instance, a standard normal. Formally, one jointly constructs a pair of random variables  $(\mathbf{X}, \mathbf{Z})$  such that  $\mathbf{X} \sim \mu$  and  $\mathbf{Z} \sim \nu$ . A special class of couplings is given by *deterministic* transformations  $S: \mathbb{R}^d \rightarrow \mathbb{R}^d$  such that  $S(\mathbf{X}) = \mathbf{Z}$  in distribution, and a transformation  $S$  that satisfies this property is called a *transport map* [67]. As a result, if the transport map is invertible, one can generate realizations of  $\mathbf{X}$  by first simulating  $\mathbf{Z}$ . If  $\mu$  and  $\nu$  have densities  $\pi$  and  $\eta$  with respect to a common base measure, one can also explicitly represent the target density  $\pi$  as a transformation of the reference density  $\eta$ . The construction of such transport maps has found numerous applications: density estimation [60, 2, 16, 20]; variational Bayesian inference [19, 48, 9]; generative modeling of images, video, and other structured objects [40, 24]; likelihood-free inference [41, 45, 31]; and beyond.

In general, there exist infinitely many transport maps between two absolutely continuous distributions on  $\mathbb{R}^d$ . In this paper, we focus on a specific, canonical choice: *triangular* transport maps [7] of the form

$$S(\mathbf{x}) = \begin{bmatrix} S_1(x_1) \\ S_2(x_1, x_2) \\ \vdots \\ S_d(x_1, \dots, x_d) \end{bmatrix}, \quad (1)$$

where each component function  $S_k$  depends only on the  $k$  variables  $\mathbf{x}_{\leq k} := (x_1, \dots, x_k)$  and is monotone increasing with respect to the last input  $x_k$ . In particular, each function  $S_k$  encodes an increasing rearrangement [67] between ordered marginal conditionals of  $\mu$  and  $\nu$ , i.e.,  $\mu(dx_k | \mathbf{x}_{<k})$  and  $\nu(dz_k | \mathbf{z}_{<k})$ . Later we will discuss properties of such triangular maps—known also as Knothe–Rosenblatt (KR) rearrangements [67, 49, 51]—more precisely, but we comment here on their utility. First of all, triangular structure facilitates computational tractability:  $S$  is easy to invert and the determinant of its Jacobian (a lower triangular matrix) is easy to evaluate (see Section 2). Triangular maps have thus been used extensively in the density estimation, inference, and generative modeling applications noted above. For example, triangular maps are the building blocks of many *normalizing flows*, popularized by the machine learning community [25, 42]; specifically, autoregressive normalizing flows define triangular maps [22] via particular structural choices and parameterizations, and compose these maps to produce more expressive transformations. More fundamentally, because triangular maps expose certain conditionals of  $\mu$ , they are particularly well suited to *conditional density estimation* and *conditional sampling*; we will describe this link in Section 2. Triangular maps also inherit *sparsity* from the conditional independence properties of  $\mu$  and  $\nu$ , as detailed in [59].

It is worth noting, of course, that other canonical choices of transport have different attractive features. For instance, optimal transport maps are invariant under relabeling of the coordinates or more general isometries on  $\mathbb{R}^d$  (unlike triangular maps), and have deep links to partial differential equations [67]. But optimal transport maps are in general more challenging to represent, evaluate, and estimate from data in the continuous setting; moreover, they do not enjoy such a direct link to conditioning or to graphical models.

Many representations and finite-dimensional parameterizations of monotone triangular maps have been proposed in recent years. These include representations based on polynomials [33, 22], radial basis functions [60], neural networks of varying capacity [17, 24], and tensor decompositions [14, 15]. A core challenge in this setting is to satisfy the monotonicity constraint  $\partial_{x_k} S_k > 0$ . For instance,

one might enforce the monotonicity constraint at a finite collection of points in the support of  $\pi$  [44], but this approach cannot in general guarantee that  $S$  is monotone over the entire support of  $\pi$ . Other efforts have sought to enforce monotonicity by construction—via the parameterization of  $S$  itself. For example, [43] employs map components with affine dependence on the last variable, i.e.,  $S_k(\mathbf{x}_{\leq k}) = \alpha(\mathbf{x}_{<k}) + \exp(\beta(\mathbf{x}_{<k}))x_k$ , where  $\alpha$  and  $\beta$  are neural networks. While  $S$  is then guaranteed to be monotone, it can only represent a restricted class of distributions  $\mu$ . (If  $\nu$  is Gaussian, then  $\mu$  can only be a product of Gaussian marginal conditionals.) Such representations therefore cannot consistently approximate the *true* KR rearrangement for general  $\mu$ . Recent work [61] has shown that a *composition* of such affine maps, interleaved with rotations and permutations, can approximate a general class of distributions, though approximation rates remain unknown. The required rotations or permutations break the overall triangular structure of the transformation, however. Alternatively, to increase the “expressiveness” of a given triangular function, [69, 21, 18] have introduced more general parametric representations of the monotone component functions  $S_k$ . For distributions with analytic densities on *bounded domains*, a complete approximation theory for the KR map was recently developed in [71, 72], using polynomial or ReLU neural network representations of range-constrained monotone triangular functions.

Despite these myriad proposals, relatively little attention has been paid to the structure and tractability of the *optimization problem* involved in learning triangular maps (e.g., in estimating a map given an i.i.d. sample  $\{\mathbf{X}^i\}$  from  $\mu$ ). Properties of this optimization problem are intimately tied to the representation and parameterization of the associated map. It is desirable to have a flexible and general representation—one that can consistently recover the KR rearrangement for a broad class of distributions  $(\mu, \nu)$ —that at the same time makes optimization tractable. It is also desirable to have adaptivity: a parameterization whose size or complexity can be adapted to properties of the target distribution and the available sample size, for good empirical statistical performance.

This paper directly addresses these desiderata. We do so by developing and analyzing a *functional framework* for representing and learning triangular maps. Our main contributions are as follows. We propose a general representation of monotone triangular functions, based on a *rectification* operator  $\mathcal{R}_k$  that transforms sufficiently smooth non-monotone functions  $f_k : \mathbb{R}^k \rightarrow \mathbb{R}$  into monotone component functions  $S_k$  of a triangular map. This operator takes the form

$$\mathcal{R}_k(f_k)(\mathbf{x}_{\leq k}) = f_k(\mathbf{x}_{<k}, 0) + \int_0^{x_k} g(\partial_k f_k(\mathbf{x}_{<k}, t)) dt, \quad (2)$$

where  $g : \mathbb{R} \rightarrow \mathbb{R}_{>0}$  is a positive function. We then analyze the infinite-dimensional optimization problem associated with learning maps from data, recast as optimization over functions  $\{f_k\}_{k=1}^d$ . Specifically, we establish conditions on the rectification operator and on the target distribution such that the resulting optimization problem is well-posed, smooth, and has *no spurious local minima*. Under further conditions on the target distribution (essentially that it has Gaussian tails), we show that the optimization problem has a *unique global minimizer* corresponding to the KR map.

These theoretical results guarantee, in practice, fast and reliable learning of monotone triangular maps given an appropriate function space  $V_k$  for each  $f_k$ . The second main contribution of our paper is algorithmic: given a hierarchical basis (e.g., polynomials or wavelets) for each  $V_k$ , we propose a greedy adaptive procedure to learn parametric representations of  $f_k$ . The procedure naturally produces map representations that are *sparse* and interpretable—in particular, it exploits and implicitly discovers conditional independence. We use these learned maps for density estimation, given an i.i.d. sample  $\{\mathbf{X}^i\}$  from  $\pi$ . Maintaining a strict triangular structure also exposes marginal conditionals of the target density, thus enabling *conditional density estimation*. Our numerical experiments show that the algorithm provides robust performance at small-to-moderate sample sizes, and constitutes a *semi-parametric* approach that naturally links map complexity to the size of the data.

The remainder of the paper is organized as follows. Section 2 recalls properties of triangular transport maps and introduces some estimation problems of interest. Our main theoretical contributions are in Section 3, which introduces a framework for representing monotone triangular maps and analyzes the resulting optimization problem. Section 4 then introduces our greedy adaptive procedure for learning maps, and Section 5 contains numerical experiments. Proofs of certain theoretical results are deferred to the appendix.

## 2. Triangular transport for density estimation and simulation

Consider the unsupervised learning problem of approximating a target probability density function  $\pi$  defined on  $\mathbb{R}^d$ , given an i.i.d. sample from  $\pi$ . Our goal is to construct a sufficiently smooth and invertible map  $S: \mathbb{R}^d \rightarrow \mathbb{R}^d$  such that the pullback density

$$S^\sharp \eta(\mathbf{x}) = \eta(S(\mathbf{x})) |\det \nabla S(\mathbf{x})|, \quad (3)$$

is a good approximation to  $\pi$ , where  $\eta$  is a simple/tractable probability density function on  $\mathbb{R}^d$ . The choice of  $\eta$  is a degree of freedom of the method, and here we take  $\eta(\mathbf{x}) \propto \exp(-\|\mathbf{x}\|^2/2)$ ; i.e.,  $\eta$  is the density of the standard normal distribution on  $\mathbb{R}^d$ , where  $\|\cdot\|$  is the canonical norm of  $\mathbb{R}^d$ .

To ensure invertibility of  $S$ , a common practice is to constrain  $S$  to be an increasing lower triangular map of the form (1), where each component  $S_k: \mathbb{R}^k \rightarrow \mathbb{R}$  is such that  $x_k \mapsto S_k(\mathbf{x}_{<k}, x_k)$  is increasing for all  $\mathbf{x}_{<k} = (x_1, \dots, x_{k-1}) \in \mathbb{R}^{k-1}$ . Such a map is easy to invert<sup>1</sup> and has a lower triangular Jacobian  $\nabla S(\mathbf{x})$  so that  $|\det \nabla S(\mathbf{x})| = \prod_{k=1}^d \partial_k S_k(\mathbf{x}_{\leq k})$  is readily computable. This structure in fact corresponds to the *Knothe–Rosenblatt* (KR) rearrangement  $S_{\text{KR}}$ : the increasing lower triangular map satisfying

$$\pi(\mathbf{x}) = S_{\text{KR}}^\sharp \eta(\mathbf{x}).$$

For a measure  $\mu$  on  $\mathbb{R}^d$  that is absolutely continuous with respect to a Gaussian measure  $\nu$  (and hence has a density  $\pi$  with respect to the Lebesgue measure), the KR rearrangement  $S_{\text{KR}}$  exists and is the *unique* map of the form (1) that *pulls back*  $\eta$  to  $\pi$  (or equivalently *pushes forward*  $\pi$  to  $\eta$ ), up to sets of measure zero [7].

A useful way to measure the difference between  $\pi$  and its approximation  $S^\sharp \eta$  is the Kullback–Leibler (KL) divergence  $\mathcal{D}_{\text{KL}}(\pi \| S^\sharp \eta) = \int \log(\pi/S^\sharp \eta) d\pi$ . As explained below, this choice has direct links to maximum likelihood estimation of  $S$ . Furthermore, the following inequality shows that convergence in the KL sense  $\mathcal{D}_{\text{KL}}(\pi \| S^\sharp \eta) \rightarrow 0$  implies convergence of  $S$  towards  $S_{\text{KR}}$  in the  $L_\pi^2$  sense. This result is a direct consequence of Corollary 3.10 in [7]; see Appendix A.1 for a proof.

**Proposition 1.** *Let  $S_{\text{KR}}$  be the KR rearrangement pushing forward a distribution with density  $\pi$  on  $\mathbb{R}^d$  to the standard normal distribution on  $\mathbb{R}^d$ , with density  $\eta$ . For any map  $S: \mathbb{R}^d \rightarrow \mathbb{R}^d$  as in (1), we have*

$$\int \|S_{\text{KR}}(\mathbf{x}) - S(\mathbf{x})\|^2 d\pi(\mathbf{x}) \leq 2\mathcal{D}_{\text{KL}}(\pi \| S^\sharp \eta). \quad (4)$$

Since the standard normal density  $\eta$  is a product of its marginal densities, the KL divergence decomposes as

$$\mathcal{D}_{\text{KL}}(\pi \| S^\sharp \eta) = \sum_{k=1}^d \mathcal{J}_k(S_k) - \mathcal{J}_k(S_{\text{KR},k}), \quad (5)$$

<sup>1</sup>For any  $\mathbf{z} \in \mathbb{R}^d$ ,  $\mathbf{x} = S^{-1}(\mathbf{z})$  can be computed recursively by  $x_k = T^k(\mathbf{x}_{<k}, z_k)$  for  $k = 1, \dots, d$ , where the function  $T^k(\mathbf{x}_{<k}, \cdot)$  is the inverse of  $x_k \mapsto S_k(\mathbf{x}_{<k}, x_k)$ .

where the functionals  $\mathcal{J}_1, \dots, \mathcal{J}_d$  are given by

$$\mathcal{J}_k(s) = \int \left( \frac{1}{2} s(\mathbf{x}_{\leq k})^2 - \log |\partial_k s(\mathbf{x}_{\leq k})| \right) \pi(\mathbf{x}) d\mathbf{x}. \quad (6)$$

Minimizing the KL divergence (5) over triangular maps  $S$  of the form (1) is therefore equivalent to *independently* minimizing each objective functional  $\mathcal{J}_k$  to find the associated map component  $S_k$  [33]. Solution of these optimization problems is thus embarrassingly parallel. This parallel structure was also exploited for Cholesky factorization via KL minimization in [52]. In addition, minimizing each objective  $s \mapsto \mathcal{J}_k(s)$  over functions  $s: \mathbb{R}^k \rightarrow \mathbb{R}$  that are strictly increasing in the last variable is a *strictly convex* optimization problem; see Lemma 10.

Given an i.i.d. sample  $\{\mathbf{X}^i\}_{i=1}^n$  from  $\pi$ , we can replace the expectation in (6) by the sample average, which yields the objective

$$\hat{\mathcal{J}}_k(s) = \frac{1}{n} \sum_{i=1}^n \left( \frac{1}{2} s(\mathbf{X}_{\leq k}^i)^2 - \log |\partial_k s(\mathbf{X}_{\leq k}^i)| \right). \quad (7)$$

Minimizing (7) under the constraint  $\partial_k s(\mathbf{x}_{\leq k}) > 0$  produces an estimator  $\hat{S}_k$  of  $S_{\text{KR},k}$ . The collection of all such component functions defines an estimator  $\hat{S} = (\hat{S}_1, \dots, \hat{S}_d)$  of the KR rearrangement, along with an estimate of the density  $\pi$  as  $\hat{\pi}(\mathbf{x}) := \hat{S}^\# \eta(\mathbf{x})$ . This  $\hat{S}$  is also the maximum likelihood estimator of  $S_{\text{KR}}$ , i.e.,

$$\hat{S} = \arg \max_{\{S \in (1), \partial_k S_k > 0, k=1, \dots, d\}} \sum_{i=1}^n \log S^\# \eta(\mathbf{X}^i),$$

with the optimization being over the space of monotone increasing and triangular maps of the form (1). This connection between maximum likelihood estimation and minimization of an empirical forward KL divergence is standard. Furthermore, convexity of the optimization objective is preserved when replacing the expectation in (6) by the sample average in (7).

A core question when estimating maps and densities by minimizing (7) is how to parameterize sufficiently expressive monotone map components  $S_k$ —i.e., maps capable of representing a wide class of distributions  $\pi$ —while ensuring that the optimization problem can be solved efficiently. As explained in the introduction, this question is intimately tied to the monotonicity constraint  $\partial_k S_k(\mathbf{x}_{\leq k}) > 0$ . For example, choosing a linear parameterization for  $S_k$  that admits only affine dependence on  $x_k$  (to easily enforce monotonicity) allows the map component be identified efficiently through the solution of a least-squares problem (see [58, Appendix A]), but such maps can only capture distributions that factor into a sequence of Gaussian marginal conditionals. On the other hand, a more complex ansatz for  $S_k$  will often yield a much more difficult optimization problem. Note that with any nonlinear parameterization of  $S_k$ , the convexity of (6) and (7) with respect to  $S_k$  does not in general yield convexity in the parameters. As we shall demonstrate later, many nonlinear parameterizations that enforce monotonicity yield optimization problems that may not even be smooth, and that have many local minima. We will address these issues in Section 3.

**Conditional density estimation and sampling** Another important feature of the triangular structure (1) is that each component of the map represents one marginal conditional density of  $\pi$ . More precisely, in the present setting where  $\eta$  is a product density,  $S_{\text{KR},k}$  pushes forward the marginal conditional  $\pi_k(x_k | \mathbf{x}_{<k})$  to the  $k$ -th marginal of the reference  $\eta_k(z_k)$  [51, Section 2.3]. Now partition  $\mathbf{x} = (\mathbf{y}, \mathbf{w})$ , where  $\mathbf{y} \in \mathbb{R}^m$  and  $\mathbf{w} \in \mathbb{R}^p$ , with  $d = m + p$ . This property of the KR map lets us estimate the conditional probability density function  $\pi(\mathbf{w} | \mathbf{y})$ , for any value of  $\mathbf{y}$ , given a sample



$\{(\mathbf{Y}^i, \mathbf{W}^i)\}_{i=1}^n$  from the joint density  $\pi(\mathbf{y}, \mathbf{w})$ . Observe that the KR map immediately has the block structure:

$$S(\mathbf{y}, \mathbf{w}) = \begin{bmatrix} S^{\mathcal{Y}}(\mathbf{y}) \\ S^{\mathcal{W}}(\mathbf{y}, \mathbf{w}) \end{bmatrix}, \quad (8)$$

where  $\mathbf{y} \mapsto S^{\mathcal{Y}}(\mathbf{y}): \mathbb{R}^m \rightarrow \mathbb{R}^m$  and  $\mathbf{w} \mapsto S^{\mathcal{W}}(\mathbf{y}^*, \mathbf{w}): \mathbb{R}^p \rightarrow \mathbb{R}^p$  are increasing lower triangular maps, the latter for any  $\mathbf{y}^* \in \mathbb{R}^m$ . Recall that the reference density  $\eta$  is the standard normal on  $\mathbb{R}^d$ , and thus is a product of two standard normals,  $\eta_1(\mathbf{z}')\eta_2(\mathbf{z})$ , where  $\mathbf{z}' \in \mathbb{R}^m$  and  $\mathbf{z} \in \mathbb{R}^p$ . Using a KR map of the form (8), we can write the marginal density of  $\mathbf{Y}$  as  $\pi(\mathbf{y}) = (S_{\text{KR}}^{\mathcal{Y}})^{\#}\eta_1(\mathbf{y})$  and, more interestingly, the conditional density of  $\mathbf{W}$  as  $\pi(\mathbf{w}|\mathbf{y}) = S_{\text{KR}}^{\mathcal{W}}(\mathbf{y}, \cdot)^{\#}\eta_2(\mathbf{w})$ .

Each of the last  $p$  components of the KR rearrangement,  $S_{\text{KR}}^{\mathcal{W},k}(\mathbf{y}, \mathbf{w}_{\leq k})$ ,  $1 \leq k \leq p$ , can be estimated from a sample  $\{(\mathbf{Y}^i, \mathbf{W}^i)\}_{i=1}^n$  by minimizing

$$\hat{\mathcal{J}}_k(s) = \frac{1}{n} \sum_{i=1}^n \left( \frac{1}{2} s(\mathbf{Y}^i, \mathbf{W}_{\leq k}^i)^2 - \log |\partial_{m+k} s(\mathbf{Y}^i, \mathbf{W}_{\leq k}^i)| \right), \quad (9)$$

under the constraints  $\partial_{m+k} s(\mathbf{y}, \mathbf{w}_{\leq k}) > 0$ . This produces an estimator  $\hat{S}^{\mathcal{W}}$  of  $S_{\text{KR}}^{\mathcal{W}}$ , which in turn yields an estimator of the conditional density  $\pi(\mathbf{w}|\mathbf{y})$  as  $\hat{\pi}(\mathbf{w}|\mathbf{y}) := \hat{S}^{\mathcal{W}}(\mathbf{y}, \cdot)^{\#}\eta_2(\mathbf{w})$ . This property has been used to perform conditional density estimation (CDE) in [33, 41, 27, 15].

One important application of CDE is likelihood-free inference, where  $\mathbf{W}$  represents a parameter to be inferred and  $\mathbf{Y}$  are data whose conditional density  $\pi(\mathbf{y}|\mathbf{w})$  is computationally intractable or unavailable in closed form. This setting arises in many applications: inference in stochastic models, or inference in the presence of high-dimensional nuisance parameters or latent variables (including parameter inference for state-space models). Given a joint sample  $\{(\mathbf{Y}^i, \mathbf{W}^i)\}_{i=1}^n$ , we can estimate the map component  $S^{\mathcal{W}}$  and use it to simulate from the estimated conditional density  $\hat{\pi}(\mathbf{w}|\mathbf{y}^*)$  given any realization of the data  $\mathbf{y}^*$ , simply by sampling  $\mathbf{z}^i \sim \eta_2$  and solving the triangular system  $\hat{S}^{\mathcal{W}}(\mathbf{y}^*, \mathbf{w}^i) = \mathbf{z}^i$  for each  $\mathbf{w}^i$ . Note that we can simulate from or evaluate the estimated conditional densities for multiple realizations of the data  $\mathbf{y}^*$ , including values that are not present in the dataset  $\{\mathbf{Y}^i\}_{i=1}^n$ . Thus, learning a single map  $S^{\mathcal{W}}$  parameterized by  $\mathbf{y}$  is said to *amortize* the cost of conditional sampling over the data.

### 3. Representing and learning continuous monotone functions

In this section, we define a general representation for components of monotone triangular maps and present our main theoretical results. The essential idea, as mentioned in the introduction, is to express each component  $S_k$  as a nonlinear transformation  $S_k = \mathcal{R}_k(f)$  of a smooth function  $f$ , where  $\mathcal{R}_k$  (2) is an operator that enforces the monotonicity constraint by construction. We then identify the map component by solving the re-parameterized optimization problem

$$\min_{f \in V_k} \mathcal{L}_k(f) = \mathcal{J}_k(\mathcal{R}_k(f)), \quad (10)$$

where  $V_k$  is a linear space of functions in which we seek  $f$ . With this nonlinear transformation, we lose the convexity of the constrained problem  $\min_{\{s: \partial_k s > 0\}} \mathcal{J}_k(s)$ , but obtain an unconstrained minimization problem instead. It turns out that, with appropriate conditions on  $\mathcal{R}_k$ , the transformed optimization problem retains many desirable properties.

In Section 3.1, we discuss the construction of  $\mathcal{R}_k$  and motivate certain critical choices therein. Then, in Section 3.2, we analyze the regularity of the *composed* objective functional  $\mathcal{L}_k$  for a certain choice of the function space  $V_k$ . In Section 3.3 we discuss the existence and uniqueness of solutions to (10), and describe conditions under which solving (10) exactly recovers the KR rearrangement.

### 3.1. The rectification operator

Recalling (2), for any sufficiently smooth  $f: \mathbb{R}^k \rightarrow \mathbb{R}$  we let  $\mathcal{R}_k(f): \mathbb{R}^k \rightarrow \mathbb{R}$  be the function defined by

$$\mathcal{R}_k(f)(\mathbf{x}_{\leq k}) = f(\mathbf{x}_{< k}, 0) + \int_0^{x_k} g(\partial_k f(\mathbf{x}_{< k}, t)) dt,$$

where  $g: \mathbb{R} \rightarrow \mathbb{R}_{>0}$  is a positive function. We call the operator  $\mathcal{R}_k: f \mapsto \mathcal{R}_k(f)$  a *rectifier* because it transforms any  $f$  into a function which is increasing in the  $k$ -th variable, i.e.,  $\partial_k \mathcal{R}_k(f)(\mathbf{x}_{\leq k}) = g(\partial_k f(\mathbf{x}_{\leq k})) > 0$ . As a simple example, functions of the form  $f(\mathbf{x}_{\leq k}) = \alpha(\mathbf{x}_{< k}) + \beta(\mathbf{x}_{< k})x_k$  are transformed into  $\mathcal{R}_k(f)(\mathbf{x}_{\leq k}) = \alpha(\mathbf{x}_{< k}) + g(\beta(\mathbf{x}_{< k}))x_k$ . Figure 1 illustrates the application of the rectifier to a nonlinear univariate function  $f$ .

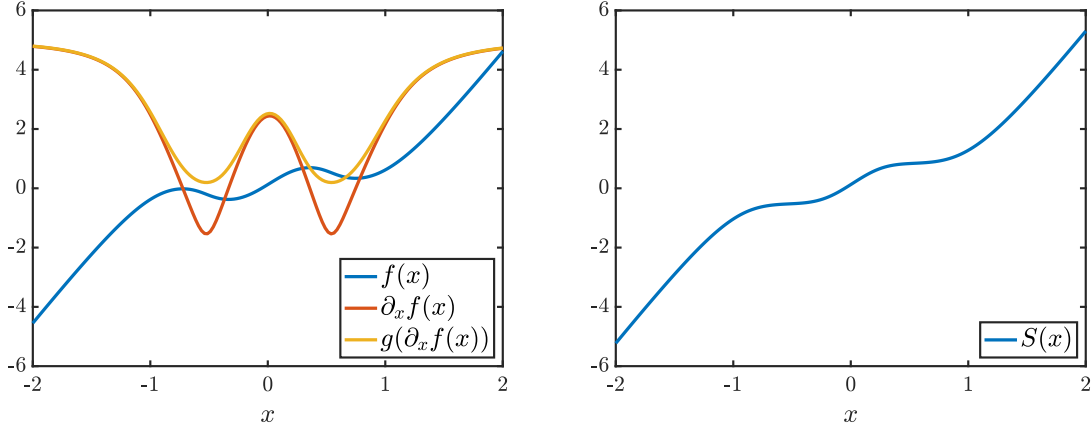


Fig 1: The rectifier (2) transforms the non-monotone function  $f$  into the monotone function  $S = \mathcal{R}(f)$ . Here we choose  $g(\cdot) = \log(1 + \exp(\cdot))$ .  $S$  is an increasing transport that pushes forward a one-dimensional mixture of Gaussians  $\pi(x) = 0.5\mathcal{N}(x; -1, 1) + 0.5\mathcal{N}(x; 1, 1)$  to the standard Gaussian reference density  $\eta$ .

The choice of the function  $g$  in (2) has a crucial impact on properties of the optimization problem (10). One possible choice, proposed in [22, 9], is the square function  $g(\xi) = \xi^2$ . While this choice permits closed-form computation of the integral in (3.1) when  $f$  is polynomial, it yields an optimization problem (10) which possesses many spurious local minima and is in fact non-smooth; see Figure 2 for an illustration. This can be explained in part by the fact that this  $g$  is not bijective.

Instead, one can choose  $g$  to be a bijective function from  $\mathbb{R}$  to  $\mathbb{R}_{>0}$ , such as the soft-plus function,

$$g(\xi) = \log(1 + \exp(\xi)), \quad (11)$$

whose inverse is  $g^{-1}(\xi) = \log(\exp(\xi) - 1)$ . Another example of a bijective function, considered in [69], is the shifted exponential linear unit (ELU),

$$g(\xi) = \begin{cases} \exp(\xi) & \xi < 0 \\ \xi + 1 & \xi \geq 0 \end{cases}, \quad (12)$$

whose inverse is  $g^{-1}(\xi) = \xi - 1$  if  $\xi \geq 1$  and  $g^{-1}(\xi) = \log(\xi)$  otherwise.

As a consequence of  $g$  being bijective, the inverse of the rectifier  $\mathcal{R}_k^{-1}(s)$  exists for any sufficiently smooth  $s: \mathbb{R}^k \rightarrow \mathbb{R}$  with  $\partial_k s(\mathbf{x}_{\leq k}) > 0$  and can be written as

$$\mathcal{R}_k^{-1}(s)(\mathbf{x}_{\leq k}) = s(\mathbf{x}_{< k}, 0) + \int_0^{x_k} g^{-1}(\partial_k s(\mathbf{x}_{< k}, t)) dt. \quad (13)$$



More importantly, the fact that  $g$  is invertible yields an objective function  $\mathcal{L}_k$  that is far better behaved than with  $g(\xi) = \xi^2$ ; see the numerical illustration in Figure 2, with both soft-plus and shifted ELU  $g$ . With these choices of  $g$ , we observe that  $\mathcal{L}_k = \mathcal{J}_k \circ \mathcal{R}_k$ , though non-convex in general, has no local minima and no saddle points. In the next sections, we will analyze these properties of the optimization problem (10) (e.g., smoothness, existence and uniqueness of solutions) and elucidate their precise dependence on  $g$ ,  $V_k$ , and  $\pi$ .

**Remark 1.** *In the finite-dimensional setting, it can easily be shown that composing a smooth and (strictly) convex objective function with a  $\mathcal{C}^1$ -diffeomorphic map preserves the (unique) global minima of the objective. Such diffeomorphic maps have been explored for accelerating optimization on finite-dimensional manifolds; see, e.g., [29]. Establishing that the operator  $\mathcal{R}_k$  in (3.1) is  $\mathcal{C}^1$ -diffeomorphic is, however, non-trivial. We will show in Section 3.3 that, under additional assumptions, our reparameterization yields an infinite-dimensional optimization problem where local minima are global minima.*

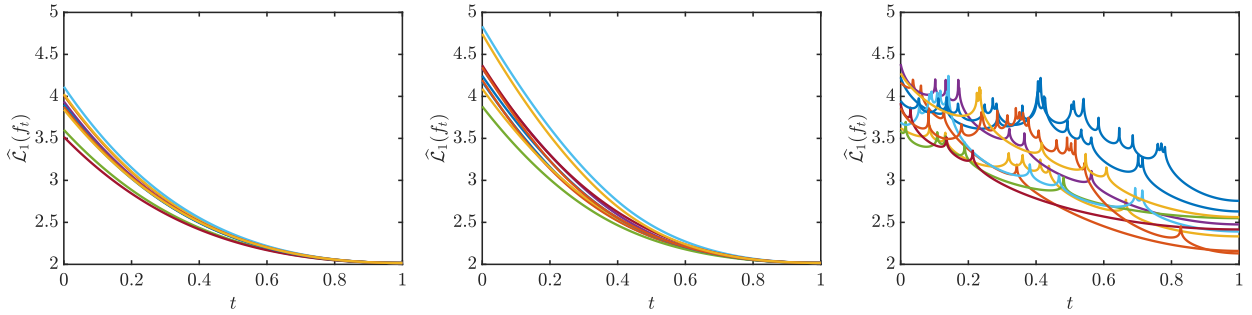


Fig 2: Objective function  $\hat{\mathcal{L}}_1 = \hat{\mathcal{J}}_1 \circ \mathcal{R}_1$  where the rectifier  $\mathcal{R}_1$  is defined using the soft-plus  $g$  (11) (left), the shifted ELU  $g$  (12) (middle), or the square function  $g(\xi) = \xi^2$  (right). Here,  $\pi(x) = 1/2\mathcal{N}(x; -2, 0.5) + 1/2\mathcal{N}(x; 2, 2)$  is a univariate Gaussian mixture, and we use  $n = 50$  to define  $\mathcal{J}_1$  with  $f_1$  represented using a linear combination of Hermite functions up to degree 10. The objective is evaluated along line segments that interpolate between random initial maps ( $t = 0$ ) and critical points resulting from a gradient-based optimization method ( $t = 1$ ). Observe that with bijective  $g$  (left and middle) the algorithm always arrives at the same optimal value, whereas with the square function  $g$  (right) the algorithm gets stuck in local minima and rarely attains the optimal value.

### 3.2. Smoothness of the optimization problem

In this section we give sufficient conditions on the function  $g$  and the target density  $\pi$  to guarantee that the objective function  $\mathcal{L}_k$  is smooth over an appropriate function space for  $f$ . We introduce this function space  $V_k$  and show continuity properties of the rectifier that hold for functions  $f \in V_k$ , before stating our main result in Proposition 4.

We begin by defining the weighted Sobolev space

$$V_k = \left\{ f: \mathbb{R}^k \rightarrow \mathbb{R} \text{ such that } \|f\|_{V_k}^2 := \int \left( |f(\mathbf{x})|^2 + |\partial_k f(\mathbf{x})|^2 \right) \eta_{\leq k}(\mathbf{x}) d\mathbf{x} < +\infty \right\}, \quad (14)$$

where  $\eta_{\leq k}$  is the standard normal density on  $\mathbb{R}^k$ . By [28, Theorem 1.11], this space is complete and thus is a Hilbert space. The space  $V_k$  has sufficient regularity for  $\partial_k f$  to exist, but also to

permit the pointwise evaluation  $f(\mathbf{x}_{<k}, 0)$ , as required in the definition (3.1) of the rectifier. This property is formalized by the following proposition, a trace theorem, which shows that for any  $f \in V_k$ , the function  $\mathbf{x}_{<k} \mapsto f(\mathbf{x}_{<k}, 0)$  is a function in  $L^2_{\eta_{<k}}$ , the  $\eta_{<k}$ -weighted space of square integrable functions.

**Proposition 2.** *There exists a constant  $C_T < \infty$  such that for any  $f \in V_k$*

$$\int f(\mathbf{x}_{<k}, 0)^2 \eta_{<k}(\mathbf{x}) d\mathbf{x} \leq C_T \|f\|_{V_k}^2. \quad (15)$$

*Proof.* See Appendix A.3. □

**Remark 2.** Notice that  $H^1_{\eta_{\leq k}} \subset V_k \subset L^2_{\eta_{\leq k}}$ , where  $L^2_{\eta_{\leq k}} = \{f: \mathbb{R}^k \rightarrow \mathbb{R} : \|f\|_{L^2_{\eta_{\leq k}}}^2 := \int f^2 d\eta_{\leq k} < \infty\}$  and  $H^1_{\eta_{\leq k}} = \{f: \mathbb{R}^k \rightarrow \mathbb{R} : \|f\|_{H^1_{\eta_{\leq k}}}^2 := \int f^2 + \|\nabla f\|^2 d\eta_{\leq k} < \infty\}$  are the standard weighted Sobolev spaces. Given that the standard normal density factorizes as  $\eta_{\leq k}(\mathbf{x}) = \prod_{i=1}^k \eta_i(x_i)$ , the space  $V_k$  admits the following tensor product structure

$$V_k = L^2_{\eta_1} \otimes \cdots \otimes L^2_{\eta_{k-1}} \otimes H^1_{\eta_k}, \quad (16)$$

and the norm  $\|\cdot\|_{V_k}$  is a product norm.<sup>2</sup> This tensor product structure will be used later in Section 4 to construct a numerical scheme for approximating the map components  $S_k$ .

The following proposition shows that, under mild assumptions on  $g$ , the rectifier  $\mathcal{R}_k$  is a Lipschitz continuous operator from  $V_k$  to itself. The proof relies on the Hardy inequality [37] and on Proposition 2.

**Proposition 3.** *Let  $g: \mathbb{R} \rightarrow \mathbb{R}_{>0}$  be a Lipschitz function, i.e., there exists a constant  $L < \infty$  so that*

$$|g(\xi) - g(\xi')| \leq L|\xi - \xi'|, \quad (17)$$

*holds for any  $\xi, \xi' \in \mathbb{R}$ . Then  $\mathcal{R}_k(f) \in V_k$  for any  $f \in V_k$ , where  $\mathcal{R}_k(f)$  is defined in (3.1). Furthermore there exists a constant  $C < \infty$  such that*

$$\|\mathcal{R}_k(f_1) - \mathcal{R}_k(f_2)\|_{V_k} \leq C \|f_1 - f_2\|_{V_k}. \quad (18)$$

*holds for any  $f_1, f_2 \in V_k$ .*

*Proof.* See Appendix A.4. □

The next result relies on Proposition 3 to show that the objective functional  $\mathcal{L}_k$  in (10), seen as a function from  $V_k$  to  $\mathbb{R}$ , is well-defined, continuous, and differentiable.

**Theorem 4.** *Let  $\pi$  be a probability density function on  $\mathbb{R}^d$  satisfying*

$$\pi(\mathbf{x}) \leq C_\pi \eta(\mathbf{x}), \quad (19)$$

*for all  $\mathbf{x} \in \mathbb{R}^d$  and a constant  $C_\pi < \infty$ , and let  $g: \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$  be an increasing function where*

$$|g(\xi) - g(\xi')| \leq L|\xi - \xi'|, \quad (20)$$

$$|\log \circ g(\xi) - \log \circ g(\xi')| \leq L|\xi - \xi'|, \quad (21)$$

---

<sup>2</sup>That is  $\|v_1 \otimes \cdots \otimes v_k\|_{V_k} = \|v_1\|_{L^2_{\eta_1}} \|v_2\|_{L^2_{\eta_2}} \cdots \|v_{k-1}\|_{L^2_{\eta_{k-1}}} \|v_k\|_{H^1_{\eta_k}}$  for any  $v_j \in L^2_{\eta_j}$  and  $v_k \in H^1_{\eta_k}$ .

holds for any  $\xi, \xi' \in \mathbb{R}$  and a constant  $L < \infty$ . Then

$$\mathcal{L}_k(f) := \mathcal{J}_k(\mathcal{R}_k(f)) < \infty,$$

for any  $f \in V_k$ , where  $\mathcal{J}_k(s) = \int (\frac{1}{2}s(\mathbf{x}_{\leq k})^2 - \log |\partial_k s(\mathbf{x}_{\leq k})|) \pi(\mathbf{x}) d\mathbf{x}$  as in (6) and where  $\mathcal{R}_k(f)$  is defined as in (3.1). Moreover, the function  $\mathcal{L}_k: V_k \rightarrow \mathbb{R}$  is continuous and differentiable at any  $f \in V_k$  and its gradient  $\nabla \mathcal{L}_k(f) \in V_k$  is given by

$$\begin{aligned} \langle \nabla \mathcal{L}_k(f), \varepsilon \rangle_{V_k} &= \int \mathcal{R}_k(f)(\mathbf{x}) \left( \varepsilon(\mathbf{x}_{< k}, 0) + \int_0^{x_k} g'(\partial_k f(\mathbf{x}_{< k}, t)) \partial_k \varepsilon(\mathbf{x}_{< k}, t) dt \right) \pi(\mathbf{x}) d\mathbf{x} \\ &\quad - \int (\log \circ g)'(\partial_k f(\mathbf{x})) \partial_k \varepsilon(\mathbf{x}) \pi(\mathbf{x}) d\mathbf{x}. \end{aligned} \quad (22)$$

for any  $\varepsilon \in V_k$ , where  $\langle \cdot, \cdot \rangle_{V_k}$  is the scalar product in  $V_k$ .

*Proof.* See Appendix A.5. □

It is useful to discuss some implications of Theorem 4. A smooth objective function permits us to use deterministic first or second-order optimization algorithms. Furthermore, in Section 4, we exploit the gradients of  $\mathcal{J}_k$  in (22) to construct an adaptive polynomial (or wavelet) basis for  $V_k$ . We also note that many functions  $g$  satisfy the conditions (20) and (21) in the proposition. For example, these include the soft-plus (11) and shifted ELU (12) functions considered in Figure 2.

**Remark 3** (Assumption (19) implies Gaussian tails). *The condition (19) implies that  $\mathbb{P}_{\mathbf{X} \sim \pi}(\|\mathbf{X}\|_2 > t) \leq C_\pi \mathbb{P}_{\mathbf{Z} \sim \eta}(\|\mathbf{Z}\|_2 > t)$  holds for any  $t \geq 0$ , and hence that  $\pi$  has sub-Gaussian tails [65]. The reverse, however, does not hold in general. For instance, a compactly supported random variable with density  $\pi(x) \propto 1/\sqrt{|x|} \mathbb{1}_{\{x \in [-1, 1]\}}$  is sub-Gaussian, but the density is unbounded at  $x = 0$  and thus does not satisfy (19). We also note that (19) can be relaxed with the (less interpretable) assumption that there exist a constant  $C_\pi < \infty$  and diagonal scalings  $d_k > 0$  so that  $\pi(\mathbf{x}) \leq C_\pi \prod_{k=1}^d \eta_k(d_k x_k)$  for all  $\mathbf{x} \in \mathbb{R}^d$ . In practice, we can always normalize the marginals of  $\pi$  to match the mean and variance of a standard Gaussian reference. Thus, without loss of generality, we will use the assumption above for the remainder of this article.*

**Remark 4.** *Under additional assumptions on  $g$ , the gradient  $f \mapsto \nabla \mathcal{L}_k(f)$  is locally Lipschitz from  $\bar{V}_k$  to  $V_k$ , where  $\bar{V}_k = \{f \in V_k, \partial_k f \in L^\infty\}$  is the space endowed with the norm  $\|f\|_{\bar{V}_k} = \|f\|_{V_k} + \|\partial_k f\|_{L^\infty}$ . More specifically, if  $g'$  and  $(\log \circ g)'$  are Lipschitz functions, then there exists a constant  $M < \infty$  such that*

$$\|\nabla \mathcal{L}_k(f_1) - \nabla \mathcal{L}_k(f_2)\|_{V_k} \leq M(1 + \|f_2\|_{V_k}) \|f_1 - f_2\|_{\bar{V}_k}, \quad (23)$$

holds for any  $f_1, f_2 \in \bar{V}_k$ . See the derivation of this result in Appendix A.6. Such local Lipschitz regularity is useful to analyze the convergence of backtracking gradient descent procedures, i.e., gradient descent with an inexact line search such as Armijo's rule, to a stationary point; see Theorem 2.1 in [63] and Proposition 2.1.1 in [4]. We leave the extension of such optimization guarantees for solving  $\min_{f \in V_k} \mathcal{L}_k(f)$  to future work.

### 3.3. Existence and uniqueness of solutions

In this section, we show that the optimization problem (10) does not admit any spurious local minima, meaning that local minimizers are in fact global minimizers. We also show that problem (10) admits a unique global minimizer which permits us to recover the KR rearrangement. To prove these results, we will need the following proposition, which provides conditions ensuring that the inverse rectifier  $\mathcal{R}_k^{-1}$  is continuous.

**Proposition 5.** *Let  $g$  be a bijective function from  $\mathbb{R}$  to  $\mathbb{R}_{>0}$  such that for any  $c > 0$  there exists a constant  $L_c < \infty$  so that*

$$|g^{-1}(\xi) - g^{-1}(\xi')| \leq L_c |\xi - \xi'|, \quad (24)$$

*holds for any  $\xi, \xi' \geq c$ . Then, for any  $s \in V_k$  such that  $\text{ess inf } \partial_k s > 0$ , we have  $\mathcal{R}_k^{-1}(s) \in V_k$  and  $\text{ess inf } \partial_k \mathcal{R}_k^{-1}(s) > -\infty$ . Furthermore for any  $c > 0$ , there exists a constant  $C_c < \infty$  such that*

$$\|\mathcal{R}_k^{-1}(s_1) - \mathcal{R}_k^{-1}(s_2)\|_{V_k} \leq C_c \|s_1 - s_2\|_{V_k}, \quad (25)$$

*holds for any  $s_1, s_2 \in V_k$  such that  $\text{ess inf } \partial_k s_i \geq c$ .*

*Proof.* See Appendix A.8. □

Note that the softplus function (11) and the shifted exponential linear unit (12) satisfy (24) with  $L_c = (1 - e^{-c})^{-1}$  and  $L_c = \max\{1/c, 1\}$ , respectively.

### 3.3.1. Local minima are global minima

The following proposition shows that, under certain conditions on the function  $g$ , the image set  $\mathcal{R}_k(V_k) = \{\mathcal{R}_k(f) : f \in V_k\}$  is convex.

**Proposition 6.** *Let  $g: \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$  be an increasing function such that  $\xi \mapsto (g^{-1}(\xi))^2$  is a convex function. Let  $\mathcal{R}_k(f)$  be defined as in (3.1) and  $V_k$  as in (14). Then  $\mathcal{R}_k(V_k) = \{\mathcal{R}_k(f) : f \in V_k\}$  is convex.*

*Proof.* See Appendix A.7. □

When  $g$  is the softplus function (11) or the exponential linear unit (12), the mapping  $\xi \mapsto (g^{-1}(\xi))^2$  is indeed convex. Interestingly, however, the exponential function  $g(\xi) = \exp(\xi)$ —which has been used to parameterize monotone maps in [33, 59] and for monotone regression in [46, 55]—does not satisfy this property, and there is no guarantee for  $\mathcal{R}_k(V_k)$  to be convex in this case.

An important consequence of Proposition 6 is that the constrained problem  $\min_{s \in \mathcal{R}_k(V_k)} \mathcal{L}_k(s)$  remains convex. Hence, this optimization problem has a unique global minimizer by the strict convexity of  $\mathcal{L}_k$  (see Appendix A.2). The next result shows that the local minima of the unconstrained problem (10) are in fact global minima.

**Theorem 7.** *Under the assumptions of Propositions 5 and 6, let  $f^* \in V_k$  be a local minimizer of  $f \mapsto \mathcal{J}_k(\mathcal{R}_k(f))$ , meaning that there exists  $\rho > 0$  such that*

$$\mathcal{J}_k(\mathcal{R}_k(f^*)) \leq \mathcal{J}_k(\mathcal{R}_k(f)), \quad \forall f \in V_k \text{ such that } \|f - f^*\|_{V_k} \leq \rho. \quad (26)$$

*If  $\text{ess inf } \partial_k f^* > -\infty$ , then  $f^*$  is a global minimizer, meaning that*

$$\mathcal{J}_k(\mathcal{R}_k(f^*)) \leq \mathcal{J}_k(\mathcal{R}_k(f)), \quad \forall f \in V_k.$$

*Proof.* Let  $f \in V_k$ . For any  $0 < t < 1$  we let  $s_t = t\mathcal{R}_k(f) + (1-t)\mathcal{R}_k(f^*)$ . By Proposition 6,  $\mathcal{R}_k(V_k)$  is convex so that  $s_t \in \mathcal{R}_k(V_k)$ . By convexity of  $\mathcal{J}_k$  (c.f. Appendix A.2), we can write  $\mathcal{J}_k(s_t) \leq t\mathcal{J}_k(\mathcal{R}_k(f)) + (1-t)\mathcal{J}_k(\mathcal{R}_k(f^*))$ , or equivalently

$$\mathcal{J}_k(s_t) - \mathcal{J}_k(\mathcal{R}_k(f^*)) \leq t(\mathcal{J}_k(\mathcal{R}_k(f)) - \mathcal{J}_k(\mathcal{R}_k(f^*))). \quad (27)$$

Next we show that there exists a sufficiently small  $t > 0$  such that the above left-hand side is positive. As a consequence, the right hand side of (27) will be positive, which will conclude the proof.

Let  $b = \text{ess inf } \partial_k f^*$  and  $c = g(b)/2 > 0$ . For any  $t \leq 1/2$ , as have  $\partial_k s_t = tg(\partial_k f) + (1-t)g(\partial_k f^*) \geq 1/2g(\partial_k f^*)$  so that  $\text{ess inf } \partial_k s_t \geq c$ . In addition, we have  $\text{ess inf } \partial_k \mathcal{R}_k(f^*) = \text{ess inf } g(\partial_k f^*) \geq c$ . Thus, by Proposition 5, there exists a constant  $C_c < \infty$  such that

$$\begin{aligned} \|\mathcal{R}_k^{-1}(s_t) - f^*\|_{V_k} &= \|\mathcal{R}_k^{-1}(s_t) - \mathcal{R}_k^{-1}(\mathcal{R}_k(f^*))\|_{V_k} \\ &\leq C_c \|s_t - \mathcal{R}_k(f^*)\|_{V_k} \\ &= tC_c \|\mathcal{R}_k(f) - \mathcal{R}_k(f^*)\|_{V_k}. \end{aligned}$$

By letting  $t = \rho/(C_c \|\mathcal{R}_k(f) - \mathcal{R}_k(f^*)\|_{V_k})$ , we have  $\|\mathcal{R}_k^{-1}(s_t) - f^*\|_{V_k} \leq \rho$ . Therefore, setting  $f = \mathcal{R}_k^{-1}(s_t)$  in (26) ensures that  $0 \leq \mathcal{J}_k(s_t) - \mathcal{J}_k(\mathcal{R}_k(f^*))$ .  $\square$

Let us remark that Theorem 7 holds for general target densities  $\pi$  without the bound assumed in Theorem 4. Applying this result, however, depends on the existence of a local minimizer in the function space  $V_k$ , which must be established for a given target density. The following section studies a relatively broad class of distributions where the  $k$ -th component of the KR rearrangement is in  $\mathcal{R}_k(V_k)$ , thereby providing a candidate solution to apply Theorem 7.

### 3.3.2. Uniqueness of the global minimizer and recovery of the KR rearrangement

As discussed in Section 2, the Knothe–Rosenblatt rearrangement  $S_{\text{KR}}$  is the unique lower triangular and monotone map such that  $\mathcal{D}_{\text{KL}}(\pi \| S_{\text{KR}}^\# \eta) = 0$ ; see [7]. The decomposition (5) of the KL divergence  $\mathcal{D}_{\text{KL}}(\pi \| S^\# \eta)$  thus permits us to write

$$\mathcal{J}_k(S_{\text{KR},k}) \leq \mathcal{J}_k(\mathcal{R}_k(f)), \quad (28)$$

for any  $f \in V_k$  and for any  $1 \leq k \leq d$ . Indeed, by letting  $S$  be the map such that  $S_k = \mathcal{R}_k(f)$  and  $S_i = S_{\text{KR},i}$  for  $i \neq k$ , (5) yields (28). Thus, if there exists a function  $f_{\text{KR},k} \in V_k$  such that  $S_{\text{KR},k} = \mathcal{R}_k(f_{\text{KR},k})$ , then  $f_{\text{KR},k}$  is a global minimizer of  $f \mapsto \mathcal{J}_k(\mathcal{R}_k(f))$  over  $f \in V_k$ . To show this, we first need the following intermediate result.

**Proposition 8.** *Let  $\pi$  be a probability density function on  $\mathbb{R}^d$  such that*

$$c\eta(\mathbf{x}) \leq \pi(\mathbf{x}) \leq C\eta(\mathbf{x}), \quad (29)$$

*for all  $\mathbf{x} \in \mathbb{R}^d$  with some constants  $0 < c \leq C < \infty$ . Then, for all  $\mathbf{x}_{<k} \in \mathbb{R}^{k-1}$  and  $k = 1, \dots, d$ ,  $S_{\text{KR},k}(\mathbf{x}_{<k}, x_k) = \mathcal{O}(x_k)$  and  $\partial_k S_{\text{KR},k}(\mathbf{x}_{<k}, x_k) = \mathcal{O}(1)$  as  $|x_k| \rightarrow \infty$ . Furthermore, we have  $S_{\text{KR},k} \in V_k$  and  $\text{ess inf } \partial_k S_{\text{KR},k} > 0$  for all  $k = 1, \dots, d$ .*

*Proof.* See Appendix A.9.  $\square$

Let us comment on the assumption (29). While the results in Section 3.2 assume only an upper bound on the joint density  $\pi$ , here we have both an upper and a lower bound. These bounds together imply that the  $k$ -th marginal conditional density, for  $1 \leq k \leq d$ , satisfies  $(c/C)\eta_k(x_k) \leq \pi(x_k | \mathbf{x}_{<k}) \leq (C/c)\eta_k(x_k)$  for all  $\mathbf{x}_{<k} \in \mathbb{R}^k$ . This means that the marginal conditionals of the target density have Gaussian tails, rather than potentially being lighter than Gaussian (as in assumption (19)). This condition guarantees that the asymptotic behavior of each map component  $S_{\text{KR},k}$  in its last variable (as  $|x_k| \rightarrow \infty$ ) is affine.

We now combine our previous results to show that under the assumption (29) on the target density  $\pi$ , we have  $f_{\text{KR},k} \in V_k$  and the existence of a unique solution to  $\min_{f \in V_k} \mathcal{J}_k(\mathcal{R}_k(f))$ .

**Corollary 9.** *Under the assumptions on  $\pi$  in Proposition 8, let  $g$  be a Lipschitz bijection from  $\mathbb{R}$  to  $\mathbb{R}_{\geq 0}$  such that  $\xi \mapsto (g^{-1}(\xi))^2$  is convex and such that, for any  $c > 0$  there exists a constant  $L_c < \infty$  such that  $|g^{-1}(\xi) - g^{-1}(\xi')| \leq L_c |\xi - \xi'|$  for any  $\xi, \xi' \geq c$ . Then, for any  $1 \leq k \leq d$ , there exists a unique function  $f_{KR,k} \in V_k$  that satisfies  $\text{ess inf } \partial_k f_{KR,k} > -\infty$  such that  $\mathcal{R}_k(f_{KR,k})$  is the  $k$ -th component of the KR rearrangement  $S_{KR}$ . As a consequence,  $\min_{f \in V_k} \mathcal{J}_k(\mathcal{R}_k(f))$  admits a unique solution and*

$$\mathcal{J}_k(\mathcal{R}_k(f_{KR,k})) = \min_{f \in V_k} \mathcal{J}_k(\mathcal{R}_k(f)).$$

*Proof.* Combining Proposition 8 and Proposition 5, there exists a function  $f_{KR,k} \in V_k$  that satisfies  $\text{ess inf } \partial_k f_{KR,k} > -\infty$  such that  $\mathcal{R}_k(f_{KR,k}) = S_{KR,k}$ . Then from Theorem 7 and inequality (28), we have that  $f_{KR,k}$  is a global minimizer of  $f \mapsto \mathcal{J}_k(\mathcal{R}_k(f))$  over  $f \in V_k$ . From the uniqueness of the KR rearrangement  $S_{KR,k}$  among monotone triangular maps and the injectivity of  $\mathcal{R}_k^{-1}$ ,  $f_{KR,k} = \mathcal{R}_k^{-1}(S_{KR,k})$  is unique.  $\square$

#### 4. Adaptive parameterization of transport maps

Given an i.i.d. sample  $\{\mathbf{X}^i\}_{i=1}^n \sim \pi$ , we now propose an adaptive algorithm to build  $f \in V_k$  that minimizes the empirical objective  $\hat{\mathcal{L}}_k := \hat{\mathcal{J}}_k \circ \mathcal{R}_k(f)$ , where  $\hat{\mathcal{J}}_k$  is as in (7). Doing so for each  $k = 1, \dots, d$  yields a monotone triangular map, as described in Section 2. For each  $k$ , we represent  $f \in V_k$  with an  $m$ -term expansion

$$f(\mathbf{x}_{\leq k}) = \sum_{\alpha \in \Lambda} c_\alpha \psi_\alpha(\mathbf{x}_{\leq k}), \quad (30)$$

where  $\Lambda$  is a set of multi-indices  $\alpha = (\alpha_1, \dots, \alpha_k) \in \mathbb{N}_0^k$  with  $\#\Lambda = m$ ,  $c_\alpha \in \mathbb{R}$  are coefficients, and  $\psi_\alpha: \mathbb{R}^k \rightarrow \mathbb{R}$  are basis functions for  $V_k$ , constructed as products of univariate functions,

$$\psi_\alpha(\mathbf{x}_{\leq k}) = \prod_{j=1}^k \psi_{\alpha_j}^j(x_j).$$

Here,  $\{\psi_\alpha^j\}_{\alpha \in \mathbb{N}_0}$  is chosen to be a basis of  $L_{\eta_j}^2$  if  $j < k$  and of  $H_{\eta_k}^1$  if  $j = k$ . Because  $V_k$  possesses a tensor product structure (16) and because its norm  $\|\cdot\|_{V_k}$  is a product norm, the basis  $\{\psi_\alpha\}_{\alpha \in \mathbb{N}_0^k}$  is orthonormal if  $\{\psi_\alpha^j\}_{\alpha \in \mathbb{N}_0}$  is an orthonormal basis for all  $j \leq k$ . Since  $f$  depends linearly on the coefficients  $c_\alpha$ , the smoothness properties of the objective  $\mathcal{L}_k$  transfer to the objective function parameterized by the coefficients  $c_\alpha$ ,  $\alpha \in \Lambda$ .

Sections 4.1 and 4.2 describe two different choices of basis: polynomials and wavelets, respectively. Then, in Section 4.3, we present a greedy algorithm for building the multi-index set  $\Lambda$ , applicable to any such hierarchical basis. In addition, we propose a cross validation procedure to determine when to stop enriching  $\Lambda$  in order to avoid over-fitting to the finite sample.

##### 4.1. Polynomial basis

Probabilists' Hermite polynomials  $\{\varphi_\alpha\}_{\alpha \in \mathbb{N}}$  form an orthogonal basis for  $L_\eta^2$  where  $\eta$  is the univariate standard Gaussian density.  $\varphi_\alpha$  is a polynomial of degree  $\alpha \geq 0$  defined as

$$\varphi_\alpha(x) = (-1)^\alpha \exp(x^2/2) \frac{d^\alpha}{dx^\alpha} \exp(-x^2/2).$$



Furthermore, we have  $\langle \varphi_\alpha, \varphi_\beta \rangle_{L_\eta^2} = \alpha! \delta_{\alpha,\beta}$  so that  $\{\varphi_\alpha / \sqrt{\alpha!}\}_{\alpha \in \mathbb{N}_0}$  forms an orthonormal basis for  $L_\eta^2$ . Similarly,  $\langle \varphi_\alpha, \varphi_\beta \rangle_{H_\eta^1} = (\alpha + 1)! \delta_{\alpha,\beta}$  so that  $\{\varphi_\alpha / \sqrt{(\alpha + 1)!}\}_{\alpha \in \mathbb{N}_0}$  forms an orthonormal basis for  $H_\eta^1$ ; see [53, Proposition 1.3]. We can thus tensorize these basis functions to obtain an orthonormal basis for the anisotropic space  $V_k$  (16).

In practice, we modify the univariate Hermite polynomials to be linear outside of an arbitrary compact domain. Following [44], we let

$$\psi_{\alpha_j}(x_j) = \frac{1}{\sqrt{Z_{\alpha_j}}} \begin{cases} \varphi_{\alpha_j}(x_j^a) + \varphi'_{\alpha_j}(x_j^a)(x_j - x_j^a) & \text{if } x_j < x_j^a \\ \varphi_{\alpha_j}(x_j) & \text{if } x_j^a \leq x_j \leq x_j^b \\ \varphi_{\alpha_j}(x_j^b) + \varphi'_{\alpha_j}(x_j^b)(x_j - x_j^b) & \text{if } x_j > x_j^b \end{cases} \quad (31)$$

where  $Z_{\alpha_j} = \alpha_j!$  for  $j = 1, \dots, k-1$  and  $Z_{\alpha_j} = (\alpha_j + 1)!$  for  $j = k$ . In our numerical experiments, we set  $x_j^a$  and  $x_j^b$  to be the 0.01 and 0.99 (empirical) quantiles, respectively, of the marginal sample  $\{X_j^i\}_{i=1}^n$ . The basis  $\{\psi_{\alpha_j}\}_{\alpha_j \in \mathbb{N}_0}$ , albeit not exactly orthonormal, is close to being orthonormal for sufficiently small  $x_j^a$  and large  $x_j^b$ . This is important for the numerical stability of the algorithm presented in the following section.

The basis  $\{\psi_\alpha\}_\alpha$  generated by the piecewise polynomials (31) has numerical and structural advantages over the Hermite polynomials. First, these functions avoid the fast growth rate of standard polynomials as  $\|\mathbf{x}\| \rightarrow \infty$ , providing a more stable extrapolation of the map to the tails of  $\pi$ , where there are typically no samples available to otherwise constrain its growth. Second, a map that reverts to linearity far from the origin yields densities in the class considered in Subsection 3.3.2; see condition (29). Indeed, each map component  $S_k = \mathcal{R}_k(f)$  with  $f \in \text{span}\{\psi_\alpha : \alpha \in \Lambda\}$  behaves linearly as  $|x_k| \rightarrow \infty$ , so that the pullback density  $S^\# \eta$  has Gaussian tails like the reference  $\eta$ .

#### 4.2. Wavelet basis

Wavelets are a popular approach for approximating functions that are not uniformly smooth [12, 66]. In particular, wavelet techniques define a multi-resolution approximation to a function that can better capture local features than, for instance, global polynomials. Given a compactly supported function  $\psi : \mathbb{R} \rightarrow \mathbb{R}$  called the mother wavelet, and indices  $(l, q) \in \mathbb{Z}^2$ , the function

$$\psi_{(l,q)} : x \mapsto 2^{l/2} \psi(2^l x - q),$$

is the  $q$ th wavelet at the level  $l$ . Common choices for the mother wavelet include the Haar, the Daubechies, and the Meyer wavelets.

In our numerical experiments, we use the continuous Mexican hat mother wavelet

$$\psi(x) = \frac{2}{\sqrt{3\sigma\pi^{1/4}}} (1 - (x/\sigma)^2) \exp(-x^2/(2\sigma^2)), \quad (32)$$

with scale parameter  $\sigma = 1$  [32], which is obtained from the second derivative of the univariate Gaussian density. The Mexican hat function is also commonly referred to as the Ricker wavelet in the geophysics community. We treat this function as having essentially compact support on  $[-6, 6]$ , up to numerical precision.

Tensorizing a univariate wavelet in the  $k$ th coordinate direction with modified Hermite polynomials in the first  $k-1$  directions yields an expansion of the form in (30), where the indices in  $\alpha = (\alpha_1, \dots, \alpha_k)$  are tuples of the form

$$\alpha_1, \dots, \alpha_{k-1} \in \mathbb{N}_0, \quad \alpha_k = (l_k, q_k) \in \mathbb{Z}^2. \quad (33)$$



In practice, we consider a truncated set of wavelets by first rescaling the mother wavelet  $\psi$  to have the same support as the interval between the 0.01 and 0.99 empirical quantiles of the marginal sample  $\{X_j^i\}_{i=1}^n$ . Then, we constrain  $l_j \geq 0$  and consider  $l_j = 0$  to be the coarsest level of the approximation. Next, for any  $l_j$ , the translation parameter  $q_j$  is bounded as  $0 \leq q_j \leq 2^{l_j} - 1$ . Thus, we can consider

$$\alpha_1, \dots, \alpha_{k-1} \in \mathbb{N}_0, \quad \alpha_k = (l_k, q_k) \in \mathbb{N}_0^2.$$

rather than (33).

Lastly, we note that any function  $f$  of the form (30) expanded with compactly supported wavelets will decay to 0 as  $|x_k| \rightarrow \infty$ . Thus, the map  $S_k = \mathcal{R}_k(f)$  will have asymptotic linear growth as  $|x_k| \rightarrow \infty$ , similarly to the modified Hermite polynomial expansions.

### 4.3. Adaptive transport map algorithm

Now we propose a greedy method for constructing the multi-index set  $\Lambda = \Lambda_t$  in (30). For simplicity, we consider only the case of single indices  $\alpha_j \in \mathbb{N}_0$ . The extension to the case  $\alpha_j \in \mathbb{N}_0^2$  (as in the wavelet construction above) is described in Appendix B. At each greedy iteration  $t$ , we add one multi-index  $\alpha_t^*$  to  $\Lambda_t$ . Starting with  $\Lambda_0 = \emptyset$  we let

$$\Lambda_{t+1} = \Lambda_t \cup \{\alpha_t^*\}, \quad (34)$$

where  $\alpha_t^* \notin \Lambda_t$  is a multi-index that yields the *best* improvement of  $\widehat{\mathcal{L}}_k$  in a sense to be defined below. Borrowing ideas from [34, 35], we constrain the sets  $\{\Lambda_t\}_{t \geq 0}$  to be *downward closed* [11, 13], meaning that they satisfy the property

$$\alpha \in \Lambda_t \text{ and } \alpha' \leq \alpha \Rightarrow \alpha' \in \Lambda_t, \quad (35)$$

where  $\alpha' \leq \alpha$  denotes  $\alpha'_i \leq \alpha_i$  for all  $1 \leq i \leq k$ . Intuitively, (35) means that  $\Lambda_t$  has a pyramidal shape containing no holes. Downward-closed sets are known to preserve good approximation properties and permit a tractable construction of  $\Lambda_t$ ; see [13]. Indeed,  $\Lambda_{t+1}$  remains downward-closed if and only if the multi-index  $\alpha_t^*$  is selected from the so-called *reduced margin* of  $\Lambda_t$ , defined by

$$\Lambda_t^{\text{RM}} = \{\alpha \notin \Lambda_t : \alpha - \mathbf{e}_i \in \Lambda_t \text{ for all } 1 \leq i \leq k \text{ with } \alpha_i \neq 0\},$$

where  $\mathbf{e}_i$  denotes the  $i$ -th canonical vector of  $\mathbb{N}^k$ . The reduced margin is a subset of the margin set  $\Lambda_t^{\text{M}}$  (i.e., multi-indices  $\alpha \notin \Lambda_t$  such that  $\exists i > 0$  where  $\alpha - \mathbf{e}_i \in \Lambda_t$ ); see Figure 3. The size of the reduced margin typically grows more slowly with respect to the dimension  $k$  than the margin itself. For instance, if  $\Lambda_t$  contains all multi-indices in a hypercube of width  $p$  in dimension  $k$ , the margin has cardinality  $(p+1)^k - p^k$ , while the reduced margin has cardinality  $k$ .

Denoting by  $f_t$  the minimizer of  $f \mapsto \widehat{\mathcal{L}}_k(f)$  over  $f \in \text{span}\{\psi_\alpha, \alpha \in \Lambda_t\}$ , we select  $\alpha_t^*$  in the reduced margin  $\Lambda_t^{\text{RM}}$  with the following heuristic

$$\alpha_t^* \in \arg \max_{\alpha \in \Lambda_t^{\text{RM}}} |\nabla_\alpha \widehat{\mathcal{L}}_k(f_t)|. \quad (36)$$

Here, the notation  $\nabla_\alpha \widehat{\mathcal{L}}_k(f_t)$  denotes the derivative of  $c_\alpha \mapsto \widehat{\mathcal{L}}_k(f_t + c_\alpha \psi_\alpha)$  evaluated at  $c_\alpha = 0$ . In other words, we select the multi-index  $\alpha_t^*$  by choosing the largest functional derivative of  $\widehat{\mathcal{L}}_k$  at  $f_t$  with respect to the basis functions not contained in the current expansion. This greedy procedure for learning each map component is presented in detail in Algorithm 1.

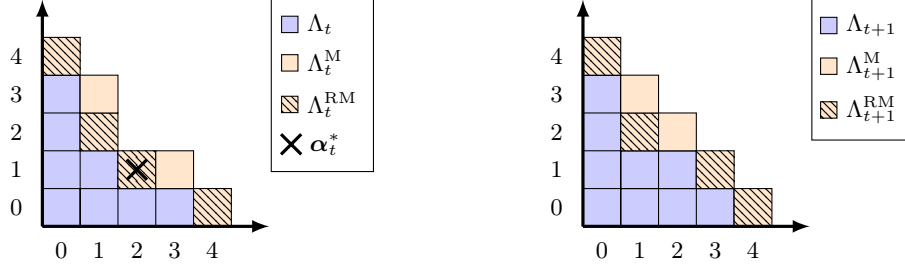


Fig 3: A  $k = 2$  dimensional downward-closed active set of multi-indices  $\Lambda_t$  with its margin  $\Lambda_t^M$  and reduced margin  $\Lambda_t^{RM}$ . The margin and reduced margins are plotted before (*left*) and after (*right*) adding to  $\alpha_t^* = (2, 1)$ , that is denoted with a cross, to  $\Lambda_t$ .

---

**Algorithm 1** Estimate map component  $S_k$ 


---

- 1: **Input:** training sample  $\{\mathbf{X}_{1:k}^i\}_{i=1}^n$ , maximum cardinality  $m$  for  $\Lambda_t$
  - 2: Initialize  $\Lambda_0 = \emptyset$ ,  $f_0 = 0$
  - 3: **for**  $t = 0, \dots, m-1$  **do**
  - 4:   Construct the reduced margin:  $\Lambda_t^{RM}$
  - 5:   Select the new multi-index:  $\alpha_t^* \in \arg \max_{\alpha \in \Lambda_t^{RM}} |\nabla_{\alpha} \hat{\mathcal{L}}_k(f_t)|$
  - 6:   Update the active set:  $\Lambda_{t+1} = \Lambda_t \cup \{\alpha_t^*\}$
  - 7:   Update the approximation:  $f_{t+1} = \arg \min_{f \in \text{span}\{\psi_{\alpha} : \alpha \in \Lambda_{t+1}\}} \hat{\mathcal{L}}_k(f)$
  - 8: **end for**
  - 9: **Output:**  $\hat{S}_k = \mathcal{R}_k(f_m)$
- 

**Remark 5.** If the objective function  $\hat{\mathcal{L}}_k$  is the linear least-squares loss, then Algorithm 1 corresponds to orthogonal matching pursuit for sparse linear regression with normalized basis functions. An alternative heuristic for selecting  $\alpha_t^*$  that requires second-order information from the objective is  $\alpha_t^* \in \arg \max_{\alpha \in \Lambda_t^{RM}} |\nabla_{\alpha} \hat{\mathcal{L}}_k(f_t)|^2 / |\nabla_{\alpha}^2 \hat{\mathcal{L}}_k(f_t)|$ . We found this criterion to perform similarly to (36) for the polynomial basis, but it led to faster convergence of the objective when working with the wavelet basis.

**Remark 6.** A potential drawback of Algorithm 1 is that the greedy enrichment procedure cannot “see” behind the reduced margin. For instance, if a relevant multi-index is located far beyond the reduced margin and the gradient  $\nabla_{\alpha} \hat{\mathcal{L}}_k(f_t)$  vanishes for all  $\alpha \in \Lambda_t^{RM}$ , then the algorithm will be stuck. [35] suggested a safeguard mechanism to avoid this behaviour: arbitrarily activate the most ancient index from the reduced margin every  $t_{sg} \in \mathbb{N}$  iterations. This modification, however, was not needed in our numerical experiments.

**Remark 7.** As pointed out in [13, 34, 35], adding multiple multi-indices at each greedy iteration could also yield better performance compared to adding only one multi-index at a time. The so-called “bulk-chasing” procedure identifies a subset  $\lambda_t^* \subset \Lambda_t^{RM}$  of multi-indices that capture a fraction of the  $L_2$ -norm of the gradient  $\nabla_{\alpha} \hat{\mathcal{L}}_k(f_t)$  along the reduced margin.

To determine the maximal cardinality of  $\Lambda_t$  (i.e., when to stop adaptation) we use  $\nu$ -fold cross-validation as in [68, 5]. For each fold, we run Algorithm 1 with  $\nu - 1$  folds of the training data for  $m = n^{\frac{\nu-1}{\nu}}$  iterations, and evaluate the objective function in (7) on the held-out test set. We then select the number of terms  $m^* \leq n$  in the expansion (30) that minimizes the test error averaged over the  $\nu$  folds. The full sample is then used to run Algorithm 1 for  $m^*$  iterations. The complete procedure for learning each map component  $S_k$ , with cross-validation, is called the Adaptive Transport Map (ATM) algorithm. This procedure is detailed in Algorithm 2. In practice, we also stop the training on

each fold early if the log-likelihood of the test set does not continue to increase for more than 20 iterations.

The cross-validation procedure produces an approximation for  $f$  with an expansion whose cardinality is *adapted* to the size of the training sample  $n$ . With a larger sample, we observe that the cross-validation procedure reliably adds more parameters, when needed, to reduce the bias in the approximation while controlling the variance of the estimated map parameters. Thus, we consider ATM a *semi-parametric* approach for approximating monotone triangular transport maps. Adaptation of the map complexity to the sample size  $n$  will be demonstrated in the numerical results to follow.

---

**Algorithm 2** ATM algorithm for learning map component  $S_k$ 


---

- 1: **Input:** Training sample  $\chi = \{\mathbf{X}_{1:k}^i\}_{i=1}^n$ , number of folds  $\nu$
  - 2: Partition data into  $\nu$  folds of equal size
  - 3: **for**  $j = 1, \dots, \nu$  **do**
  - 4:   Partition  $\chi$ :  $\chi_j^{\text{test}}$  is the  $j$ th subset of  $\chi$  and  $\chi_j^{\text{train}} = \chi \setminus \chi_j^{\text{test}}$
  - 5:   Estimate  $f_t$  using Algorithm 1 for  $m = |\chi_j^{\text{train}}|$  iterations with sample  $\chi_j^{\text{train}}$
  - 6:   Store iterates  $\widehat{S}_k^{j,1} := \mathcal{R}_k(f_1), \dots, \widehat{S}_k^{j,m} := \mathcal{R}_k(f_m)$
  - 7:   Evaluate log-likelihood  $\log(\widehat{S}_k^{j,t})^\# \eta$  for iterations  $t = 1, \dots, m$  with sample  $\chi_j^{\text{test}}$
  - 8: **end for**
  - 9: Define  $m^* \in \{1, \dots, m\}$  as the minimizer of negative log-likelihood  $t \mapsto \sum_{j=1}^\nu -\log(\widehat{S}_k^{j,t})^\# \eta(\chi_j^{\text{test}})$
  - 10: Estimate map  $f^*$  using Algorithm 1 for  $m^*$  iterations with sample  $\chi$
  - 11: **Output:**  $\widehat{S}_k = \mathcal{R}_k(f^*)$
- 

## 5. Numerical experiments

In this section, we evaluate the performance of the ATM algorithm for learning monotone triangular transport maps associated with a variety of target distributions. Sections 5.1 and 5.2 evaluate and visualize the approximation power of the proposed method for strongly non-Gaussian targets. Section 5.3 illustrates the benefits of adaptivity over a range of sample sizes, and Section 5.4 demonstrates how the estimated maps reveal and exploit conditional independence structure in the target distribution. Section 5.5 presents joint and conditional density estimation results for a suite of UCI datasets. Code for reproducing all of these numerical results is available online.<sup>3</sup>

In all of the numerical experiments, we pre-process the data by standardizing each marginal, i.e., subtracting the empirical mean and dividing each component of  $\mathbf{X}$  by its empirical standard deviation. Within the rectifier  $\mathcal{R}_k$ , we employ the modified soft-plus function  $g(\xi) = \log(1 + 2^\xi)/\log(2)$ , which satisfies the conditions of Propositions 4, 5, and 6. This function also satisfies  $g(0) = 1$ , so that  $f(\mathbf{x}_{1:k}) = 0$  is transformed into  $\mathcal{R}_k(f)(\mathbf{x}_{1:k}) = x_k$ . We evaluate the integral in (3.1) numerically using an adaptive quadrature method based on Clenshaw-Curtis points with a relative error of  $10^{-3}$ . At each iteration of the ATM algorithm, we optimize over the coefficients  $c_\alpha$  for  $\alpha \in \Lambda_t$  using a BFGS (quasi-Newton) method [38].

### 5.1. One-dimensional examples

We first consider a one-dimensional mixture of Gaussians with density  $\pi(x) = 0.5\mathcal{N}(x; -2, 0.5) + 0.5\mathcal{N}(x; 2, 2)$ . To estimate  $\pi$ , we generate  $n = 10^4$  training samples  $\{X^i\}_{i=1}^n$  and use them to estimate the map  $S_{\text{KR}}$  that pushes forward  $\pi$  to  $\eta$ . In one dimension, the KR rearrangement (which coincides with the optimal transport map [51]) is simply the increasing function  $S_{\text{KR}} = F_\eta^{-1} \circ F_\pi$ , where

---

<sup>3</sup><https://github.com/baptistar/ATM>

$F_\eta$  and  $F_\pi$  denote the cumulative distribution functions of the reference and target distributions, respectively. We can thus compare our estimated maps to this exact expression.

We approximate  $S$  as the transformation  $\mathcal{R}(f)$  of a smooth non-monotone function  $f$ , seeking  $f$  in the finite-dimensional space  $V_1^p \subset V_1$  spanned by polynomials of degree at most  $p$ . Figure 4a plots the estimated pullback densities  $\hat{\pi} = \hat{S}^\# \eta$  for increasing polynomial degree  $p$ , or equivalently an increasing number of terms  $\#\Lambda$  in  $f$  (30). Figure 4b shows convergence of this approximation in KL divergence, again for increasing  $p$ . We also observe convergence of the estimated map  $\hat{S}$  towards the true KR rearrangement in the  $L_\pi^2$  norm, as expected from the upper bound in (4). Figure 5 illustrates the approximation to the map  $S$  and the associated non-monotone function  $f$  for different choices of basis  $\{\psi_\alpha\}_\alpha$ . Here we observe that approximating  $f$  using the modified Hermite polynomials (31) very closely tracks the KR rearrangement  $S_{\text{KR}}$  and the function  $f_{\text{KR}} = \mathcal{R}^{-1}(S_{\text{KR}})$  as  $x \rightarrow \pm\infty$ , as compared to standard Hermite polynomials or Hermite functions [8]. Enforcing linear asymptotic behavior of the transport map ensures that the approximate distribution has Gaussian tails in regions where there are few samples.

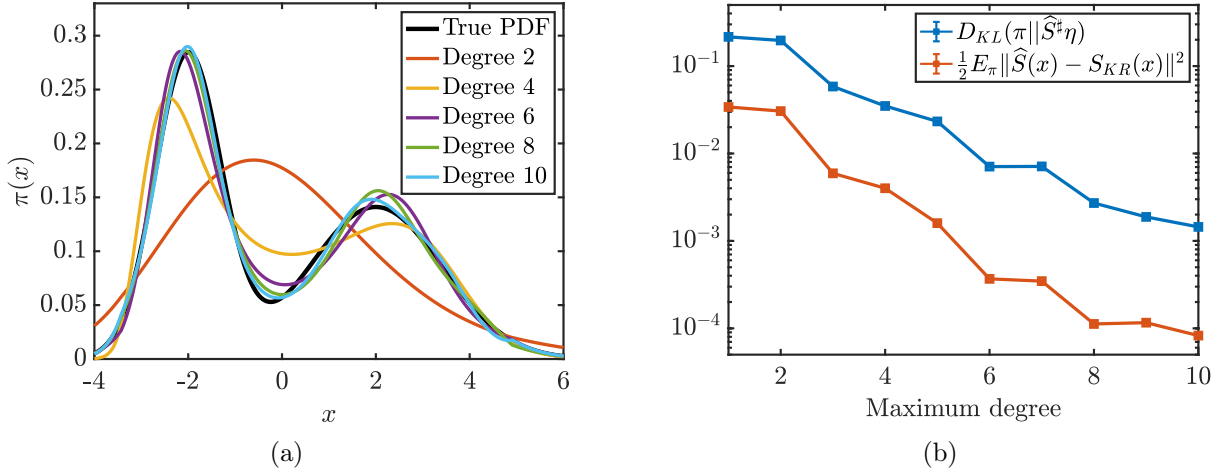


Fig 4: (a) The pullback density  $\hat{S}^\# \eta$  approaches the target Gaussian mixture density  $\pi$  when increasing the maximum polynomial degree  $p$  of the space  $V_1^p \subset V_1$ . (b) With increasing  $p$ , the pullback density converges to  $\pi$  in KL divergence, and the estimated map converges to  $S_{\text{KR}}$  in  $L_\pi^2$ .

Next we consider a Gaussian mixture with density  $\pi(x) = 0.5\mathcal{N}(x; 0, 1) + 0.5\mathcal{N}(x; 0, 0.025)$ . This mixture is a common test case for sampling and density estimation methods, as it evaluates their ability to capture densities with multiple scales [57]. Given  $n = 10^4$  samples, we compute the approximate map using either modified Hermite polynomials (Section 4.1) or Mexican hat wavelets (Section 4.2). Figure 6(a) plots approximations to the target density using an  $m = 15$  term expansion (30) for  $f$ , while Figure 6(b) shows convergence in KL divergence for both the polynomial and wavelet bases. We observe that the polynomial approximation suffers from oscillations, while the wavelets better capture localized features. This results in a much faster convergence of the KL divergence for wavelets than for polynomials.

## 5.2. Two-dimensional datasets

To demonstrate the expressiveness of the modified Hermite polynomial basis (31), we use the ATM algorithm to approximate the Knothe–Rosenblatt rearrangement for several two-dimensional distributions with strongly non-Gaussian geometries: the “banana,” “funnel,” “cosine,” “mixture

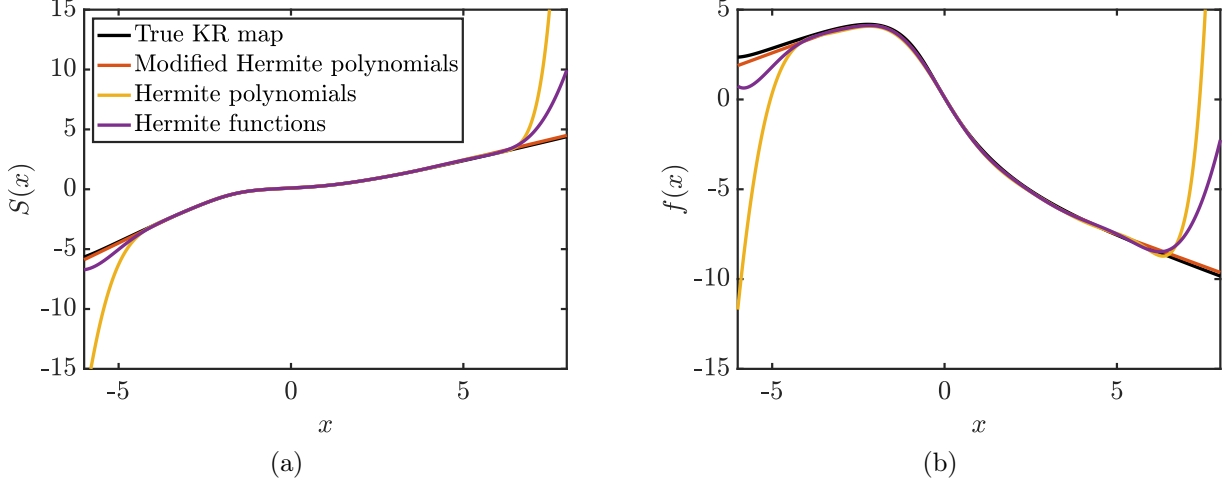


Fig 5: (a) The approximate transport maps  $\hat{S}$  compared to  $S_{\text{KR}}$  (black), and (b) the corresponding non-monotone functions  $f$  compared to  $f_{\text{KR}} := \mathcal{R}^{-1}(S_{\text{KR}})$ . Both subfigures illustrate different choices of basis  $\{\psi_\alpha\}_\alpha$ , for the Gaussian mixture target of Figure 4. The modified Hermite polynomial basis provides the closest approximation to  $S_{\text{KR}}$  and  $f_{\text{KR}}$ .

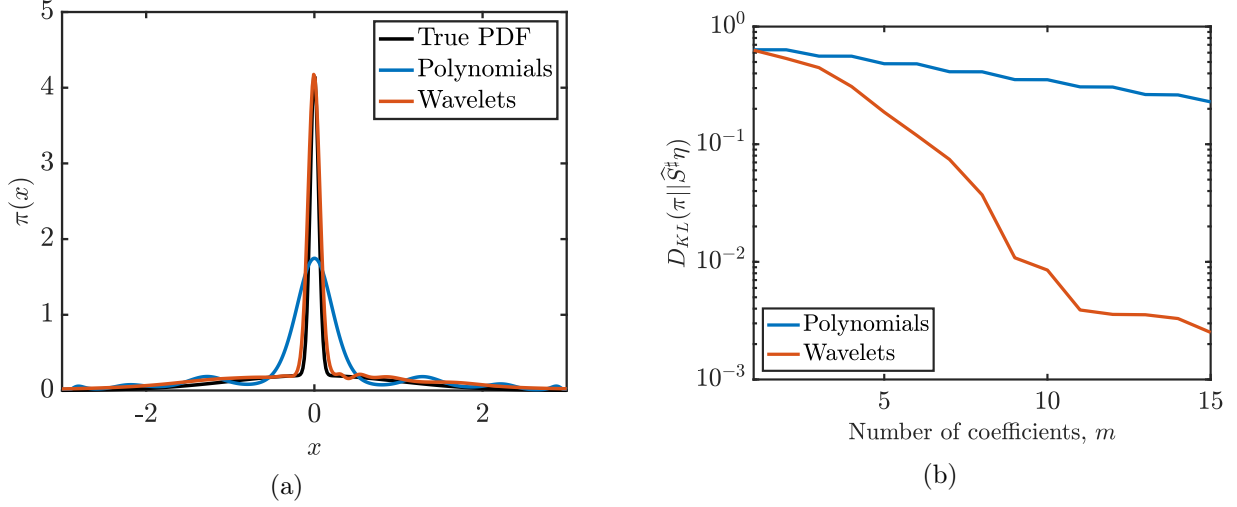


Fig 6: Approximation of a scale mixture of Gaussians, using Hermite polynomials or Ricker wavelet expansions. (a) A wavelet expansion with  $m = 15$  coefficients approximates the target density better than a polynomial expansion of the same size, and (b) yields significantly lower KL divergence for this example.

of Gaussians (MoG),” and “ring” distributions  $\pi$  considered in [70, 22]. For each distribution, we generate an i.i.d. sample of size  $n = 10^4$  from  $\pi$  and apply a random rotation to the data using an angle uniformly distributed in  $[0, \pi/2]$ . Figure 7 plots the true densities  $\pi$  and the approximate densities  $\hat{\pi} := \hat{S}^\# \eta$  found using ATM. We use 5-fold cross-validation to identify the optimal number of elements in the multi-index set for each map component:  $S_1$  and  $S_2$ . Figure 7 indicates the total number of coefficients  $\#\Lambda_{m^*}$  across both map components, underscoring the parsimony of the approximation.

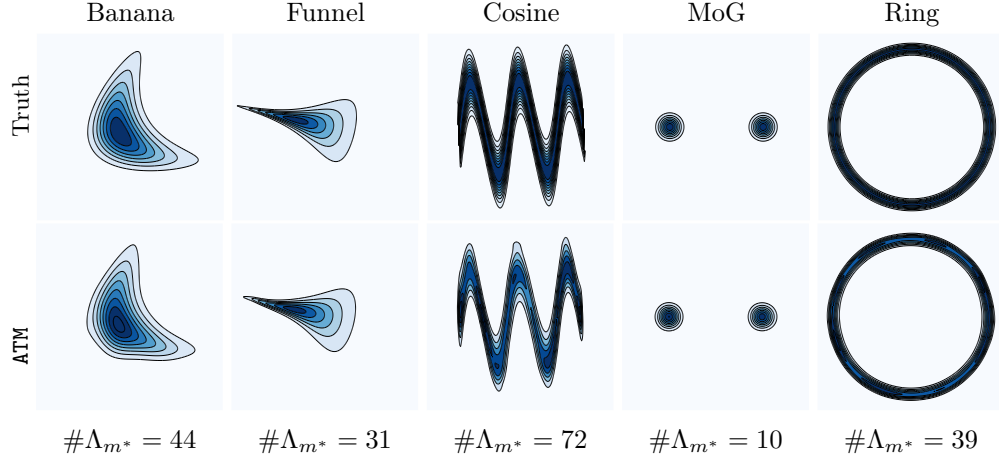


Fig 7: Five different densities  $\pi$  (top row) and their approximations  $\hat{\pi}$  (bottom row) built using the ATM algorithm with a Hermite polynomial basis for  $V_k$ , given a sample of size  $n = 10^4$  from  $\pi$ . The total number of coefficients  $\#\Lambda_{m^*}$  in both map components is indicated below each density.

### 5.3. Random mixture of Gaussians

In this example, we evaluate the stability and accuracy of the ATM algorithm for learning transport maps in the small sample regime. We consider a three-dimensional distribution  $\pi$  defined as a mixture of Gaussians centered at the eight vertices of the hypercube  $[-4, 4]^3$ , each with covariance matrix  $I_d$ . The weights of the mixture components are randomly sampled from a uniform distribution  $\mathcal{U}([0, 1])$  and then normalized so that  $\int_{\mathbb{R}^3} d\pi = 1$ . To estimate  $\pi$ , we generate a training sample  $\{\mathbf{X}^i\}_{i=1}^n \sim \pi$  and use the ATM algorithm to estimate the KR rearrangement that pushes forward  $\pi$  to  $\eta$  using 5-fold cross validation.

Figure 8a plots the KL divergence of  $\hat{\pi}$  from  $\pi$  averaged over 10 experiments. Here, the training sets used to build  $\hat{\pi} = \hat{S}^\# \eta$  are of varying size  $n \in [10^1, 10^4]$  and the reported KL divergence is computed on a common test set of  $10^4$  samples. Figure 8b plots the total number of coefficients  $\#\Lambda_{m^*}$  identified by the ATM algorithm. The performance of ATM is compared to a non-adaptive method where  $\Lambda = \Lambda(p) := \{\alpha \in \mathbb{N}_0^k, \|\alpha\|_1 \leq p\}$  is arbitrarily fixed, for  $p = 1, 3, 5$ . Note that  $\Lambda(p)$  corresponds to polynomial  $f$  with total degree  $p$ . For each sample size  $n$ , ATM consistently finds a better estimator of  $\pi$  (in the sense of KL divergence) than a non-adaptive method with a fixed degree  $p$ , by adaptively identifying the basis  $\{\psi_\alpha\}_{\alpha \in \Lambda_t}$  to represent the map components. In addition, the ATM estimator achieves smaller KL divergence with fewer total map coefficients, as seen in Figure 8b.

### 5.4. Stochastic volatility

Next we consider data from a Markov process that describes the volatility of the return on a financial asset over time. The model has two hyperparameters  $\mu$  and  $\phi$ , with a state  $(Z_k)_k$  that represents the log-volatility at times  $k = 1, \dots, T$ . The two hyperparameters are drawn from the distributions  $\mu \sim \mathcal{N}(0, 1)$  and  $\phi = 2 \exp(\phi^*) / (1 + \exp(\phi^*))$  for  $\phi^* \sim \mathcal{N}(3, 1)$ . The log-volatility obeys the order-one autoregressive process  $Z_{k+1} = \mu + \phi(Z_k - \mu) + \epsilon_k$  for all  $k > 1$ , where  $\epsilon_k \sim \mathcal{N}(0, 1)$  is independent of all other variables and  $Z_0 | \mu, \phi \sim \mathcal{N}(\mu, \frac{1}{1-\phi^2})$ . While the states are Gaussian conditioned on the hyperparameters, the joint distribution of

$$\mathbf{X} = (\mu, \phi, Z_1, \dots, Z_T),$$

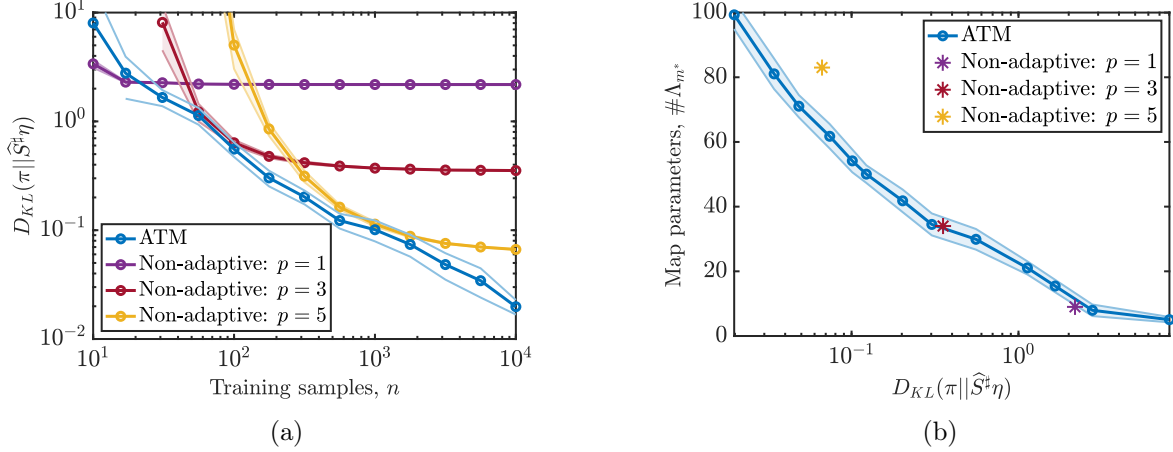


Fig 8: (a) KL divergence over 10 sets of training samples is lower for **ATM** than for non-adapted expansions. (b) Trade-off between the approximation quality and number of coefficients using the adaptive algorithm for different  $n$ . For any given number of map coefficients, the **ATM** algorithm finds a representation with similar or lower KL divergence than the lowest achievable error of each non-adaptive approximation (indicated with stars in plot (b)).

is non-Gaussian. In this example, the dimension of  $\mathbf{X}$  is made arbitrarily large by increasing  $T$ .

A key property of triangular transport maps is that they inherit sparsity from the conditional independence structure of  $\pi$ . [59] showed that the Markov structure of  $\pi$  yields a lower bound on the sparsity of the map  $S$  (i.e., the functional dependence of each component  $S^k$  on the input variables  $\mathbf{x}_{1:k-1}$ ). This sparsity was exploited to learn the structure of undirected probabilistic graph models from data in [36, 3]. From the conditional independence properties of the stochastic volatility model described above, we know that the KR rearrangement  $S_{KR}$  between the joint distribution of  $\mathbf{X}$  and a standard Gaussian reference  $\eta$  is sparse. Moreover, the exact sparsity of  $S_{KR}$  can be derived from Theorem 5.1 in [59]. We now show that the **ATM** algorithm is capable of detecting and exploiting this structure—without knowing in advance that it exists.

Figure 9a compares the variable dependence of the true KR rearrangement  $S_{KR}$  and the map  $\hat{S}$  learned by the **ATM** algorithm for a distribution with  $T = 40$  using a sample of size  $n = 10^3$  from  $\pi$ ; a non-filled entry  $(j, k)$  entry of the plot indicates that  $k$ -th map component does not depend on variable  $x_j$ . The dependence of component  $S_k$  on  $(x_k, x_{k-1})$  shows that each state  $Z_k$  strongly depends on the previous state in time. Most of the map components also show dependence on the hyperparameters  $(\mu, \phi)$ . The estimated sparsity closely matches the exact sparsity of the KR rearrangement. Furthermore, the sparse variable dependence of the  $k$ th map component  $\hat{S}_k$  on parent nodes  $\text{Pa}(k) \subseteq \{1, \dots, k-1\}$  produces an approximation to the  $k$ th marginal conditional density given by

$$\hat{\pi}(x_k | \mathbf{x}_{<k}) = \hat{\pi}(x_k | \mathbf{x}_{\text{Pa}(k)}).$$

By identifying the parent nodes  $\text{Pa}(k)$  of each variable  $k$ , we also learn a sparse Bayesian network or directed acyclic graphical (DAG) model representing the target distribution [26]. As a result, we can see the **ATM** algorithm as a technique for learning DAGs from samples, given a prescribed variable ordering.

Next, we consider the approximation of  $S_{KR}$  for Markov models of *increasing* state dimension  $T$  and hence increasing map dimension  $d = T + 2$ ; these experiments use a fixed sample size  $n = 10^3$  as before. Figure 9b plots the KL divergence from the **ATM** approximation to the joint density of



$\mathbf{X} \in \mathbb{R}^d$ , for increasing dimensions  $d$ . We compare ATM with a variant of ATM where the exact sparsity pattern of the the KR rearrangement is provided to the algorithm in advance. This variant, labeled “sparsity-aware ATM” in Figure 9b, differs from the ATM algorithm in that it only activates multi-indices  $\alpha_t^*$  which match the sparsity of  $S_{KR}$  (meaning of the form  $\alpha = (\alpha_1, \alpha_2, 0, \dots, 0, \alpha_{k-1}, \alpha_k)$ ). We also compare ATM with non-adaptive maps of degree  $p = 1$  and  $p = 2$ ; these maps do not exploit the conditional independence structure of  $\pi$ , and hence depend on all input variables. For low dimensions, we see that degree  $p = 2$  maps can better capture the non-Gaussian target distributions than  $p = 1$  linear maps. As the dimension  $d$  increases, however, the growing number of coefficients in the  $p = 2$  maps results in higher-variance estimators, outweighing their smaller bias, and hence we see a rapidly increasing KL divergence as the dimension  $d$  increases (with crossover around  $d = 21$ ). In contrast, ATM outperforms the non-adaptive maps for all  $d$  and, more importantly, achieves a KL divergence that grows linearly with  $d$ ; indeed, it performs similarly to the best case “oracle sparsity” of sparsity-aware ATM.

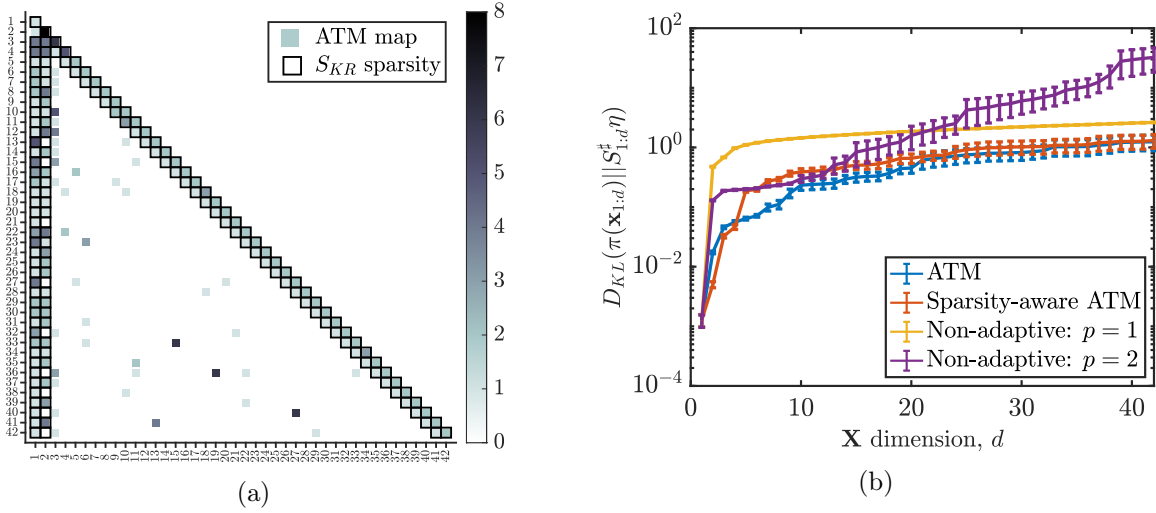


Fig 9: (a) The sparsity of an ATM map  $\hat{S}$  and of the KR rearrangement  $S_{KR}$  for the stochastic volatility model with  $T = 40$  time steps. (b) KL divergence using ATM with increasing dimension  $d$ , compared to adaptive map estimators with known sparsity and non-adaptive maps estimators.

### 5.5. Tabular datasets

Lastly, we evaluate ATM’s performance for density estimation on a suite of UCI datasets [30] that were also considered in [64]. These datasets have dimensionalities between  $d = 10$  and 15 and sample sizes between  $n = 506$  and 5875. We pre-process each dataset to eliminate discrete-valued variables and one variable in every pair that has Pearson correlation coefficient greater than 0.98, following the procedure in [64]. We consider 10 splits of the data. For each split, we use one fold (i.e., 10% of the data) as a test set and the remaining 9 folds (i.e., 90% of the data) as the training set to build  $\hat{S}$ . To assess the quality of our estimated map  $\hat{S}$ , we evaluate the negative log-likelihood of the pullback density  $\hat{\pi} = \hat{S}^\# \eta$  on the test set. The negative log-likelihood is an empirical estimator of  $-\mathbb{E}_\pi[\log \hat{S}^\# \eta] = \mathcal{D}_{KL}(\pi || \hat{S}^\# \eta) - \int \log \pi d\pi$  and is, up to the unknown constant  $-\int \log \pi d\pi$ , the KL divergence from  $\hat{S}^\# \eta$  to  $\pi$ . Table 1 presents the mean of the negative log-likelihoods over the 10 splits and a 95% confidence interval for the mean. For all datasets, we observe an improvement using ATM

over non-adaptive  $p = 2$  maps and multivariate Gaussian approximations (i.e.,  $p = 1$  maps), with a lower number of total coefficients.

TABLE 1

Mean negative log-likelihood for UCI datasets over 10 sets of training samples.  
The estimator with best performance (lowest negative log-likelihood) is highlighted in bold.

Dataset	$(d, n)$	Gaussian	ATM	# coefficients	$p = 2$ maps	# coefficients
Boston	(10, 506)	$11.3 \pm 0.5$	<b><math>3.1 \pm 0.6</math></b>	$228 \pm 7$	$6.5 \pm 0.4$	285
Red Wine	(11, 1599)	$13.2 \pm 0.3$	<b><math>9.8 \pm 0.4</math></b>	$289 \pm 9$	$10.5 \pm 0.2$	363
White Wine	(11, 4898)	$13.2 \pm 0.5$	<b><math>11.0 \pm 0.2</math></b>	$342 \pm 26$	$12.0 \pm 1.0$	363
Parkinsons	(15, 5875)	$10.8 \pm 0.4$	<b><math>2.8 \pm 0.4</math></b>	$783 \pm 17$	$5.1 \pm 0.4$	815

Lastly, we evaluate ATM’s performance for *conditional* density estimation on a different suite of UCI datasets. We follow a similar procedure as above to pre-process each dataset. Each dataset has a one-dimensional predictor variable  $X_d$  and covariates  $\mathbf{X}_{<d}$  of varying dimension. To approximate the conditional density  $\pi_d(x_d|\mathbf{x}_{<d})$ , we estimate one map component  $S_d$  as a function of the predictor and the covariates using joint samples  $\{(\mathbf{X}_{<k}^i, X_k^i)\}_{i=1}^n$ . Table 2 presents the mean negative conditional log-likelihoods over 10 splits of the training data. The negative conditional log-likelihood is an empirical estimator of  $-\mathbb{E}_\pi[\log \hat{S}_d^\# \eta] = \mathbb{E}_{\pi(\mathbf{x}_{<d})}[\mathcal{D}_{\text{KL}}(\pi_d(x_d|\mathbf{x}_{<d})||\hat{S}_d^\# \eta)] - \int \log \pi_d(x_d|\mathbf{x}_{<d}) d\pi_d$  and is, up to the unknown constant  $-\int \log \pi_d(x_d|\mathbf{x}_{<d}) d\pi_d$ , the expected KL divergence from  $\hat{S}_d^\# \eta$  to the marginal conditional  $\pi_d(x_d|\mathbf{x}_{<d})$ .

In Table 2, we compare ATM to conditional kernel density estimation (CKDE) [56],  $\epsilon$ -neighborhood kernel density estimation (NKDE), and kernel mixture networks (KMN) [1] using the implementation provided by [50]. We also include results for high-capacity parametric conditional density estimation methods that rely on neural networks: mixture density networks (MDN) [6] and conditional normalizing flows (NF) [62]; implementation details of these methods are provided in Appendix C. For all datasets, we observe that ATM has performance comparable to the neural-network based methods and improved performance over the nonparametric approaches (e.g., on the Concrete and Yacht datasets).

In Table 3, we compare the number of coefficients in the ATM and MDN approximations and the computational time required to learn each of these conditional density representations. The runtime of the greedy procedure in Algorithm 1 (sequentially, for all  $d$  map components) is generally less than that of MDN, achieving similar approximation performance with about two orders of magnitude fewer coefficients. It is more realistic, however, to consider ATM with the cross-validation procedure of Algorithm 2 and to compare this with the hyperparameter tuning required to achieve the reported performance of MDN. ATM only requires performing cross-validation for a single hyper-parameter, i.e., the total number of coefficients  $m^*$  (see Algorithm 2). In comparison, tuning several hyperparameters in the neural-network based methods (e.g., the learning rate, and the structure of the networks) by performing cross-validation over a tensor product grid increases runtimes by more than an order of magnitude relative to ATM, as presented below.

TABLE 2

Mean negative conditional log-likelihood for UCI datasets over 10 sets of training samples.  
The method with best performance from both categories is highlighted in bold.

Dataset	$(d, n)$	ATM	CKDE	NKDE	MDN	KMN	NF
Boston	(12, 506)	<b><math>2.6 \pm 0.2</math></b>	<b><math>2.6 \pm 0.2</math></b>	$3.1 \pm 0.2$	<b><math>2.4 \pm 0.2</math></b>	$2.7 \pm 0.2$	<b><math>2.4 \pm 0.1</math></b>
Concrete	(9, 1030)	<b><math>3.1 \pm 0.1</math></b>	$3.2 \pm 0.1$	$3.9 \pm 0.1$	<b><math>2.9 \pm 0.1</math></b>	$3.5 \pm 0.1$	$3.2 \pm 0.2$
Energy	(10, 768)	$1.5 \pm 0.1$	<b><math>1.0 \pm 0.1</math></b>	$2.1 \pm 0.2$	<b><math>1.2 \pm 0.1</math></b>	$1.7 \pm 0.1$	$1.7 \pm 0.3$
Yacht	(7, 308)	<b><math>0.5 \pm 0.2</math></b>	$1.1 \pm 0.3$	$3.8 \pm 0.2$	<b><math>0.7 \pm 0.2</math></b>	$1.8 \pm 0.2$	$1.3 \pm 0.5$

TABLE 3

The number of coefficients in the *ATM* and *MDN* approximations, as well as the runtimes (in seconds) to identify coefficients via optimization and to perform cross-validation of hyperparameters.

Dataset	ATM			MDN		
	# coefficients	runtime (s)	runtime w/CV	# coefficients	runtime (s)	runtime w/CV
Boston	$31 \pm 3$	$6 \pm 1$	$88 \pm 14$	$(5 \pm 2) \times 10^3$	$18 \pm 8$	$1584 \pm 13$
Concrete	$42 \pm 5$	$14 \pm 3$	$230 \pm 18$	$(2.6 \pm 0.3) \times 10^3$	$11 \pm 1$	$1846 \pm 10$
Energy	$41 \pm 7$	$12 \pm 5$	$193 \pm 29$	$(8 \pm 3) \times 10^3$	$32 \pm 12$	$1758 \pm 17$
Yacht	$28 \pm 4$	$7 \pm 2$	$99 \pm 20$	$(6 \pm 2) \times 10^3$	$23 \pm 6$	$1453 \pm 15$

## 6. Conclusions

This paper has presented and analyzed a functional framework for learning monotone triangular transport maps. Our approach represents monotone component functions of a triangular map through the action of an invertible *rectification* operator  $\mathcal{R}_k$  on smooth, generally non-monotone, functions. Imposing appropriate structure on this operator and the function space  $V_k$  on which it acts, along with conditions on the target density  $\pi$ , yields an unconstrained optimization problem for learning the map components that has many desirable and useful properties. First, under certain assumptions on  $\mathcal{R}_k$  and  $\pi$ , we show that the optimization objective is bounded, continuous, and differentiable for all  $f \in V_k$ . This permits the use of deterministic gradient-based optimization methods to find the minimizers. Next, under certain conditions on  $\mathcal{R}_k$ , we show that the optimization problem has *no spurious local minima* on  $V_k$ . In practice, this yields important robustness to the initial guess and other parameters of the optimization algorithms. Finally, under the same (cumulative) conditions on  $\mathcal{R}_k$  and some additional assumptions on the target density, we show that the optimization problem has a unique global minimizer on  $V_k$  that corresponds to the canonical KR rearrangement.

Our functional framework also enables the construction of novel transport map estimators, based on maximum likelihood, given a finite sample from the target. The procedure we develop here is semi-parametric, making use of different hierarchical bases for  $V_k$ . In particular, we propose a greedy, adaptive algorithm that identifies a sparse set of basis functions, automatically tailored to the target density and to the sample size. Sparsity in the basis selection reflects conditional independence structure in the target. More generally, the maps we build can capture the structure of complex probability distributions with sufficient data, but also are robust in settings with few observations. This is crucial to deploying these algorithms within large-scale applications, such as data assimilation, where triangular transport maps must be learned online [58]. We demonstrate good computational performance of our algorithm—and the parsimonious representations it produces—in a variety of joint and conditional density estimation problems. We outline some interesting future research directions below.

**KR recovery for broader classes of target distributions** While our result guaranteeing no spurious local minima (Theorem 7) makes no explicit assumptions on  $\pi$ , our subsequent results on recovery of the KR rearrangement (see Corollary 9) make particular use of Gaussian tail assumptions, relating the target  $\pi$  to the Gaussian reference  $\eta$ . In principle, however, the space  $V_k$  contains functions  $\mathcal{R}_k^{-1}(S_{\text{KR},k})$  corresponding to a broader class of target densities;  $V_k$  includes functions that grow up to sub-exponentially fast, as a result of its Gaussian weighting. Hence, it includes the inverse images of KR rearrangements for lighter-tailed  $\pi$ , whose component functions grow faster than linear as  $x_k \rightarrow \pm\infty$ . It may be possible to obtain unique recovery results analogous to Corollary 9 for such densities, although some modifications to our analysis would be needed.

**Approximation theory and statistical guarantees** An open research topic, to the best of our knowledge, is to analyze the finite-dimensional approximation of triangular transport maps on unbounded domains. (See [71, 72] for an approximation theory of triangular transports between analytic densities  $\pi, \eta$  on bounded domains.) These results would show how approximate maps converge to the KR rearrangement with different parameterizations (e.g., polynomial spaces or neural networks of increasing size), and we expect that our rectification framework could be a useful ingredient of such analyses. A parallel line of work is to develop non-asymptotic statistical convergence results for triangular transports (e.g., in the context of density estimation), to understand how the quality of density estimates and of map estimates depends on the sample size  $n$ . Combining these lines of work could provide lower bounds on the number of samples required to learn maps of a given complexity, and also some analytical guidance for how to navigate the bias-variance tradeoff of a map estimator given finite samples.

**Other sources of low-dimensional structure** The ATM algorithm in this paper identifies maps with sparse variable dependence by exploiting conditional independence structure in the target distribution [59]. In future work it will be interesting to investigate other notions of low-dimensional structure and how they could facilitate the learning of maps from small samples. One promising notion is when the target density  $\pi$  departs from the reference  $\eta$  only along a low-dimensional subspace. In this case, the triangular map  $S$  pushing forward  $\pi$  to  $\eta$  can be written as a low-dimensional perturbation of the identity map, after a variable rotation. See [9] for a recent contribution in this direction applied to variational inference.

**Map ordering** One disadvantage of triangular maps is that the approximation depends on the choice of variable ordering. In particular, each variable ordering yields a different factorization of the target density and a different Knothe–Rosenblatt rearrangement  $S_{\text{KR}}$ . Thus, it is of interest to develop variable ordering algorithms that minimize the finite-sample error of the estimated transport map  $\hat{S}$ , the estimated pullback density  $\hat{\pi} = \hat{S}^\# \eta$ , or other goal-oriented metrics. For target distributions that induce sparse variable dependence in the map, one approach is to find the permutation that maximizes the sparsity of  $\hat{S}$ . This is equivalent to finding the sparsest Bayesian network or directed acyclic graphical (DAG) model for a distribution from samples. While this problem is in general NP-complete, effective algorithms [47] have been proposed. This algorithm reduces to finding a sparse maximum likelihood estimator for the Cholesky factor of the inverse covariance matrix of  $\pi$  in the Gaussian setting. Since a linear transport map is precisely this Cholesky factor for Gaussian  $\pi$  and standard normal  $\eta$  (see [3, Section 3]), we expect that such sparse permutation algorithms could be generalized to the nonlinear transport map setting.

**Nonparametric methods** Instead of finding a particular finite-dimensional basis for the map components  $S_k = \mathcal{R}_k(f)$ , nonparametric methods do not limit the functional form of  $f$  (e.g., as linear combinations of multivariate polynomials). A broadly useful nonparametric approach involves seeking  $f$  in a reproducing kernel Hilbert space (RKHS). Thanks to the representer theorem [54], the optimal  $f$  in the RKHS can still be identified by solving a finite-dimensional optimization problem. Choosing the kernel so that the RKHS is a suitable weighted Sobolev space (see, e.g., [39]) will yield map representations that fall within the framework of Section 3.2. A related approach uses Gaussian process representations for the map components [23]. It will be interesting to compare the finite-sample performance of such methods to the semi-parametric procedure proposed in this work.

## Acknowledgments

RB, YM, and OZ gratefully acknowledge support from the INRIA associate team Unquestionable. RB and YM are also grateful for support from the AFOSR Computational Mathematics program (MURI award FA9550-15-1-0038) and the US Department of Energy AEOLUS center. RB acknowledges support from an NSERC PGSD-D fellowship. OZ also acknowledges support from the ANR JCJC project MODENA (ANR-21-CE46-0006-01).

## Appendix A: Proofs and theoretical details

### A.1. Proof of Proposition 1

*Proof.* Recall that the KR rearrangement  $S_{\text{KR}}$  is a transport map that satisfies  $S_{\text{KR}}^\# \eta = \pi$ , where  $\eta$  is the density of the standard Gaussian measure on  $\mathbb{R}^d$  and  $\pi$  is the target density. Corollary 3.10 in [7] states that for any PDF  $\varrho$  on  $\mathbb{R}^d$  of the form  $\varrho(\mathbf{x}) := f(\mathbf{x})\eta(\mathbf{x})$  with  $f \log f \in L_\eta^1$ , the inequality

$$\int \|\mathbf{x} - T(\mathbf{x})\|^2 \eta(\mathbf{x}) d\mathbf{x} \leq 2 \int f(\mathbf{x}) \log f(\mathbf{x}) \eta(\mathbf{x}) d\mathbf{x}, \quad (37)$$

holds, where  $T$  is the KR rearrangement such that  $T_\# \eta = \varrho$ . Let  $S$  be an increasing lower triangular map as in (1) and let  $\varrho = S_\# \pi$ . Thus we have  $T = S \circ S_{\text{KR}}^{-1}$  and so the left-hand side of (37) becomes

$$\int \|\mathbf{x} - T(\mathbf{x})\|^2 \eta(\mathbf{x}) d\mathbf{x} = \int \|\mathbf{x} - S \circ S_{\text{KR}}^{-1}(\mathbf{x})\|^2 \eta(\mathbf{x}) d\mathbf{x} = \int \|S_{\text{KR}}(\mathbf{x}) - S(\mathbf{x})\|^2 \pi(\mathbf{x}) d\mathbf{x},$$

and the right-hand side becomes

$$2 \int f(\mathbf{x}) \log f(\mathbf{x}) \eta(\mathbf{x}) d\mathbf{x} = 2\mathcal{D}_{\text{KL}}(\varrho \parallel \eta) = 2\mathcal{D}_{\text{KL}}(\pi \parallel S^\# \eta),$$

which yields (4). □

### A.2. Convexity of $s \mapsto \mathcal{J}_k(s)$

**Lemma 10.** *The optimization problem  $\min_{\{s: \partial_k s > 0\}} \mathcal{J}_k(s)$  is strictly convex.*

*Proof.* Let  $s_1, s_2: \mathbb{R}^k \rightarrow \mathbb{R}$  be two strictly increasing functions with respect to  $x_k$ , i.e.,  $\partial_k s_1(\mathbf{x}_{\leq k}) > 0$  and  $\partial_k s_2(\mathbf{x}_{\leq k}) > 0$ . For any  $0 < t < 1$ , the function  $s_t = ts_1 + (1-t)s_2$  is also strictly increasing with respect to  $x_k$ . Finally, because both  $\xi \mapsto \frac{1}{2}\xi^2$  and  $\xi \mapsto -\log(\xi)$  are strictly convex functions, we have

$$\begin{aligned} \mathcal{J}_k(s_t) &\stackrel{(6)}{=} \int \left( \frac{1}{2} s_t(\mathbf{x}_{\leq k})^2 - \log \partial_k s_t(\mathbf{x}_{\leq k}) \right) \pi(\mathbf{x}) d\mathbf{x} \\ &< \int \left( t \frac{1}{2} s_1(\mathbf{x}_{\leq k})^2 + (1-t) \frac{1}{2} s_2(\mathbf{x}_{\leq k})^2 \right) - \left( t \log \partial_k s_1(\mathbf{x}_{\leq k}) + (1-t) \log \partial_k s_2(\mathbf{x}_{\leq k}) \right) \pi(\mathbf{x}) d\mathbf{x} \\ &= t \mathcal{J}_k(s_1) + (1-t) \mathcal{J}_k(s_2), \end{aligned}$$

which shows that  $\mathcal{J}_k$  is strictly convex. □

### A.3. Proof of Proposition 2

To prove Proposition 2 we need the following lemma.

**Lemma 11.** *Let*

$$H^1([0, 1]) = \left\{ f : [0, 1] \rightarrow \mathbb{R} \text{ such that } \|f\|_{H^1([0,1])}^2 := \int_0^1 f(t)^2 + f'(t)^2 dt \right\}.$$

Then

$$|f(0)| \leq \sqrt{2} \|f\|_{H^1([0,1])}, \quad (38)$$

holds for any  $f \in H^1([0, 1])$ .

*Proof.* Because  $\mathcal{C}^\infty([0, 1])$  is dense in  $H^1([0, 1])$ , it suffices to show (38) for any  $f \in \mathcal{C}^\infty([0, 1])$ . By the mean value theorem, there exists  $0 \leq z \leq 1$  such that

$$f(z) = \frac{1}{1-0} \int_0^1 f(t) dt.$$

Thus we can write

$$\begin{aligned} |f(0)|^2 &\leq 2|f(z) - f(0)|^2 + 2|f(z)|^2 \\ &= 2 \left| \int_0^z f'(t) dt \right|^2 + 2 \left| \int_0^1 f(t) dt \right|^2 \\ &\leq 2 \int_0^1 |f'(t)|^2 dt + 2 \int_0^1 |f(t)|^2 dt. \end{aligned}$$

This concludes the proof.  $\square$

We now prove Proposition 2.

*Proof.* For any  $f \in V_k$ , Lemma 11 permits us to write

$$\begin{aligned} \int |f(\mathbf{x}_{<k}, 0)|^2 \eta_{<k}(\mathbf{x}_{<k}) d\mathbf{x}_{<k} &\stackrel{(38)}{\leq} \int \left( 2 \int_0^1 |f(\mathbf{x}_{<k}, t)|^2 + |\partial_k f(\mathbf{x}_{<k}, t)|^2 dt \right) \eta_{<k}(\mathbf{x}_{<k}) d\mathbf{x}_{<k} \\ &\leq C_T \int \int_0^1 \left( |f(\mathbf{x}_{<k}, t)|^2 + |\partial_k f(\mathbf{x}_{<k}, t)|^2 \right) \eta_{<k}(\mathbf{x}_{<k}) \eta_1(t) d\mathbf{x}_{<k} dt \\ &\leq C_T \int \int_{-\infty}^{+\infty} \left( |f(\mathbf{x}_{<k}, t)|^2 + |\partial_k f(\mathbf{x}_{<k}, t)|^2 \right) \eta_{\leq k}(\mathbf{x}_{<k}, t) d\mathbf{x}_{<k} dt \\ &= C_T \|f\|_{V_k}^2, \end{aligned}$$

where  $C_T = 2 \sup_{0 \leq t \leq 1} \eta_1(t)^{-1}$ .  $\square$

### A.4. Proof of Proposition 3

The proof relies on Proposition 2 and on the following generalized integral Hardy inequality, see [37].

**Lemma 12.** *Let  $\eta_{\leq k}$  be the standard Gaussian density on  $\mathbb{R}^k$ . Then there exists a constant  $C_H$  such that for any  $v \in L_\eta^2(\mathbb{R}^k)$ ,*

$$\int \left( \int_0^{x_k} v(\mathbf{x}_{<k}, t) dt \right)^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \leq C_H \int v(\mathbf{x})^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x}. \quad (39)$$

*Proof of Lemma 12.* Let us recall the integral Hardy inequality [37].

**Theorem 13** (from [37]). *For weight  $\rho: \mathbb{R}_+ \rightarrow \mathbb{R}_+$  and  $u \in L^2_\rho(\mathbb{R})$ , there exists a constant  $C_H < \infty$  such that*

$$\int_0^{+\infty} \left( \int_0^x u(t) dt \right)^2 \rho(x) dx \leq C_H \int_0^{+\infty} u(x)^2 \rho(x) dx \quad (40)$$

*if and only if*

$$\sup_{x>0} \left( \int_x^{+\infty} \rho(t) dt \right)^{1/2} \left( \int_0^x \rho(t)^{-1} dt \right)^{1/2} < +\infty. \quad (41)$$

We apply Theorem 13 with the one-dimensional standard Gaussian density  $\rho = \eta$  for  $x > 0$ . In order to check condition (41), we need to show that

$$D(x) := \left( \int_x^{+\infty} \rho(t) dt \right)^{1/2} \left( \int_0^x \rho(t)^{-1} dt \right)^{1/2} = \left( \int_x^{+\infty} e^{-t^2/2} dt \right)^{1/2} \left( \int_0^x e^{t^2/2} dt \right)^{1/2},$$

is bounded. Since  $x \mapsto D(x)$  is a continuous function with a finite limit as  $x \rightarrow 0$ , it is sufficient to show that  $D(x)$  has a finite limit when  $x \rightarrow \infty$ . For  $x > 1$ ,  $\int_x^{+\infty} e^{-t^2/2} dt \leq e^{-x^2/2}$  and  $D(x)^2 \leq e^{-x^2/2} \int_0^x e^{t^2/2} dt$ . Furthermore, using integration-by-parts we have  $\int_0^x e^{t^2/2} dt = \int_0^1 e^{t^2/2} dt + e^{x^2/2}/x - \sqrt{e} + \int_1^x e^{t^2/2}/t^2 dt$ . As  $x \rightarrow \infty$  the dominating term in the sum is  $e^{x^2/2}/x$ . Thus,  $e^{-x^2/2} \int_0^x e^{t^2/2} dt$  behaves asymptotically as  $\mathcal{O}(\frac{1}{x})$ , so that  $D(x) \rightarrow 0$  when  $x \rightarrow \infty$ . Thus, condition (41) is satisfied.

Then, by the Hardy inequality in (40) for  $u \in L^2_\eta(\mathbb{R})$  we have

$$\int_0^{+\infty} \left( \int_0^{x_k} u(t) dt \right)^2 \eta(x_k) dx_k \leq C_H \int_0^{+\infty} u(x_k)^2 \eta(x_k) dx_k. \quad (42)$$

For the symmetric density  $\eta(x_k) = \eta(-x_k)$  we also have

$$\int_{-\infty}^0 \left( \int_0^{x_k} u(t) dt \right)^2 \eta(x_k) dx_k \leq C_H \int_{-\infty}^0 u(x_k)^2 \eta(x_k) dx_k. \quad (43)$$

Combining the results in (42) and (43) we have

$$\int_{-\infty}^{+\infty} \left( \int_0^{x_k} u(t) dt \right)^2 \eta(x_k) dx_k \leq C_H \int_{-\infty}^{+\infty} u(x_k)^2 \eta(x_k) dx_k.$$

Setting  $u(t) = v(\mathbf{x}_{<k}, t)$  and integrating both sides over  $\mathbf{x}_{<k} \in \mathbb{R}^{k-1}$  with the standard Gaussian weight function  $\eta(\mathbf{x}_{<k})$  gives the result.  $\square$

We now prove Proposition 3.

*Proof.* By Proposition 2, Lemma 12 and by the Lipschitz property of  $g$ , we can write

$$\begin{aligned} \|\mathcal{R}_k(f_1) - \mathcal{R}_k(f_2)\|_{L^2_{\eta_{\leq k}}}^2 &\leq 2 \int \left( f_1(\mathbf{x}_{<k}, 0) - f_2(\mathbf{x}_{<k}, 0) \right)^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\ &\quad + 2 \int \left( \int_0^{x_k} g(\partial_k f_1(\mathbf{x}_{<k}, t)) - g(\partial_k f_2(\mathbf{x}_{<k}, t)) dt \right)^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\ &\leq 2C_T \|f_1 - f_2\|_{V_k}^2 + 2C_H \|g(\partial_k f_1) - g(\partial_k f_2)\|_{L^2_{\eta_{\leq k}}}^2 \\ &\leq 2C_T \|f_1 - f_2\|_{V_k}^2 + 2C_H L^2 \|\partial_k f_1 - \partial_k f_2\|_{L^2_{\eta_{\leq k}}}^2 \\ &\leq 2(C_T + C_H L^2) \|f_1 - f_2\|_{V_k}^2, \end{aligned} \quad (44)$$



for any  $f_1, f_2 \in V_k$ . Furthermore, using the Lipschitz property of  $g$  we have

$$\begin{aligned} \|\partial_k \mathcal{R}_k(f_1) - \partial_k \mathcal{R}_k(f_2)\|_{L^2_{\eta_{\leq k}}}^2 &= \int \left( g(\partial_k f_1(\mathbf{x}_{<k}, t)) - g(\partial_k f_2(\mathbf{x}_{<k}, t)) \right)^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\ &\leq L^2 \int \left( \partial_k f_1(\mathbf{x}_{<k}, t) - \partial_k f_2(\mathbf{x}_{<k}, t) \right)^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\ &\leq L^2 \|f_1 - f_2\|_{V_k}^2. \end{aligned} \quad (45)$$

Combining (44) with (45), we obtain (18) with  $C = \sqrt{2(C_T + C_H L^2) + L^2}$ .

It remains to show that  $\|\mathcal{R}_k(f)\|_{V_k} < \infty$  for any  $f \in V_k$ . Letting with  $f_1 = f$  and  $f_2 = 0$  in (18), the triangle inequality yields

$$\|\mathcal{R}_k(f)\|_{V_k} \leq \|\mathcal{R}_k(0)\|_{V_k} + C\|f\|_{V_k}$$

Because  $\mathcal{R}_k(0)$  is the affine function  $\mathbf{x} \mapsto g(0)x_k$ , we have that  $\|\mathcal{R}_k(0)\|_{L^2_{\eta_{\leq k}}}^2 = g(0)^2 \int x_k^2 \eta(\mathbf{x}) d\mathbf{x}$  and  $\|\partial_k \mathcal{R}_k(0)\|_{L^2_{\eta_{\leq k}}}^2 = g(0)^2$  are finite and so is  $\|\mathcal{R}_k(0)\|_{V_k}$ . Thus,  $\mathcal{R}_k(f) \in V_k$  for all  $f \in V_k$ .  $\square$

#### A.5. Proof of Proposition 4

*Proof.* For any  $f \in V_k$  we have

$$\begin{aligned} |\mathcal{L}_k(f)| &= \left| \int \left( \frac{1}{2} \mathcal{R}_k(f)^2 - \log(\partial_k \mathcal{R}_k(f)) \right) d\pi \right| \\ &\stackrel{(19)}{\leq} \frac{C_\pi}{2} \|\mathcal{R}_k(f)\|_{L^2_{\eta_{\leq k}}}^2 + C_\pi \int |\log(g(\partial_k f))| d\eta_{\leq k} \\ &\leq \frac{C_\pi}{2} \|\mathcal{R}_k(f)\|_{L^2_{\eta_{\leq k}}}^2 + C_\pi \int |\log(g(0))| + |\log(g(\partial_k f)) - \log(g(0))| d\eta_{\leq k} \\ &\stackrel{(21)}{\leq} \frac{C_\pi}{2} \|\mathcal{R}_k(f)\|_{L^2_{\eta_{\leq k}}}^2 + C_\pi |\log(g(0))| + C_\pi L \int |\partial_k f - 0| d\eta_{\leq k} \\ &\leq \frac{C_\pi}{2} \|\mathcal{R}_k(f)\|_{L^2_{\eta_{\leq k}}}^2 + C_\pi |\log(g(0))| + C_\pi L \|f\|_{V_k}^2. \end{aligned}$$

Because Proposition 3 ensures  $\mathcal{R}_k(f) \in V_k \subset L^2_{\eta_{\leq k}}$ , we have that  $\mathcal{L}_k(f)$  is finite for any  $f \in V_k$ . Now, for any  $f_1, f_2 \in V_k$ , we can write

$$\begin{aligned} |\mathcal{L}_k(f_1) - \mathcal{L}_k(f_2)| &= \left| \int \left( \frac{1}{2} \mathcal{R}_k(f_1)^2 - \frac{1}{2} \mathcal{R}_k(f_2)^2 - \log(\partial_k \mathcal{R}_k(f_1)) + \log(\partial_k \mathcal{R}_k(f_2)) \right) d\pi \right| \\ &\stackrel{(19)}{\leq} C_\pi \int \frac{1}{2} \left| \mathcal{R}_k(f_1)^2 - \mathcal{R}_k(f_2)^2 \right| + \left| \log(g(\partial_k f_1)) - \log(g(\partial_k f_2)) \right| d\eta \\ &\stackrel{(21)}{\leq} \frac{C_\pi}{2} \|\mathcal{R}_k(f_1) + \mathcal{R}_k(f_2)\|_{L^2_{\eta_{\leq k}}} \|\mathcal{R}_k(f_1) - \mathcal{R}_k(f_2)\|_{L^2_{\eta_{\leq k}}} + C_\pi L \|\partial_k f_1 - \partial_k f_2\|_{L^2_{\eta_{\leq k}}} \\ &\stackrel{(18)}{\leq} C_\pi \frac{\|\mathcal{R}_k(f_1)\|_{L^2_{\eta_{\leq k}}} + \|\mathcal{R}_k(f_2)\|_{L^2_{\eta_{\leq k}}}}{2} C \|f_1 - f_2\|_{V_k} + C_\pi L \|f_1 - f_2\|_{V_k}. \end{aligned}$$

This shows that  $\mathcal{L}_k: V_k \rightarrow \mathbb{R}$  is continuous. To show that  $\mathcal{L}_k$  is differentiable, we let  $f, \varepsilon \in V_k$  so that

$$\begin{aligned} \mathcal{L}_k(f + \varepsilon) &= \int \left( \frac{1}{2} \mathcal{R}_k(f + \varepsilon)^2 - \log(\partial_k \mathcal{R}_k(f + \varepsilon)) \right) d\pi \\ &= \int \left( \frac{1}{2} \left( f(\mathbf{x}_{<k}, 0) + \varepsilon(\mathbf{x}_{<k}, 0) + \int_0^{x_k} g(\partial_k(f + \varepsilon)(\mathbf{x}_{<k}, t)) dt \right)^2 - \log \circ g(\partial_k(f + \varepsilon)(\mathbf{x})) \right) \pi(\mathbf{x}) d\mathbf{x} \\ &= \int \left( \frac{1}{2} \mathcal{R}_k(f)(\mathbf{x})^2 + \mathcal{R}_k(f)(\mathbf{x}) \left( \varepsilon(\mathbf{x}_{<k}, 0) + \int_0^{x_k} g'(\partial_k f(\mathbf{x}_{<k}, t)) \partial_k \varepsilon(\mathbf{x}_{<k}, t) dt \right) \right) \pi(\mathbf{x}) d\mathbf{x} \\ &\quad - \int \log \circ g(\partial_k f) + (\log \circ g)'(\partial_k f) \partial_k \varepsilon d\pi + \mathcal{O}(\|\varepsilon\|_{V_k}^2) \\ &= \mathcal{L}_k(f) + \ell(\varepsilon) + \mathcal{O}(\|\varepsilon\|_{V_k}^2) \end{aligned}$$

where  $\ell: V_k \rightarrow \mathbb{R}$  is the linear form defined by

$$\ell(\varepsilon) = \int \mathcal{R}_k(f)(\mathbf{x}) \left( \varepsilon(\mathbf{x}_{<k}, 0) + \int_0^{x_k} g'(\partial_k f(\mathbf{x}_{<k}, t)) \partial_k \varepsilon(\mathbf{x}_{<k}, t) dt \right) - (\log \circ g)'(\partial_k f(\mathbf{x})) \partial_k \varepsilon(\mathbf{x}) \pi(\mathbf{x}) d\mathbf{x}.$$

If  $\ell$  is continuous, meaning if there exists a constant  $C_\ell$  such that  $|\ell(\varepsilon)| \leq C_\ell \|\varepsilon\|_{V_k}$  for any  $\varepsilon \in V_k$ , then the Riesz representation theorem states that there exists a vector  $\nabla \mathcal{L}_k(f) \in V_k$  such that  $\ell(\varepsilon) = \langle \nabla \mathcal{L}_k(f), \varepsilon \rangle_{V_k}$ . This proves  $\mathcal{L}_k$  is differentiable everywhere.

To show that  $\ell$  is continuous, we write

$$\begin{aligned} |\ell(\varepsilon)| &\stackrel{(19)}{\leq} C_\pi \int \left| \mathcal{R}_k(f)(\mathbf{x}) \left( \varepsilon(\mathbf{x}_{<k}, 0) + \int_0^{x_k} g'(\partial_k f(\mathbf{x}_{<k}, t)) \partial_k \varepsilon(\mathbf{x}_{<k}, t) dt \right) \right| \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\ &\quad + C_\pi \int \left| (\log \circ g)'(\partial_k f(\mathbf{x})) \partial_k \varepsilon(\mathbf{x}) \right| \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\ &\stackrel{(21)}{\leq} C_\pi \|\mathcal{R}_k(f)\|_{L_{\eta_{\leq k}}^2} \sqrt{\int \left| \varepsilon(\mathbf{x}_{<k}, 0) + \int_0^{x_k} g'(\partial_k f(\mathbf{x}_{<k}, t)) \partial_k \varepsilon(\mathbf{x}_{<k}, t) dt \right|^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x}} + C_\pi L \|\partial_k \varepsilon\|_{L_{\eta_{\leq k}}^2} \\ &\stackrel{(20)}{\leq} C_\pi \|\mathcal{R}_k(f)\|_{L_{\eta_{\leq k}}^2} \sqrt{2C_T \|\varepsilon\|_{V_k}^2 + 2C_H L^2 \int |\partial_k \varepsilon(\mathbf{x})|^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x}} + C_\pi L \|\partial_k \varepsilon\|_{L_{\eta_{\leq k}}^2} \\ &\leq C_\pi \left( \|\mathcal{R}_k(f)\|_{L_{\eta_{\leq k}}^2} \sqrt{2C_T + 2C_H L^2} + L \right) \|\varepsilon\|_{V_k}, \end{aligned}$$

where the second last inequality also uses Proposition 2 and Lemma 12. This concludes the proof.  $\square$

#### A.6. Proof of the local Lipschitz regularity (23)

**Proposition 14.** *In addition to the assumptions of Theorem 4, we further assume there exists a constant  $L < \infty$  such that for all  $\xi, \xi' \in \mathbb{R}$  we have*

$$|g'(\xi) - g'(\xi')| \leq L|\xi - \xi'| \tag{46}$$

$$|(\log \circ g)'(\xi) - (\log \circ g)'(\xi')| \leq L|\xi - \xi'|. \tag{47}$$

Then there exists  $M < \infty$  such that

$$\|\nabla \mathcal{L}_k(f_1) - \nabla \mathcal{L}_k(f_2)\|_{V_k} \leq M(1 + \|\mathcal{R}_k(f_2)\|_{V_k}) \|f_1 - f_2\|_{\bar{V}_k},$$

for any  $f_1, f_2 \in \bar{V}_k$ , where  $\bar{V}_k = \{f \in V_k, \partial_k f \in L^\infty\}$  is the space endowed with the norm  $\|f\|_{\bar{V}_k} = \|f\|_{V_k} + \|\partial_k f\|_{L^\infty}$ .

*Proof.* Recall the definition (22) of  $\nabla \mathcal{L}_k(f)$

$$\begin{aligned} \langle \nabla \mathcal{L}_k(f), \varepsilon \rangle_{V_k} &= \int \mathcal{R}_k(f)(\mathbf{x}) \left( \varepsilon(\mathbf{x}_{<k}, 0) + \int_0^{x_k} g'(\partial_k f(\mathbf{x}_{<k}, t)) \partial_k \varepsilon(\mathbf{x}_{<k}, t) dt \right) \pi(\mathbf{x}) d\mathbf{x} \\ &\quad - \int (\log \circ g)'(\partial_k f(\mathbf{x})) \partial_k \varepsilon(\mathbf{x}) \pi(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

Then for any  $f_1, f_2 \in \bar{V}_k$ , we can write

$$\langle \nabla \mathcal{L}_k(f_1) - \nabla \mathcal{L}_k(f_2), \varepsilon \rangle_{V_k} = A + B + C + D,$$

where

$$\begin{aligned} A &= \int \left( \mathcal{R}_k(f_1)(\mathbf{x}) - \mathcal{R}_k(f_2)(\mathbf{x}) \right) \varepsilon(\mathbf{x}_{<k}, 0) \pi(\mathbf{x}) d\mathbf{x} \\ B &= \int \left( \mathcal{R}_k(f_1)(\mathbf{x}) - \mathcal{R}_k(f_2)(\mathbf{x}) \right) \left( \int_0^{x_k} g'(\partial_k f_1(\mathbf{x}_{<k}, t)) \partial_k \varepsilon(\mathbf{x}_{<k}, t) dt \right) \pi(\mathbf{x}) d\mathbf{x} \\ C &= \int \mathcal{R}_k(f_2)(\mathbf{x}) \left( \int_0^{x_k} \left( g'(\partial_k f_1(\mathbf{x}_{<k}, t)) - g'(\partial_k f_2(\mathbf{x}_{<k}, t)) \right) \partial_k \varepsilon(\mathbf{x}_{<k}, t) dt \right) \pi(\mathbf{x}) d\mathbf{x} \\ D &= \int \left( (\log \circ g)'(\partial_k f_1(\mathbf{x})) - (\log \circ g)'(\partial_k f_2(\mathbf{x})) \right) \partial_k \varepsilon(\mathbf{x}) \pi(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

For the first term  $A$  we write

$$\begin{aligned} |A| &\stackrel{(19)}{\leq} C_\pi \int \left| \mathcal{R}_k(f_1)(\mathbf{x}) - \mathcal{R}_k(f_2)(\mathbf{x}) \right| |\varepsilon(\mathbf{x}_{<k}, 0)| \eta(\mathbf{x}) d\mathbf{x} \\ &\leq C_\pi \|\mathcal{R}_k(f_1) - \mathcal{R}_k(f_2)\|_{V_k} \left( \int |\varepsilon(\mathbf{x}_{<k}, 0)|^2 \eta(\mathbf{x}) d\mathbf{x} \right)^{1/2} \\ &\stackrel{(15)}{\leq} C_\pi \sqrt{C_T} \|\mathcal{R}_k(f_1) - \mathcal{R}_k(f_2)\|_{V_k} \|\varepsilon\|_{V_k} \\ &\stackrel{(18)}{\leq} C_\pi \sqrt{C_T} C \|f_1 - f_2\|_{V_k} \|\varepsilon\|_{V_k}. \end{aligned}$$

For the second term  $B$  we write

$$\begin{aligned} |B| &\stackrel{(19)}{\leq} C_\pi \|\mathcal{R}_k(f_1) - \mathcal{R}_k(f_2)\|_{V_k} \left( \int \left( \int_0^{x_k} g'(\partial_k f_1(\mathbf{x}_{<k}, t)) \partial_k \varepsilon(\mathbf{x}_{<k}, t) dt \right)^2 \eta(\mathbf{x}) d\mathbf{x} \right)^{1/2} \\ &\stackrel{(39)}{\leq} C_\pi \sqrt{C_H} C \|f_1 - f_2\|_{V_k} \left( \int \left( g'(\partial_k f_1(\mathbf{x}_{\leq k})) \partial_k \varepsilon(\mathbf{x}_{\leq k}) \right)^2 \eta(\mathbf{x}) d\mathbf{x} \right)^{1/2} \\ &\stackrel{(20)}{\leq} C_\pi \sqrt{C_H} C L \|f_1 - f_2\|_{V_k} \left( \int \left( \partial_k \varepsilon(\mathbf{x}_{\leq k}) \right)^2 \eta(\mathbf{x}) d\mathbf{x} \right)^{1/2} \\ &\leq C_\pi \sqrt{C_H} C L \|f_1 - f_2\|_{V_k} \|\varepsilon\|_{V_k}. \end{aligned}$$

For the third term  $C$  we write

$$\begin{aligned}
|C| &\stackrel{(19)}{\leq} C_\pi \|\mathcal{R}_k(f_2)\|_{V_k} \left( \int \left( \int_0^{x_k} \left( g'(\partial_k f_1(\mathbf{x}_{<k}, t)) - g'(\partial_k f_2(\mathbf{x}_{<k}, t)) \right) \partial_k \varepsilon(\mathbf{x}_{<k}, t) dt \right)^2 \eta(\mathbf{x}) d\mathbf{x} \right)^{1/2} \\
&\stackrel{(39)}{\leq} C_\pi \sqrt{C_H} \|\mathcal{R}_k(f_2)\|_{V_k} \left( \int \left( \left( g'(\partial_k f_1(\mathbf{x}_{\leq k})) - g'(\partial_k f_2(\mathbf{x}_{\leq k})) \right) \partial_k \varepsilon(\mathbf{x}_{\leq k}) \right)^2 \eta(\mathbf{x}) d\mathbf{x} \right)^{1/2} \\
&\leq C_\pi \sqrt{C_H} \|\mathcal{R}_k(f_2)\|_{V_k} \left( \text{ess sup} \left| g' \circ \partial_k f_1 - g' \circ \partial_k f_2 \right| \right) \left( \int (\partial_k \varepsilon(\mathbf{x}_{\leq k}))^2 \eta(\mathbf{x}) d\mathbf{x} \right)^{1/2} \\
&\stackrel{(46)}{\leq} C_\pi \sqrt{C_H} L \|\mathcal{R}_k(f_2)\|_{V_k} \left( \text{ess sup} \left| \partial_k f_1 - \partial_k f_2 \right| \right) \|\varepsilon\|_{V_k} \\
&\leq C_\pi \sqrt{C_H} L \|\mathcal{R}_k(f_2)\|_{V_k} \|f_1 - f_2\|_{\bar{V}_k} \|\varepsilon\|_{V_k}.
\end{aligned}$$

For the last term  $D$  we write

$$\begin{aligned}
|D| &\stackrel{(19)}{\leq} C_\pi \left( \int \left( (\log \circ g)'(\partial_k f_1(\mathbf{x})) - (\log \circ g)'(\partial_k f_2(\mathbf{x})) \right)^2 \eta(\mathbf{x}) d\mathbf{x} \right)^{1/2} \|\varepsilon\|_{V_k} \\
&\stackrel{(47)}{\leq} C_\pi L \left( \int \left( \partial_k f_1(\mathbf{x}) - \partial_k f_2(\mathbf{x}) \right)^2 \eta(\mathbf{x}) d\mathbf{x} \right)^{1/2} \|\varepsilon\|_{V_k} \\
&\leq C_\pi L \|f_1 - f_2\|_{V_k} \|\varepsilon\|_{V_k}.
\end{aligned}$$

Thus, because  $\|f_1 - f_2\|_{V_k} \leq \|f_1 - f_2\|_{\bar{V}_k}$  we obtain

$$\begin{aligned}
\frac{|\langle \nabla \mathcal{L}_k(f_1) - \nabla \mathcal{L}_k(f_2), \varepsilon \rangle_{V_k}|}{\|\varepsilon\|_{V_k}} &\leq C_\pi \left( \sqrt{C_T} C + \sqrt{C_H} C L + \sqrt{C_H} L \|\mathcal{R}_k(f_2)\|_{V_k} + L \right) \|f_1 - f_2\|_{\bar{V}_k} \\
&\leq M (1 + \|\mathcal{R}_k(f_2)\|_{V_k}) \|f_1 - f_2\|_{\bar{V}_k},
\end{aligned}$$

where

$$M = C_\pi \max\{\sqrt{C_T} C + \sqrt{C_H} C L + L; \sqrt{C_H} L\}.$$

This concludes the proof.  $\square$

### A.7. Proof of Proposition 6

*Proof.* To show that  $\mathcal{R}_k(V_k) = \{\mathcal{R}(f) : f \in V_k\}$  is convex, let  $f_1, f_2 \in V_k$  and  $0 \leq \alpha \leq 1$ . We need to show that there exists  $f_\alpha \in V_k$  such that  $\mathcal{R}(f_\alpha) = S_\alpha$  where

$$S_\alpha := \alpha \mathcal{R}_k(f_1) + (1 - \alpha) \mathcal{R}_k(f_2).$$

Let

$$f_\alpha(\mathbf{x}_{\leq k}) := \mathcal{R}^{-1}(S_\alpha)(\mathbf{x}_{\leq k}) = S_\alpha(\mathbf{x}_{<k}, 0) + \int_0^{x_k} g^{-1}(\partial_k S_\alpha(\mathbf{x}_{<k}, t)) dt.$$

It remains to show that  $f_\alpha \in V_k$ , meaning that  $f_\alpha \in L^2_{\eta_{\leq k}}$  and  $\partial_k f_\alpha \in L^2_{\eta_{\leq k}}$ . By convexity of  $\xi \mapsto g^{-1}(\xi)^2$ , we have

$$\begin{aligned} \|\partial_k f_\alpha\|_{L^2_{\eta_{\leq k}}}^2 &= \int g^{-1}(\alpha \partial_k \mathcal{R}_k(f_1) + (1-\alpha) \partial_k \mathcal{R}_k(f_2))^2 d\eta_{\leq k} \\ &= \int g^{-1}(\alpha g(\partial_k f_1) + (1-\alpha) g(\partial_k f_2))^2 d\eta_{\leq k} \\ &\leq \int \alpha g^{-1}(g(\partial_k f_1))^2 + (1-\alpha) g^{-1}(g(\partial_k f_2))^2 d\eta_{\leq k} \\ &= \alpha \|\partial_k f_1\|_{L^2_{\eta_{\leq k}}}^2 + (1-\alpha) \|\partial_k f_2\|_{L^2_{\eta_{\leq k}}}^2. \end{aligned} \quad (48)$$

Thus  $\partial_k f_\alpha \in L^2_{\eta_{\leq k}}$ . Furthermore we have

$$\begin{aligned} \|f_\alpha\|_{L^2_{\eta_{\leq k}}}^2 &= \int \left( S_\alpha(\mathbf{x}_{<k}, 0) + \int_0^{x_k} g^{-1}(\partial_k S_\alpha(\mathbf{x}_{<k}, t)) dt \right)^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\ &\leq 2 \int S_\alpha(\mathbf{x}_{<k}, 0)^2 \eta_{<k}(\mathbf{x}_{<k}) d\mathbf{x} + 2 \int \left( \int_0^{x_k} g^{-1}(\partial_k S_\alpha(\mathbf{x}_{<k}, t)) dt \right)^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

To show that the above quantity is finite, Proposition 2 permits us to write

$$\begin{aligned} \int S_\alpha(\mathbf{x}_{<k}, 0)^2 \eta_{<k}(\mathbf{x}_{<k}) d\mathbf{x} &= \int \left( \alpha f_1(\mathbf{x}_{<k}, 0) + (1-\alpha) f_2(\mathbf{x}_{<k}, 0) \right)^2 \eta_{<k}(\mathbf{x}_{<k}) d\mathbf{x} \\ &\leq C_T \|\alpha f_1 + (1-\alpha) f_2\|_{V_k}^2, \end{aligned}$$

which is finite. Finally, because  $g^{-1}(\partial_k S_\alpha) = \partial_k f_\alpha \in L^2_{\eta_{\leq k}}$  by (48), Lemma 12 yields

$$\begin{aligned} \int \left( \int_0^{x_k} g^{-1}(\partial_k S_\alpha(\mathbf{x}_{<k}, t)) dt \right)^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} &\leq C_H \int g^{-1}(\partial_k S_\alpha(\mathbf{x}_{\leq k}))^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\ &= C_H \|\partial_k f_\alpha\|_{L^2_{\eta_{\leq k}}}^2, \end{aligned}$$

which is finite. We deduce that  $f_\alpha \in L^2_{\eta_{\leq k}}$  and therefore that  $f_\alpha \in V_k$ .  $\square$

### A.8. Proof of Proposition 5

*Proof.* Let  $s_1, s_2 \in V_k$  be strictly increasing functions with respect to  $x_k$  that satisfy  $\partial_k s_i(\mathbf{x}_{\leq k}) \geq c$  for  $i = 1, 2$  and all  $\mathbf{x}_{\leq k} \in \mathbb{R}^k$ . By the Lipschitz property of  $g^{-1}$  on the domain  $[c, \infty)$  with constant  $L_c$ , we can write

$$\begin{aligned} \|\partial_k \mathcal{R}_k^{-1}(s_1) - \partial_k \mathcal{R}_k^{-1}(s_2)\|_{L^2_{\eta_{\leq k}}}^2 &= \int (g^{-1}(\partial_k s_1(\mathbf{x}_{\leq k})) - g^{-1}(\partial_k s_2(\mathbf{x}_{\leq k})))^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\ &\leq L_c^2 \int (\partial_k s_1(\mathbf{x}_{\leq k}) - \partial_k s_2(\mathbf{x}_{\leq k}))^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\ &\leq L_c^2 \|s_1 - s_2\|_{V_k}^2. \end{aligned} \quad (49)$$

Applying Proposition 2 to  $s_1, s_2 \in V_k$  and Lemma 12 to  $\partial_k \mathcal{R}_k^{-1}(s_i) = g^{-1}(\partial_k s_i) \in L^2_{\eta_{\leq k}}$  for  $i = 1, 2$

we have

$$\begin{aligned}
\|\mathcal{R}_k^{-1}(s_1) - \mathcal{R}_k^{-1}(s_2)\|_{L_{\eta_{\leq k}}^2}^2 &\leq 2 \int (s_1(\mathbf{x}_{<k}, 0) - s_2(\mathbf{x}_{<k}, 0))^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\
&\quad + 2 \int \left( \int_0^{x_k} g^{-1}(\partial_k s_1) - g^{-1}(\partial_k s_2) dt \right)^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\
&\leq 2C_T \|s_1 - s_2\|_{V_k}^2 \\
&\quad + 2C_H \int (g^{-1}(\partial_k s_1(\mathbf{x}_{\leq k})) - g^{-1}(\partial_k s_2(\mathbf{x}_{\leq k})))^2 \eta_{\leq k}(\mathbf{x}) d\mathbf{x} \\
&\leq (2C_T + 2C_H L_c^2) \|s_1 - s_2\|_{V_k}^2,
\end{aligned} \tag{50}$$

where the last inequality follows from (49).

It remains to show that  $\|\mathcal{R}_k^{-1}(s)\|_{V_k} < \infty$  for any  $s \in V_k$  such that  $\text{ess inf } \partial_k s > 0$ . Letting  $s_1 = s$  and  $s_2 = g(0)x_k$ , the triangle inequality combined with (25) yields

$$\|\mathcal{R}_k^{-1}(s)\|_{V_k} \leq \|\mathcal{R}_k^{-1}(g(0)x_k)\|_{V_k} + C_c \|s - g(0)x_k\|_{V_k}.$$

The function  $\mathcal{R}_k^{-1}(g(0)x_k)$  is zero. Therefore,  $\|\mathcal{R}_k^{-1}(s)\|_{V_k} \leq C_c \|s - g(0)x_k\|_{V_k} \leq C_c (\|s\|_{V_k} + \|g(0)x_k\|_{V_k})$ . For a linear function,  $\|g(0)x_k\|_{V_k}^2 = \|g(0)x_k\|_{L_{\eta_{\leq k}}^2}^2 + \|g(0)\|_{L_{\eta_{\leq k}}^2}^2 = 2g(0)^2$  is finite, and so  $\|\mathcal{R}_k^{-1}(s)\|_{V_k} < \infty$  for  $s \in V_k$ . Furthermore, if  $\partial_k s \geq c > 0$ , then  $\partial_k \mathcal{R}_k^{-1}(s) = g^{-1}(\partial_k s) \geq g^{-1}(c) > -\infty$  and so  $\text{ess inf } \mathcal{R}_k^{-1}(s) > -\infty$ .  $\square$

### A.9. Proof for the KR rearrangement

*Proof.* Let  $S_{\text{KR},k}$  be the  $k$ th component of the KR rearrangement, given by composing the inverse CDF of the standard Gaussian marginal  $F_{\eta,k}(x_k)$  with the CDF of the target's  $k$ th marginal conditional  $F_{\pi_k}(x_k | \mathbf{x}_{<k})$ . That is,

$$S_{\text{KR},k}(\mathbf{x}_{\leq k}) = F_{\eta_k}^{-1} \circ F_{\pi_k}(x_k | \mathbf{x}_{<k}). \tag{51}$$

The goal is to show  $S_{\text{KR},k} \in V_k$ , that is,  $S_{\text{KR},k} \in L_{\eta_{\leq k}}^2$  and  $\partial_k S_{\text{KR},k} \in L_{\eta_{\leq k}}^2$ .

First we show  $S_{\text{KR},k} \in L_{\eta_{\leq k}}^2$ . From condition (29), we have  $F_{\eta_k}^{-1}(c_k F_{\eta_k}(x_k)) \leq S_{\text{KR},k}(x_k | \mathbf{x}_{<k}) \leq F_{\eta_k}^{-1}(C_k F_{\eta_k}(x_k))$  so that

$$S_{\text{KR},k}(x_k | \mathbf{x}_{<k})^2 \leq \max\{F_{\eta_k}^{-1}(c_k F_{\eta_k}(x_k))^2; F_{\eta_k}^{-1}(C_k F_{\eta_k}(x_k))^2\}, \tag{52}$$

for all  $\mathbf{x}_{<k} \in \mathbb{R}^{k-1}$ . To show that  $S_{\text{KR},k} \in L_{\eta_{\leq k}}^2$ , it is sufficient to prove that any function of the form  $x_k \mapsto F_{\eta_k}^{-1}(CF_{\eta_k}(x_k))$  is in  $L_{\eta_{\leq k}}^2$  for any  $C > 0$ . From Theorems 1 and 2 in [10], there exists strictly positive constants  $\alpha_i, \beta_i > 0$  for  $i = 1, 2$  such that

$$\alpha_1 \exp(-\beta_1 x_k^2) \leq F_{\eta_k}(x_k) \leq \alpha_2 \exp(-\beta_2 x_k^2), \tag{53}$$

for  $x_k < 0$ . With a change of variable  $u = F_{\eta_k}(x_k)$  we obtain  $F_{\eta_k}^{-1}(u)^2 \leq 1/\beta_2 \log(\alpha_2/u)$  for all  $u < F_{\eta_k}(0) = 1/2$ . Letting  $u = CF_{\eta_k}(x_k)$  yields

$$\begin{aligned}
F_{\eta_k}^{-1}(CF_{\eta_k}(x_k))^2 &\leq \frac{1}{\beta_2} \log \left( \frac{\alpha_2}{CF_{\eta_k}(x_k)} \right), \\
&\stackrel{(53)}{\leq} \frac{1}{\beta_2} \log \left( \frac{\alpha_2}{C\alpha_1 \exp(-\beta_1 x_k^2)} \right) \\
&= \frac{1}{\beta_2} \log \left( \frac{\alpha_2}{C\alpha_1} \right) + \frac{\beta_1}{\beta_2} x_k^2
\end{aligned}$$

for all  $x_k < \min\{F_{\eta_k}^{-1}(1/(2C)), 0\}$ . Using the same argument, we obtain a similar bound on  $F_{\eta_k}^{-1}(CF_{\eta_k}(x_k))^2$  for all  $x_k$  larger than a certain value. Together with the continuity of  $x_k \mapsto F_{\eta_k}^{-1}(CF_{\eta_k}(x_k))^2$  these bounds ensure that  $x_k \mapsto F_{\eta_k}^{-1}(CF_{\eta_k}(x_k))^2$  is in  $L^2_{\eta_{\leq k}}$  for any  $C$ . Then  $S_{\text{KR},k} \in L^2_{\eta_{\leq k}}$ . Furthermore, we have  $S_{\text{KR},k}(\mathbf{x}_{\leq k}) = \mathcal{O}(x_k)$  as  $|x_k| \rightarrow \infty$ .

Now we show that  $\partial_k S_{\text{KR},k} \in L^2_{\eta_{\leq k}}$  by showing  $\partial_k S_{\text{KR},k}$  is a continuous and bounded function. From the absolute continuity of  $\mu$  and  $\nu$ , we have that

$$\partial_k S_{\text{KR},k}(\mathbf{x}_{\leq k}) = \frac{\pi_k(x_k|\mathbf{x}_{<k})}{\eta_k(S_{\text{KR},k}(\mathbf{x}_{\leq k}))} = \frac{\pi_k(F_{\pi_k}^{-1}(F_{\pi_k}(x_k|\mathbf{x}_{<k})|\mathbf{x}_{<k})|\mathbf{x}_{<k})}{\eta_k(F_{\eta_k}^{-1} \circ F_{\pi_k}(x_k|\mathbf{x}_{<k}))}, \quad (54)$$

is continuous, where  $F_{\pi_k}^{-1}(\cdot|\mathbf{x}_{<k})$  denotes the inverse of the map  $x_k \mapsto F_{\pi_k}(x_k|\mathbf{x}_{<k})$  for each  $\mathbf{x}_{<k} \in \mathbb{R}^{k-1}$ . Hence, it is sufficient to show that  $\partial_k S_{\text{KR},k}$  goes to a finite limit as  $|x_k| \rightarrow \infty$ . For the left-hand limit, we can write

$$\begin{aligned} \lim_{x_k \rightarrow -\infty} \partial_k S_{\text{KR},k}(\mathbf{x}_{\leq k}) &= \lim_{u \rightarrow 0^+} \frac{\pi_k(F_{\pi_k}^{-1}(u|\mathbf{x}_{<k})|\mathbf{x}_{<k})}{\eta_k(F_{\eta_k}^{-1}(u))} \\ &= \lim_{u \rightarrow 0^+} \frac{(F_{\eta,k}^{-1})'(u)}{(F_{\pi_k}^{-1})'(u|\mathbf{x}_{<k})} \\ &= \lim_{u \rightarrow 0^+} \frac{F_{\eta_k}^{-1}(u)}{F_{\pi_k}^{-1}(u|\mathbf{x}_{<k})}, \end{aligned} \quad (55)$$

where in the second equality we used the implicit function theorem and the third equality follows from l'Hôpital's rule. To analyze the ratio  $F_{\eta_k}^{-1}(u)/F_{\pi_k}^{-1}(u|\mathbf{x}_{<k})$ , we combine the lower bound in (29) and the bounds in (53) to get

$$\frac{F_{\eta_k}^{-1}(u)}{F_{\pi_k}^{-1}(u|\mathbf{x}_{<k})} \leq \frac{F_{\eta_k}^{-1}(u)}{F_{\eta_k}^{-1}(u/c_k)} \leq \sqrt{\frac{\beta_2(\log \alpha_1 - \log(u))}{\beta_1(\log \alpha_2 - \log(u/c_k))}}.$$

Similarly, from the upper bound in (29) and the bounds in (53), we have

$$\frac{F_{\eta_k}^{-1}(u)}{F_{\pi_k}^{-1}(u|\mathbf{x}_{<k})} \geq \frac{F_{\eta_k}^{-1}(u)}{F_{\eta_k}^{-1}(u/C_k)} \geq \sqrt{\frac{\beta_1(\log \alpha_2 - \log(u))}{\beta_2(\log \alpha_1 - \log(u/C_k))}}.$$

Thus,  $\partial_k S_{\text{KR},k}(\mathbf{x}_{\leq k}) = \mathcal{O}(1)$  as  $x_k \rightarrow -\infty$ , and we have  $\partial_k S_{\text{KR},k} \in L^2_{\eta_k}$ .

Lastly, taking the limit in (55) we have  $\lim_{x_k \rightarrow -\infty} \partial_k S_{\text{KR},k}(\mathbf{x}_{\leq k}) \geq \sqrt{\beta_1/\beta_2}$ . For a target distribution  $\pi$  with full support, all marginal conditional densities satisfy  $\pi_k(x_k|\mathbf{x}_{<k}) > 0$  for each  $\mathbf{x}_{<k} \in \mathbb{R}^k$ . Given that the  $\partial_k S_{\text{KR},k}$  does not approach zero as  $|x_k| \rightarrow \infty$ , we can find a strictly positive constant  $c_k > 0$  such that  $\partial_k S_{\text{KR},k}(\mathbf{x}_{\leq k}) \geq c_k$  for all  $\mathbf{x}_{\leq k} \in \mathbb{R}^k$ . This shows that  $\text{ess inf } \partial_k S_{\text{KR},k} > 0$ .  $\square$

## Appendix B: Multi-index refinement for the wavelet basis

In this section we show how to greedily enrich the index set  $\Lambda_t$  for a one-dimensional wavelet basis parameterized by the tuple of indices  $(l, q)$  representing the level  $l$  and translation  $q$  of each wavelet  $\psi_{(l,q)}$ . To define the allowable indices, we construct a binary tree where each node is indexed by  $(l, q)$  and has two children with indices  $(l+1, 2q)$  and  $(l+1, 2q+1)$ . The root of the tree has index  $(0, 0)$  and corresponds to the mother wavelet  $\psi$ . Analogously to the downward closed property for



polynomial indices, we only add nodes to the tree (i.e., indices in  $\Lambda_t$ ) if its parents have already been added. Given any set  $\Lambda_t$ , we define its reduced margin as

$$\Lambda_t^{\text{RM}} = \left\{ \alpha = (l, q) \notin \Lambda_t \text{ such that } \begin{array}{ll} (l-1, q/2) \in \Lambda_t & \text{for odd } q \\ (l-1, (q-1)/2) \in \Lambda_t & \text{for even } q \end{array} \right\}.$$

Then, the ATM algorithm with a wavelet basis follows from Algorithm 1 with this construction for the reduced margin at each iteration.

## Appendix C: Architecture details of alternative methods

In this section we present the details of the alternative methods to ATM that we consider in Section 5.

For each normalizing flow model, we use the recommended stochastic gradient descent optimizer with a learning rate of  $10^{-3}$ . We partition 10% of the samples in each training set to be validation samples and use the remaining samples for training the model. We select the optimal hyper-parameters for each dataset by fitting the density with the training data and choosing the parameters that minimize the negative log-likelihood of the approximate density on the validation samples. We also use the validation samples to set the termination criteria during the optimization.

We follow the implementation of [50] to define the architectures of these models. The hyper-parameters we consider for the neural networks in the MDN and NF models are: 2 hidden layers, 32 hidden units in each layer,  $\{5, 10, 20, 50, 100\}$  centers or flows, weight normalization, and a dropout probability of  $\{0, 0.2\}$  for regularizing the neural networks during training. For CKDE and NKDE we select the bandwidth of the kernel estimators using 5-fold cross-validation.

## References

- [1] AMBROGIONI, L., GÜÇLÜ, U., VAN GERVEN, M. A. and MARIS, E. (2017). The kernel mixture network: A nonparametric method for conditional density estimation of continuous random variables. *arXiv preprint arXiv:1705.07111*.
- [2] ANDERES, E. and CORAM, M. (2012). A general spline representation for nonparametric and semiparametric density estimates using diffeomorphisms. *arXiv preprint arXiv:1205.5314*.
- [3] BAPTISTA, R., MARZOUK, Y., MORRISON, R. E. and ZAHM, O. (2021). Learning non-Gaussian graphical models via Hessian scores and triangular transport. *arXiv preprint arXiv:2101.03093*.
- [4] BERTSEKAS, D. P. (1997). Nonlinear programming. *Journal of the Operational Research Society* **48** 334–334.
- [5] BIGONI, D., MARZOUK, Y., PRIEUR, C. and ZAHM, O. (2022). Nonlinear dimension reduction for surrogate modeling using gradient information. *Information and Inference: A Journal of the IMA*.
- [6] BISHOP, C. M. (1994). Mixture density networks Technical Report No. Neural Computing Research Group report: NCRG/94/004, Aston University.
- [7] BOGACHEV, V. I., KOLESNIKOV, A. V. and MEDVEDEV, K. V. (2005). Triangular transformations of measures. *Sbornik: Mathematics* **196** 309.
- [8] BOYD, J. P. (1984). Asymptotic coefficients of Hermite function series. *Journal of Computational Physics* **54** 382–410.
- [9] BRENNAN, M., BIGONI, D., ZAHM, O., SPANTINI, A. and MARZOUK, Y. (2020). Greedy inference with structure-exploiting lazy maps. *Advances in Neural Information Processing Systems* **33**.
- [10] CHANG, S.-H., COSMAN, P. C. and MILSTEIN, L. B. (2011). Chernoff-type bounds for the Gaussian error function. *IEEE Transactions on Communications* **59** 2939–2944.
- [11] CHKIFA, A., COHEN, A. and SCHWAB, C. (2015). Breaking the curse of dimensionality in sparse polynomial approximation of parametric PDEs. *Journal de Mathématiques Pures et Appliquées* **103** 400–428.
- [12] COHEN, A. (2003). *Numerical analysis of wavelet methods*. Elsevier.

- [13] COHEN, A. and MIGLIORATI, G. (2018). Multivariate approximation in downward closed polynomial spaces. In *Contemporary Computational Mathematics-A celebration of the 80th birthday of Ian Sloan* 233–282. Springer.
- [14] CUI, T. and DOLGOV, S. (2021). Deep composition of Tensor-Trains using squared inverse Rosenblatt transports. *Foundations of Computational Mathematics* 1–60.
- [15] CUI, T., DOLGOV, S. and ZAHM, O. (2021). Conditional Deep Inverse Rosenblatt Transports. *arXiv preprint arXiv:2106.04170*.
- [16] DINH, L., SOHL-DICKSTEIN, J. and BENGIO, S. (2017). Density estimation using Real NVP. In *International Conference on Learning Representations*.
- [17] DINH, L., SOHL-DICKSTEIN, J. and BENGIO, S. (2017). Density estimation using Real NVP. In *International Conference on Learning Representations*.
- [18] DURKAN, C., BEKASOV, A., MURRAY, I. and PAPAMAKARIOS, G. (2019). Neural spline flows. In *Advances in Neural Information Processing Systems* 7509–7520.
- [19] EL MOSELHY, T. A. and MARZOUK, Y. M. (2012). Bayesian inference with optimal maps. *Journal of Computational Physics* **231** 7815–7850.
- [20] HUANG, C.-W., CHEN, R. T., TSIRIGOTIS, C. and COURVILLE, A. (2020). Convex Potential Flows: Universal Probability Distributions with Optimal Transport and Convex Optimization. In *International Conference on Learning Representations*.
- [21] HUANG, C.-W., KRUEGER, D., LACOSTE, A. and COURVILLE, A. (2018). Neural Autoregressive Flows. In *International Conference on Machine Learning* 2083–2092.
- [22] JAINI, P., SELBY, K. A. and YU, Y. (2019). Sum-of-squares polynomial flow. In *International Conference on Machine Learning* 3009–3018.
- [23] KATZFUSS, M. and SCHÄFER, F. (2021). Scalable Bayesian transport maps for high-dimensional non-Gaussian spatial fields. *arXiv preprint arXiv:2108.04211*.
- [24] KINGMA, D. P. and DHARIWAL, P. (2018). Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems* 10215–10224.
- [25] KOBYZEV, I., PRINCE, S. and BRUBAKER, M. (2020). Normalizing flows: An introduction and review of current methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- [26] KOLLER, D. and FRIEDMAN, N. (2009). *Probabilistic graphical models: principles and techniques*. MIT press.
- [27] KOVACHKI, N., BAPTISTA, R., HOSSEINI, B. and MARZOUK, Y. (2020). Conditional Sampling With Monotone GANs. *arXiv preprint arXiv:2006.06755*.
- [28] KUFNER, A. and OPIC, B. (1984). How to define reasonably weighted Sobolev spaces. *Commentationes Mathematicae Universitatis Carolinae* **25** 537–554.
- [29] LEZCANO CASADO, M. (2019). Trivializations for gradient-based optimization on manifolds. *Advances in Neural Information Processing Systems* **32** 9157–9168.
- [30] LICHMAN, M. (2013). UCI Machine Learning Repository. <http://archive.ics.uci.edu/ml>.
- [31] LUECKMANN, J.-M., BOELTS, J., GREENBERG, D., GONCALVES, P. and MACKE, J. (2021). Benchmarking simulation-based inference. In *International Conference on Artificial Intelligence and Statistics* 343–351. PMLR.
- [32] MALLAT, S. (1999). *A wavelet tour of signal processing*. Elsevier.
- [33] MARZOUK, Y., MOSELHY, T., PARNO, M. and SPANTINI, A. (2016). *Sampling via Measure Transport: An Introduction* In *Handbook of Uncertainty Quantification* 1–41. Springer International Publishing.
- [34] MIGLIORATI, G. (2015). Adaptive polynomial approximation by means of random discrete least squares. In *Numerical Mathematics and Advanced Applications-ENUMATH 2013* 547–554. Springer.
- [35] MIGLIORATI, G. (2019). Adaptive approximation by optimal weighted least-squares methods. *SIAM Journal on Numerical Analysis* **57** 2217–2245.
- [36] MORRISON, R., BAPTISTA, R. and MARZOUK, Y. (2017). Beyond normality: Learning sparse probabilistic graphical models in the non-Gaussian setting. In *Advances in Neural Information Processing Systems* 2359–2369.
- [37] MUCKENHOUT, B. (1972). Hardy’s inequality with weights. *Studia Mathematica* **44** 31–38.
- [38] NOCEDAL, J. and WRIGHT, S. (2006). *Numerical optimization*. Springer Science & Business Media.
- [39] NOVAK, E., ULLRICH, M., WOŹNIAKOWSKI, H. and ZHANG, S. (2018). Reproducing kernels of Sobolev spaces on  $\mathbb{R}^d$  and applications to embedding constants and tractability. *Analysis and Applications* **16**

- 693–715.
- [40] OORD, A. V. D., LI, Y., BABUSCHKIN, I., SIMONYAN, K., VINYALS, O., KAVUKCUOGLU, K., DRIESCHKE, G. V. D., LOCKHART, E., COBO, L. C., STIMBERG, F. et al. (2017). Parallel WaveNet: Fast high-fidelity speech synthesis. *arXiv preprint arXiv:1711.10433*.
  - [41] PAPAMAKARIOS, G. and MURRAY, I. (2016). Fast  $\varepsilon$ -free inference of simulation models with Bayesian conditional density estimation. In *Advances in Neural Information Processing Systems* 1028–1036.
  - [42] PAPAMAKARIOS, G., NALISNICK, E., REZENDE, D. J., MOHAMED, S. and LAKSHMINARAYANAN, B. (2021). Normalizing flows for probabilistic modeling and inference. *Journal of Machine Learning Research* **22** 1–64.
  - [43] PAPAMAKARIOS, G., PAVLAKOU, T. and MURRAY, I. (2017). Masked autoregressive flow for density estimation. In *Advances in Neural Information Processing Systems* 2338–2347.
  - [44] PARNO, M. D. and MARZOUK, Y. M. (2018). Transport map accelerated Markov chain Monte Carlo. *SIAM/ASA Journal on Uncertainty Quantification* **6** 645–682.
  - [45] RADEV, S. T., MERTENS, U. K., VOSS, A., ARDIZZONE, L. and KÖTHE, U. (2020). BayesFlow: Learning complex stochastic models with invertible neural networks. *IEEE transactions on neural networks and learning systems*.
  - [46] RAMSAY, J. O. (1998). Estimating smooth monotone functions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **60** 365–375.
  - [47] RASKUTTI, G. and UHLER, C. (2018). Learning directed acyclic graph models based on sparsest permutations. *Stat* **7** e183.
  - [48] REZENDE, D. J. and MOHAMED, S. (2015). Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*.
  - [49] ROSENBLATT, M. (1952). Remarks on a multivariate transformation. *The Annals of Mathematical Statistics* **23** 470–472.
  - [50] ROTHFUSS, J., FERREIRA, F., WALTHER, S. and ULRICH, M. (2019). Conditional Density Estimation with Neural Networks: Best Practices and Benchmarks. *arXiv:1903.00954*.
  - [51] SANTAMBROGIO, F. (2015). *Optimal Transport for Applied Mathematicians*. Springer International Publishing.
  - [52] SCHÄFER, F., KATZFUSS, M. and OWHADI, H. (2021). Sparse Cholesky Factorization by Kullback–Leibler Minimization. *SIAM Journal on Scientific Computing* **43** A2019–A2046.
  - [53] SCHMULAND, B. (1992). Dirichlet forms with polynomial domain. *Math. Japon* **37** 1015–1024.
  - [54] SCHÖLKOPF, B., HERBRICH, R. and SMOLA, A. J. (2001). A generalized representer theorem. In *International conference on computational learning theory* 416–426. Springer.
  - [55] SHIN, Y. E., ZHOU, L. and DING, Y. (2022). Joint estimation of monotone curves via functional principal component analysis. *Computational Statistics & Data Analysis* **166** 107343.
  - [56] SILVERMAN, B. W. (1982). On the estimation of a probability density function by the maximum penalized likelihood method. *The Annals of Statistics* 795–810.
  - [57] SISSON, S. A., FAN, Y. and TANAKA, M. M. (2007). Sequential Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences* **104** 1760–1765.
  - [58] SPANTINI, A., BAPTISTA, R. and MARZOUK, Y. (In press). Coupling techniques for nonlinear ensemble filtering. *SIAM Review*.
  - [59] SPANTINI, A., BIGONI, D. and MARZOUK, Y. (2018). Inference via low-dimensional couplings. *The Journal of Machine Learning Research* **19** 2639–2709.
  - [60] TABAK, E. G. and TURNER, C. V. (2013). A family of nonparametric density estimation algorithms. *Communications on Pure and Applied Mathematics* **66** 145–164.
  - [61] TESHIMA, T., ISHIKAWA, I., TOJO, K., OONO, K., IKEDA, M. and SUGIYAMA, M. (2020). Coupling-based invertible neural networks are universal diffeomorphism approximators. In *Advances in Neural Information Processing Systems* **33** 3362–3373.
  - [62] TRIPPE, B. L. and TURNER, R. E. (2018). Conditional density estimation with Bayesian normalising flows. In *Bayesian Deep Learning: NIPS 2017 Workshop*.
  - [63] TRUONG, T. T. and NGUYEN, H.-T. (2021). Backtracking Gradient Descent Method and Some Applications in Large Scale Optimisation. Part 2: Algorithms and Experiments. *Applied Mathematics & Optimization* **84** 2557–2586.
  - [64] URIA, B., MURRAY, I. and LAROCHELLE, H. (2013). RNADE: The real-valued neural autoregressive

- density-estimator. *arXiv preprint arXiv:1306.0186*.
- [65] VERSHYNIN, R. (2018). *High-dimensional probability: An introduction with applications in data science* **47**. Cambridge university press.
  - [66] VIDAKOVIC, B. (2009). *Statistical modeling by wavelets* **503**. John Wiley & Sons.
  - [67] VILLANI, C. (2008). *Optimal transport: old and new* **338**. Springer Science & Business Media.
  - [68] WASSERMAN, L. (2013). *All of statistics: a concise course in statistical inference*. Springer Science & Business Media.
  - [69] WEHENKEL, A. and LOUPPE, G. (2019). Unconstrained monotonic neural networks. In *Advances in Neural Information Processing Systems* 1543–1553.
  - [70] WENLIANG, L., SUTHERLAND, D., STRATHMANN, H. and GRETTON, A. (2019). Learning deep kernels for exponential family densities. In *International Conference on Machine Learning* 6737–6746.
  - [71] ZECH, J. and MARZOUK, Y. (2022). Sparse approximation of triangular transports. Part II: the infinite dimensional case. *Constructive Approximation* **55** 987–1036.
  - [72] ZECH, J. and MARZOUK, Y. (2022). Sparse Approximation of Triangular Transports. Part I: The Finite-Dimensional Case. *Constructive Approximation* **55** 919–986.