# LANDING WITH THE SCORE: RIEMANNIAN OPTIMIZATION THROUGH DENOISING

## **Anonymous authors**

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

034

035

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

## **ABSTRACT**

Under the data manifold hypothesis, high-dimensional data are concentrated near a low-dimensional manifold. We study the problem of Riemannian optimization over such manifolds when they are given only implicitly through the data distribution, and the standard manifold operations required by classical algorithms are unavailable. This formulation captures a broad class of data-driven design problems that are central to modern generative AI. Our key idea is to introduce a link function that connects the data distribution to the geometric operations needed for optimization. We show that this function enables the recovery of essential manifold operations, such as retraction and Riemannian gradient computation. Moreover, we establish a direct connection between our construction and the score function in diffusion models of the data distribution. This connection allows us to leverage well-studied parameterizations, efficient training procedures, and even pretrained score networks from the diffusion model literature to perform optimization. Building on this foundation, we propose two efficient inference-time algorithms—Denoising Landing Flow (DLF) and Denoising Riemannian Gradient Descent (DRGD)—and provide theoretical guarantees for both feasibility (approximate manifold adherence) and optimality (small Riemannian gradient norm). Finally, we demonstrate the effectiveness of our approach on finite-horizon reference tracking tasks in data-driven control, highlighting its potential for practical generative and design applications.

#### 1 Introduction

Riemannian optimization Boumal (2023); Absil et al. (2008); Hu et al. (2020) considers minimizing an objective function  $f: \mathbb{R}^d \to \mathbb{R}$  over an *explicitly known* embedded submanifold  $\mathcal{M} \subseteq \mathbb{R}^d$ ,

$$\min_{x \in \mathcal{M}} f(x). \tag{1}$$

Problem (1) is ubiquitous in fields of machine learning and control and encompasses problems such as independent component analysis Nishimori (1999), low-rank matrix completion Vandereycken (2013), training of orthogonally normalized neural networks Bansal et al. (2018), the control of rigid bodies Duong et al. (2024), as well as sensor network localization Patwari & Hero (2004) and many others. Compared to general constrained optimization, Riemannian optimization promises the advantage of exploiting the natural geometry of the problem, producing feasible iterates and increased numerical robustness Boumal (2023).

In contrast to the above classical setup, in this work we focus on the setting where the manifold  $\mathcal{M}$  is given *implicitly* through a measure  $\mu = \mu_{\text{data}}$  supported on  $\mathcal{M}$ . This perspective is especially relevant in view of the *manifold hypothesis* (Loaiza-Ganem et al., 2024), which posits that many real-world data sets lie (approximately) on a manifold with dimension much smaller than that of the ambient space Fefferman et al. (2016). Importantly, such data manifolds are not just low-dimensional geometric structures—they also capture rich *semantic meaning*. For instance, the image manifold corresponds to photo-realistic images Pope et al., the system behavior manifold to dynamically feasible input-output trajectories Willems & Polderman (1997), while the manifold of airfoils represents aerodynamically viable shapes Zheng et al. (2025). The optimization problem (1) in this setting thus encompasses a broad class of modern tasks, including airfoil design (Chen et al., 2025), reinforcement learning (Lee & Choi, 2025), and Bayesian inverse problems (Chung et al., 2022b).

Driven by the growing demand in generative AI and data-driven design, this calls for a paradigm shift: moving from optimization over explicitly known manifolds, where classical Riemannian optimization applies, to optimization over data manifolds that are accessible only implicitly through samples. In this data-driven regime, methods from classical Riemannian optimization *cannot* be applied directly, since they rely on explicit manifold operations such as tangent-space projection, retraction, or exponential maps (Boumal et al., 2019; Boumal, 2023).

While there has been extensive research on manifold learning Meilă & Zhang (2024); Lin & Zha (2008); Cayton et al. (2005); Belkin & Niyogi (2005) in the past, none of these works addressed (1) from an optimization point of view and focused instead on learning the manifold geometry. Moreover, the non-parametric nature of these methods limits their applicability to very low-dimensional manifolds and cannot exploit the successful inductive bias and the exceptional representation power of modern neural networks.

In this work, we propose a data-driven approach to recover the fundamental operations needed for optimization on manifolds. Starting from the data distribution  $\mu_{\rm data}$ , we smooth it with a Gaussian kernel to obtain

$$p_{\sigma} = \mathcal{N}(0, \sigma^2 I) * \mu_{\text{data}}, \tag{2}$$

and define the associated link function

$$\ell_{\sigma}(x) = \frac{1}{2} ||x||^2 - \sigma^2 \log p_{\sigma}(x).$$
 (3)

We show that, as the smoothing parameter  $\sigma$  decreases, the gradient  $\nabla \ell_{\sigma}$  recovers the projection back to the manifold, while the Hessian  $\nabla^2 \ell_{\sigma}$  recovers the projection onto its tangent space. These results reveal that core ingredients of Riemannian optimization—such as retraction and gradient computation—can be implemented directly from the derivative information of  $\ell_{\sigma}$ , even when the manifold itself is only given implicitly through data.

Given the above novel theoretical findings, a key practical challenge is how to access (even approximately) the gradient  $\nabla \ell_{\sigma}$  and Hessian  $\nabla^2 \ell_{\sigma}$  when only samples from  $\mu_{\rm data}$  are available. A crucial observation is that  $p_{\sigma}$  coincides with the marginal distribution of the Variance-Exploding SDE (VE-SDE), with analogous arguments for VP-SDE and DDPM formulations (Ho et al., 2020; Song et al., 2020). In this setting, the gradient of  $\log p_{\sigma}$  – commonly referred to as the score function - plays a central role. In diffusion models, this score is parameterized by a neural network and learned directly from samples of  $\mu_{\rm data}$ , typically via denoising score matching (Vincent, 2011; Song et al., 2020). It is therefore natural to adopt the same methodology here: a neural network can be trained to represent  $\nabla \ell_{\sigma}$ , while the Hessian  $\nabla^2 \ell_{\sigma}$  can be recovered by computing its Jacobian. This approach offers two key advantages: (i) strong inductive biases can be incorporated through the neural network parameterization, and (ii) efficient, well-established training techniques from the diffusion model literature can be directly leveraged. Taken together, the theoretical link between  $\ell_{\sigma}$  and manifold geometry, and the practical machinery for learning the score, form the foundation of the paradigm shift: from classical Riemannian optimization with explicit manifold knowledge to a data-driven framework where geometry is recovered from samples—thus enabling principled manifold optimization in generative and design-driven applications.

Building on these insights, we propose two algorithms for optimization over data manifolds – denoising landing flow (DLF) and denoising Riemannian gradient descent (DRGD) – to the best of our knowledge, the first such framework in the literature. The prerequisite is access to an approximate function  $s^{\sigma}_{\theta} \simeq \nabla \ell_{\sigma}$ , which can be learned from data using standard diffusion model parameterization and training given samples from  $\mu_{\rm data}$ .

- Our first method directly employs  $s_{\theta}^{\sigma}$  and its Jacobian (with  $\sigma$  fixed as a small constant) to substitute the corresponding manifold operations in Riemannian gradient descent: each iteration computes a Riemannian gradient step and then applies retraction to (approximately) return to the manifold.
- Since learning the score is inevitably imperfect when  $\sigma$  is small, we also design a complementary landing-type algorithm that avoids staying close to the feasible set at every step. Instead, it augments the target function f with the regularization term  $\ell_{\sigma}$  and performs gradient descent on this regularized objective.

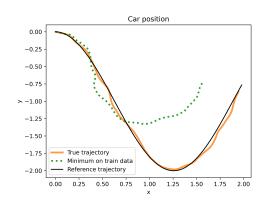


Figure 1: Optimized trajectory (orange) on the system trajectory manifold for the unicycle car model that is desired to track the reference trajectory (black, thin). The closest tracking trajectory in the set of available manifold samples is given in (green, dotted). See Section 6 for more details.

For both algorithms, we establish non-asymptotic convergence guarantees to approximate stationary points, measured by small *Riemannian gradient* norm, as  $\sigma \to 0$ . Importantly, our methods require only inference of the neural network and gradients with respect to its inputs—not with respect to the network parameters. Thus, if a pretrained score network is already available for a given task, no additional training is required to enable Riemannian optimization on the corresponding data manifold. Viewed from this perspective, our approach can also be interpreted as an *inference-time algorithm*, aligning with a growing trend in modern machine learning research.

## 1.1 OUR CONTRIBUTIONS

We summarize our contributions and give the outline of the paper as follows.

- Data-driven recovery of manifold operations. By introducing the link function  $\ell_{\sigma}$ , we develop strategies to recover core manifold operations directly from data. In Section 3, by leveraging the score function from diffusion models to train and parameterize  $\nabla \ell_{\sigma}$  and  $\nabla^2 \ell_{\sigma}$ , we bridge the gap between the requirements of classical Riemannian optimization (which assumes explicit manifold knowledge) and the emerging need to optimize over data manifolds that are only implicitly available.
- First algorithms for optimization over data manifolds. We propose two algorithms denoising landing flow (DLF) via (8) in Section 4 and denoising Riemannian gradient descent (DRGD) via (12) in Section 5 to the best of our knowledge, the first in the literature that exploit operations enabled by a pretrained score function. Our methods require only inexpensive inference-time queries and back-propagation with respect to the input of the neural network, making them computationally efficient and readily applicable when pretrained scores are available.
- **Rigorous non-asymptotic guarantees.** We provide the first non-asymptotic convergence analysis for optimization over data manifolds. Our main results (Theorem 3 and Theorem 5) ensure both approximate feasibility (output is close to the manifold) and approximate optimality (small *Riemannian gradient* norm). A key technical contribution is a novel non-asymptotic analysis of how  $\nabla \ell_{\sigma}$  and  $\nabla^2 \ell_{\sigma}$  approximate the true manifold operations when  $\sigma$  is small but non-zero (Theorem 1).

Finally in Section 6 we validate our findings by numerically simulating the proposed flow on well-known manifolds from Riemannian optimization literature and data-driven optimal control for reference trajectory tracking. Our finding are that we can generate feasible points on the manifold with objective values far lower than the ones available in the training set – see Figure 1 – and shows the effectiveness of exploiting the strong inductive biases of modern deep learning for the classical field of constrained optimization.

## 1.2 RELATED WORK

**Riemannian Optimization.** Riemannian optimization (RO) was originally developed under the assumption that the constraint manifold is explicitly known, either through closed-form descriptions such as matrix manifolds or via nonlinear equality constraints (Sato, 2021; Boumal, 2023). A prototypical algorithm in this setting is *Riemannian gradient descent*, which updates by taking a gradient step along the tangent space and then retracting back to the manifold (Boumal et al., 2019). Such algorithms guarantee that the iterates remain on the manifold at every iteration and are thus referred to as *feasible methods*. In contrast, *infeasible methods*, where the iterates are not constrained to stay on  $\mathcal{M}$ , have also been studied—for example, augmented Lagrangian approaches (Xie & Wright, 2021). However, these typically involve solving complicated optimization subproblems in each inner loop, making them computationally expensive in practice. More recently, a new class of *landing-type algorithms* has emerged: instead of enforcing feasibility at every step, they regularize the objective function with the distance to the manifold and then perform gradient flow or descent on this regularized objective. Such methods have demonstrated strong empirical performance, offering a promising alternative to classical approaches (Ablin & Peyré, 2022; Schechtman et al., 2023).

**Manifold Learning.** High-dimensional data in modern machine learning often exhibit an intrinsic lower-dimensional structure. Such structure is of central importance: it enables more efficient representation and compression of data, facilitates interpretability by revealing meaningful semantic organization, and provides a foundation for designing algorithms that exploit geometry rather than ambient dimensionality. The task of uncovering this structure is commonly referred to as *manifold estimation* or *manifold learning*, with a large body of work devoted to this goal. Classical approaches include Isomap (Tenenbaum et al., 2000), Laplacian Eigenmaps (Belkin & Niyogi, 2003), and Locally Linear Embedding (LLE) (Roweis & Saul, 2000), Diffusion Map (Coifman & Lafon, 2006), among many subsequent developments, e.g. (Zhou et al., 2020). We defer the discussion on diffusion-model-related manifold learning literature to the next paragraph.

**Diffusion models.** Diffusion models have achieved remarkable success in generative modeling, where the goal is to generate new samples consistent with an underlying data distribution  $\mu_{\rm data}$ . A central ingredient of these methods is the learning of the *score function*—the gradient of the log-density of the diffused distribution  $p_{\sigma}$  Tang & Zhao (2025); Song et al. (2020). In practice, the score is parameterized by a neural network whose architectural design has been extensively studied, and its training is carried out using well-established techniques such as denoising score matching (Song et al., 2020). This combination of principled theory and mature practice has made diffusion models one of the most effective tools for data-driven generative modeling.

When the data distribution resides on a manifold, a recent observation in the diffusion model community is that the score is asymptotically orthogonal to the manifold surface Stanczuk et al. (2024). This observation has been exploited for the estimation of the manifold dimension. See (Kamkari et al., 2024) for the same task. Furthermore, Ventura et al. (2024) has shown that in the case of linear manifolds (i.e. affine subspaces), the Jacobian of the score scaled by the diffusion temperature asymptotically approximates the projection of the manifold onto the normal space and used it to study the geometric phases of diffusion models.

Recent work on the statistical complexity of diffusion models under the manifold hypothesis provides indirect evidence that these models capture geometric information about the underlying data manifold (Oko et al., 2023; Tang & Yang, 2024). In particular, the sample complexity required to learn the data distribution  $\mu_{\rm data}$  depends only on the intrinsic dimension of the manifold, rather than on the ambient dimension.

**Two Formulations: Optimization and Posterior Sampling.** To avoid possible confusion, we stress the distinction between our optimization formulation and the posterior sampling literature. In short, our optimization formulation in eq. (1) enforces the manifold constraint directly. This ensures (approximate) feasibility at the final step and guarantees that the optimization process remains semantically meaningful, which is not guaranteed by the sampling formulation, as discussed below.

Across the literatures Classifier(-Free) Guidance in Diffusion Models (Dhariwal & Nichol, 2021; Ho & Salimans, 2022) and the Plug-and-Play Framework in Bayesian Inverse Problems (Venkatakrishnan et al., 2013; Laumont et al., 2022; Pesme et al., 2025; Graikos et al., 2022; Chung et al., 2022a),

a unifying perspective is the task of sampling from a posterior distribution of the form

$$p_{\rm post} \propto p_{\rm pre} \exp\left(-\frac{r}{\alpha}\right)$$

where r denotes a cost function to be minimized—corresponding to our objective f—and  $\alpha>0$  is a temperature-like parameter (Domingo-Enrich et al., 2024). In Classifier(-Free) Guidance, r encodes a classifier signal (e.g., specified by a prompt) and in Bayesian inverse problems, r is the negative log-likelihood of observations. Meanwhile,  $p_{\rm pre}$  serves as a prior distribution, typically derived from a large-scale pretrained generative model, which is expected to capture the semantic structure of the data manifold. Sampling from  $p_{\rm post}$  thus aims to balance semantic plausibility with low cost.

If  $p_{\text{pre}}$  were *exactly* supported on the data manifold  $\mathcal{M}$  and  $\alpha \to 0$ , this formulation would reduce to the constrained optimization problem eq. (1). In practice, however, the situation is very different: the distribution induced by a pretrained diffusion model is not concentrated on a low-dimensional manifold but rather has support of non-zero Lebesgue measure in the ambient space—indeed, in many cases, essentially the full space (due to the noisy generation process of  $p_{\text{pre}}$ ). As a result, when  $\alpha$  is set too small, the posterior  $p_{\text{post}}$  becomes dominated by the exponential tilt  $\exp(-r/\alpha)$ , pushing samples into regions far from the true data manifold  $\mathcal M$  and thereby *losing* semantic meaning. Consequently, these sampling-based frameworks must carefully tune  $\alpha$  to trade off semantic fidelity (staying close to  $\mathcal M$ ) against optimization quality (achieving low r).

# 2 PRELIMINARIES

Here we introduce some preliminary notation and concepts to state our results. We refer the reader to the Appendices A and C for more details.

#### 2.1 Manifolds and distance functions

Let  $\mathcal{M}\subseteq\mathbb{R}^d$  be a k-dimensional embedded compact  $C^2$ -submanifold without boundary. For any point  $p\in\mathcal{M}$  we denote by  $\mathrm{T}_p\,\mathcal{M}\subseteq\mathbb{R}^d$  and  $\mathrm{N}_p\,\mathcal{M}\subseteq\mathbb{R}^d$  the tangent and normal spaces of  $\mathcal{M}$  at p, respectively, and their orthogonal projections by  $\mathrm{P}_{\mathrm{T}_p\,\mathcal{M}}$  and  $\mathrm{P}_{\mathrm{N}_p\,\mathcal{M}}$ . The squared distance function is defined by  $\mathrm{d}(x)=\inf_{p\in\mathcal{M}}\frac{1}{2}\|x-p\|^2$ . For a  $C^2$ -submanifold, then there exists a radius  $\tau_{\mathcal{M}}>0$  such that every point in  $x\in\mathcal{T}(\tau)=\bigcup_{p\in\mathcal{M}}B_{\tau}(p)$  has a unique projection  $\pi(x)\in\mathcal{M}$  such that  $\frac{1}{2}\|x-\pi(x)\|^2=\mathrm{d}(x)$  and  $x-\pi(x)\in\mathrm{N}_{\pi(x)}\,\mathcal{M}$ . Sets of the form  $\mathcal{T}(\tau)$  will be called tubular neighborhoods of  $\mathcal{M}$  with closure denoted by  $\overline{\mathcal{T}}(\tau)$ . One can prove that  $\mathrm{d}$  and  $\pi$  are differentiable on  $\mathcal{T}(\tau_{\mathcal{M}})$  with  $\mathrm{d}'(x)=x-\pi(x)$  and  $\pi'(x)$  given explicitly in Appendix A. The maximal principal curvature of a manifold will be denoted by  $\kappa_{\mathcal{M}}$ . It always holds that  $\tau_{\mathcal{M}}\leq 1/\kappa_{\mathcal{M}}$ . Finally, for a function  $f:\mathcal{M}\to\mathbb{R}$  we denote the Riemannian gradient by  $\mathrm{grad}_{\mathcal{M}}\,f(p)$ , which in the case of  $f:\mathbb{R}^d\to\mathbb{R}$  is given by  $\mathrm{grad}_{\mathcal{M}}\,f(p)=\mathrm{P}_{\mathrm{T}_p\,\mathcal{M}}\,\nabla f(p)$ . We write

# 2.2 THE STEIN SCORE FUNCTION AND SCORE-BASED DIFFUSION MODELS

For a Borel probability measure we denote its Gaussian blurring by  $p_{\sigma} = \mathcal{N}(0, \sigma^2 I) * \mu$  for  $\sigma > 0$  with  $p_0 = \mu$ . The score function  $\nabla \log p_{\sigma}$  of  $p_{\sigma}$  and its Jacobian allow for the following interpretation (see Jaffer & Gupta (1972), also Appendix B):

$$x + \sigma^2 \nabla \log p_{\sigma}(x) = \nabla \ell_{\sigma}(x) = \mathbb{E}\nu_{x,\sigma}, \quad I + \sigma^2 \nabla^2 \log p_{\sigma}(x) = \nabla^2 \ell_{\sigma}(x) = \frac{1}{\sigma^2} \operatorname{Cov}(\nu_{x,\sigma}), \quad (4)$$

where  $\ell_{\sigma}$  is the link function (3) and  $\nu_{x,\sigma}$  is the posterior distribution observing x under the noise model  $p_{\sigma}$  and prior  $\mu$ . The representation for  $\mathbb{E}\nu_{x,\sigma}$  has also been known under the name of Tweedie's formula Robbins (1992); Efron (2011). The score function has gained recent attention due to its use in score-based diffusion models in the field of generative modelling. Specifically in the so-called *variance exploding* (VE) diffusion scheme one seeks to learn  $\nabla \log p_{\sigma}$  for different noise scales  $\sigma$  via a neural network  $s(\cdot,\sigma)$  by minimizing the conditional score matching loss  $L_{\text{CSM}}(s(\cdot,\sigma))$ , which attains its unique minimum in  $s(\cdot,\sigma) = \nabla \log p_{\sigma}$ . Then  $s(\cdot,\sigma)$  is used for sampling from  $\mu$  by following a particular reverse-time SDE or ODE flow in the noise scale  $\sigma$  (see Appendix C for more details).

## 3 SCORE AS RETRACTION AND TANGENT SPACE PROJECTION

In this section we give another interpretation to the score function  $\nabla \log p_{\sigma}(x)$  and its Jacobian  $\nabla^2 \log p_{\sigma}(x)$ . Namely, it has been already observed in Stanczuk et al. (2024) that the score function is for  $\sigma \to 0$  asymptotically orthogonal to the tangent space of the data manifold. Our first contribution is showing that when the support of the distribution  $\mu$  is a manifold  $\mathcal M$  and  $\mu$  is absolutely continuous w.r.t. its volume measure, then both quantities (4) approximate the projection operator  $\pi(x)$  and its Jacobian  $\pi'(x)$  uniformly on tubular neighborhoods of  $\mathcal M$ . Formally we establish the following

**Theorem 1** (Main). Let  $\mathcal{M} \subseteq \mathbb{R}^d$  be a compact, embedded  $C^3$ -submanifold and  $\mu \in \mathcal{P}(\mathbb{R}^d)$  a Borel probability measure with supp  $\mu = \mathcal{M}$  and  $\mu \ll \operatorname{Vol}_{\mathcal{M}}$  such that  $\frac{\mathrm{d}\,\mu}{\mathrm{d}\,\operatorname{Vol}_{\mathcal{M}}} \in C^3(\mathcal{M})$ . Then for any  $\tau \in (0, \tau_{\mathcal{M}})$  there exist some constants  $K = K(\tau, \mathcal{M}, \mu) > 0$  and  $\overline{\sigma} = \overline{\sigma}(\tau, \mathcal{M}, \mu) > 0$  depending on  $\tau$ ,  $\mathcal{M}$  and  $\mu$  such that

$$\|\mathbb{E}\nu_{x,\sigma} - \pi(x)\| \le K\sigma |\log(\sigma)|^3 \text{ and } \left\| \frac{1}{\sigma^2} \operatorname{Cov}(\nu_{x,\sigma}) - \pi'(x) \right\| \le K\sigma |\log(\sigma)|^3$$
 (5)

for all  $\sigma \in (0, \overline{\sigma})$  and  $x \in \mathcal{T}(\tau)$ .

 The proof is deferred to Appendix D and is based on a careful non-asymptotic estimate of the Laplace integral method. As a consequence and together with fact that  $\pi'(x)$  coincides with  $P_{T_x \mathcal{M}}$  for  $x \in \mathcal{M}$  (see Appendix A) we obtain the following result.

**Corollary 2.** Let  $\mathcal{M}$  and  $\mu$  be as in Theorem 1 and suppose that  $x \in \mathcal{M}$ . Then

$$\lim_{\sigma \to 0} I + \sigma^2 \nabla^2 \log p_{\sigma}(x) = P_{T_x \mathcal{M}}.$$

In view of Theorem 1 let us abbreviate (4) into

$$d_{\sigma}(x) = -\sigma^2 \log p_{\sigma}(x), \quad \pi_{\sigma}(x) = x + \sigma^2 \nabla \log p_{\sigma}(x), \quad P_{\sigma}(x) = I + \sigma^2 \nabla^2 \log p_{\sigma}(x), \quad (6)$$

with the limiting cases  $d_0 = d$ ,  $\pi_0 = \pi$  and  $P_0(x) = \pi'(x)$ . We stress that the expressions in (6) are defined for all  $x \in \mathbb{R}^d$ , whereas d,  $\pi$  and  $\pi'$  in are only sensible in a tubular neighborhood  $\mathcal{T}$  of  $\mathcal{M}$ . Thus, in theory, a well-trained diffusion model score  $s(\cdot,\sigma)$  and its Jacobian  $s'(\cdot,\sigma)$  allow us to approximate the closest-point projection  $\pi(x)$  and the tangent space projection  $P_{T_x \mathcal{M}}$  as  $\sigma \to 0$  arbitrarily well via its Tweedie score  $x + \sigma^2 s(x,\sigma)$  and its Jacobian, respectively.

## 4 DENOISING RIEMANNIAN GRADIENT FLOW WITH LANDING

In this section we show how to use the score function from Section 3 for Riemannian optimization of (1) for some smooth  $f \in C^1(\mathbb{R}^d)$ . Assuming that we have access to a sufficiently accurate estimate of the Tweedie score in form of a vector function  $s \in C^1(\mathbb{R}^d; \mathbb{R}^d)$  such that  $s(x) = x + \sigma^2 s(x, \sigma)$  and

$$||s(x) - \pi_{\sigma}(x)|| \le \epsilon \text{ and } ||s'(x) - P_{\sigma}(x)|| \le \epsilon \text{ for } x \in \mathcal{T}(\tau)$$
 (7)

for some  $\tau \in (0, \tau_{\mathcal{M}})$ , we propose for  $\eta \geq 0$  the denoising landing flow (DLF)

$$\dot{x} = -s'(x)\nabla f(s(x)) + \eta(s(x) - x). \tag{8}$$

In the exact case  $s(x) = \pi_{\sigma}(x)$  and  $s'(x) = P_{\sigma}(x)$  (i.e.  $\epsilon = 0$  in (7)), flow (8) is the gradient flow

$$\dot{x} = -\nabla F_{\sigma}^{\eta}(x) = -P_{\sigma}(x)\nabla f(\pi_{\sigma}(x)) + \eta(\pi_{\sigma}(x) - x) \text{ with } F_{\sigma}^{\eta}(x) = f(\pi_{\sigma}(x)) + \eta \,\mathrm{d}_{\sigma}(x) \quad (9)$$

and the dynamics in (9) consists of two parts: An approximate projection  $P_{\sigma}(x)\nabla f(\pi_{\sigma}(x))$  of the gradient  $\nabla f(\pi_{\sigma}(x))$  and an approximate landing term  $\eta(\pi_{\sigma}(x)-x)$  corresponding to the penalty function  $\eta \, \mathrm{d}_{\sigma}(x)$ . In the further case of  $\sigma=0$  and  $x(0)\in\mathcal{M}$  the flow (9) reduces to the ordinary Riemannian gradient flow, which has been extensively studied Helmke & Moore (2012); Ambrosio et al. (2005). Interestingly, when  $\sigma=0$ , but only  $x(0)\in\mathcal{T}(\tau)$ , then (9) reduces to

$$\dot{x} = -H_x^{-1} \operatorname{grad}_{\mathcal{M}} f(\pi(x)) + \eta(\pi(x) - x),$$
 (10)

with  $H_x^{-1}$  a linear operator on  $T_{\pi(x)} \mathcal{M}$  given in Appendix A. In particular the two terms in (10) belong to  $T_{\pi(x)} \mathcal{M}$  and  $N_{\pi(x)} \mathcal{M}$ , respectively, and are orthogonal to each other, implying that the distance between x and  $\mathcal{M}$  is non-increasing, which allows for perfect landing on the manifold via similar arguments as in Ablin & Peyré (2022); Schechtman et al. (2023), see Theorem 20 in Appendix E.3. For  $\sigma > 0$  or  $\epsilon > 0$ , the two summands in (9) and (8) in general not perpendicular to each other and the landing is not exact. However, using Theorem 1 we show the following.

**Theorem 3.** Consider the flow (8) and  $\tau \in (0, \tau_{\mathcal{M}})$ . Set  $C = \|\nabla f|_{\mathcal{T}(\tau)}\|_{\infty}$  and  $L = \operatorname{Lip}(\nabla f)$ . Suppose that for some  $\epsilon > 0$  and  $\sigma \in (0, \overline{\sigma}(\tau, \mathcal{M}, \mu))$  with  $\epsilon + K(\tau, \mathcal{M}, \mu)\sigma|\log(\sigma)|^3 \leq \min\{\tau, \frac{2\tau}{1+C/\eta}\}$  the function s satisfies (7). Then for any  $x(0) \in \mathcal{T}(\tau)$  the solution x(t) to (8) exists for all  $t \geq 0$  and is contained in  $\mathcal{T}(\tau)$ . Moreover, every accumulation point  $x_*$  of this flow satisfies

$$\operatorname{dist}_{\mathcal{M}}(x_*) \leq \tau_0 := \frac{1}{2} \left( \frac{C}{\eta} + 1 \right) \left( \epsilon + K(\tau, \mathcal{M}, \mu) \sigma |\log(\sigma)|^3 \right),$$

and for the projection  $p_* = \pi(x_*)$  it holds that

$$\|\operatorname{grad}_{\mathcal{M}} f(p_*)\| \le \left(2(L+C+2\eta) + \frac{(1+C/\eta)/\tau_{\mathcal{M}}}{1-\tau/\tau_{\mathcal{M}}}C\right) \left(\epsilon + K(\tau, \mathcal{M}, \mu)\sigma|\log(\sigma)|^3\right). \tag{11}$$

Thus, Theorem 3 shows that one can still use the flow (8) for a fixed  $\sigma > 0$  to converge to approximate critical points of the objective f, at which the approximation error and norm of the Riemannian gradient are both  $\widetilde{O}(\sigma)$  plus the score error  $\epsilon$ .

**Remark 4.** We can evaluate the right hand side of the flow (8) in a single forward-backward pass of the network s. Namely, given an input x, we compute and store p = s(x) by a forward pass of s, while keeping the computational graph of s(x). Then we evaluate  $v = \nabla f(p)$  and build the computational graph of  $y = \langle s(x), v \rangle$ , while detaching v. Finally we backpropagate on x in y to obtain  $s'(x)v = s'(x)\nabla f(s(x))$ .

# 5 DENOISING RIEMANNIAN GRADIENT DESCENT

In a practical implementation one has to consider a discretized version of the flow (8). A natural alternative is to study the following approximate version of the Riemannian gradient descent Absil et al. (2008); Boumal (2023)

$$x_{k+1} = s(x_k - \gamma_k s'(x_k) \nabla f(x_k)), \qquad (12)$$

which we term the denoising Riemannian gradient descent (DRGD). Here s acts as an approximate retraction and s as an approximate projection onto the tangent space. We obtain the following convergence result for this algorithm.

**Theorem 5.** Let  $\tau \in (0, \tau_{\mathcal{M}}/2)$  and set  $C = \|\nabla f|_{\mathcal{T}(\tau)}\|_{\infty}$  and  $L = \operatorname{Lip}(\nabla f)$ ,  $D = \|f|_{\mathcal{T}(\tau)}\|_{\infty}$  and

$$L_0 = 8C \left( 2\left(\frac{3}{\tau_{\mathcal{M}}} + \tau M\right) + \frac{1}{\tau_{\mathcal{M}}} \right) + 2L.$$

Suppose that for some  $\sigma \in (0, \overline{\sigma}(\tau, \mathcal{M}, \mu))$  and  $\epsilon > 0$  with  $\epsilon' := \epsilon + K(\tau, \mathcal{M}, \mu)\sigma |\log(\sigma)|^3 \le \tau/2$  the function s satisfies (7) and that the step-size  $\gamma_k$  is constrained by  $\gamma_k \in [\gamma_{\min}, \gamma_{\max}]$  with

$$0 < \gamma_{\min} < \gamma_{\max} < \min \left\{ \frac{2}{L_0}, \frac{\tau}{C(4+\tau)} \right\}$$
.

Then for any  $x_0 \in \mathcal{T}(\tau/2)$  the iterates  $x_k$  of (12) satisfy  $\{x_k\}_{k=1}^{\infty} \subseteq \mathcal{T}(\epsilon') \subseteq \mathcal{T}(\tau/2)$  and for the projection  $p_k = \pi(x_k)$  the following average-of-gradient-norm condition holds:

$$\frac{1}{N} \sum_{k=0}^{N} \|\operatorname{grad}_{\mathcal{M}} f(p_k)\|^2 \le \frac{4D/N + (8C^2 \epsilon' / \tau_{\mathcal{M}}^2 + 2(2C + L_0 \gamma_{\max}(C + L)))\epsilon'}{\gamma_{\min}(1 - \frac{L_0}{2} \gamma_{\max})}.$$

In particular there exists at least one accumulation point  $x_* \in \overline{\mathcal{T}}(\epsilon)$  of  $\{x_k\}_{k=0}^{\infty}$  such that its projection  $p_* = \pi(x_*)$  satisfies

$$\|\operatorname{grad}_{\mathcal{M}} f(p_*)\|^2 \le \frac{(8C^2\epsilon'/\tau_{\mathcal{M}}^2 + 2(2C + L_0\gamma_k(C+L)))\epsilon'}{\gamma_{\min}(1 - \frac{L_0}{2}\gamma_{\max})}$$

**Remark 6.** Note that both Theorem 3 and Theorem 5 require the rather strong  $L^{\infty}$ -approximation assumption (7) on the Tweedie score s and its Jacobian. The analysis under a weaker  $L^2$ -bound is out of score for this paper and left for future work.

# 6 NUMERICAL EXPERIMENTS AND APPLICATIONS

In this section we provide some numerical results for our proposed algorithms, namely the denoising landing flow (8) (more precisely, the discretized version (48)) and the denoising Riemannian gradient descent (12).

## 6.1 Optimization on orthogonal group O(n)

In this section we evaluate flow (8) on a synthetic example. In order to compare the error of our method to classical Riemannian optimization techniques, we consider distributions supported on manifolds that are the focus of study in Absil et al. (2008); Boumal (2023). Specifically we consider the orthogonal group manifold  $\mathcal{M} = \mathrm{O}(n) \subseteq \mathbb{R}^{n \times n}$  with  $\mu = \mathrm{Vol}_{\mathcal{M}}$  being the uniform volume measure and Brockett's cost function Helmke & Moore (2012) defined by

$$\min_{X \in \mathcal{O}(n)} f(X) := \operatorname{tr}(AXQX^{\top}),$$

where  $A,Q\in\mathbb{S}^{n\times n}$  are given. We consider the cases  $n\in\{10,20\}$  and assume that we are given a set  $\mathcal{D}_{\text{train}}\subseteq\mathcal{M}$  of  $N_{\text{data}}=20000$  data points from  $\mu$  and train the score function s with denoising score matching (see Appendix C for diffusion models and Appendix G.1.1 for implementation details). In Figure 3 (left) (in Appendix G.1.2 for space reasons) we compare the evolution of the objective value of our approximation of (8) to the exact landing flow (10) for different noise levels  $\sigma>0$  and dimension n. We observe that we can obtain objective values with cost lower than the best possible point in the training set and that the accuracy improves as  $\sigma\to0$ .

## 6.2 Reference tracking via data-driven control

**Problem definition:** In this example we consider applying our method to the control of discrete-time dynamical systems on a finite horizon. Specifically we assume that we are given a discrete-time state-space system

$$x_{k+1} = f(x_k, u_k), \quad y_{k+1} = g(x_k, u_k), \quad k = 0, \dots, N_h - 1,$$
 (13)

on a finite time horizon  $N_h$  with state  $x_k \in \mathbb{R}^{n_x}$ , input  $u_k \in \mathbb{R}^{n_u}$ , output  $y_k \in \mathbb{R}^{n_y}$  and a fixed initial state  $x_0 = 0 \in \mathbb{R}^{n_x}$ . The task is to find inputs  $\boldsymbol{u} = (u_0, \dots, u_{N_h-1})$  such that the corresponding outputs  $\boldsymbol{y} = (y_0, \dots, y_{N_h})$  closely track a prespecified reference trajectory  $\boldsymbol{r} = (r_0, \dots, r_{N_h})$  by solving the optimal control problem

$$\min_{(\boldsymbol{u},\boldsymbol{y})\in\mathcal{M}_{1O}} f(\boldsymbol{u},\boldsymbol{y}) \tag{14}$$

with tracking objective

$$f(\boldsymbol{u}, \boldsymbol{y}) = \sum_{k=0}^{N_{\rm h}-1} u_k^{\top} R u_k + (y_k - r_k)^{\top} Q (y_k - r_k) + (y_{N_h} - r_{N_h})^{\top} Q (y_{N_h} - r_{N_h})$$
(15)

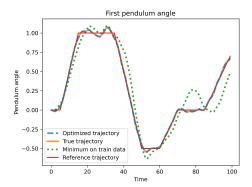
for positive-definite weight matrices  $R \in \mathbb{S}^{n_u \times n_u}$  and  $Q \in \mathbb{S}^{n_y \times n_y}$  and feasible input-output set

$$\mathcal{M}_{\mathrm{IO}} = \left\{ (\boldsymbol{u}, \boldsymbol{y}) \in (\mathbb{R}^{n_u})^{N_{\mathrm{h}}} \times (\mathbb{R}^{n_y})^{N_{\mathrm{h}}+1} \mid \begin{array}{c} \mathrm{exists} \ \boldsymbol{x} = (x_0, \dots, x_{N_h}) \in (\mathbb{R}^{n_x})^{N_h+1} \\ \mathrm{with} \ (\boldsymbol{u}, \boldsymbol{x}, \boldsymbol{y}) \ \mathrm{satisfying} \ (13) \end{array} \right\} \ .$$

Under certain smoothness assumptions on the dynamics f and g the set  $\mathcal{M}_{IO}$  is a (non-compact) embedded smooth submanifold of  $(\mathbb{R}^{n_u})^{N_h} \times (\mathbb{R}^{n_y})^{N_h+1}$ . The problem (14) is ubiquitous in receding horizon control applications such as model predictive control (MPC) and used for e.g. autonomous driving Vu et al. (2021), motion planning Cohen et al. (2020), optimizing HVAC system energy efficiency Serale et al. (2018) and inventory control Kostić (2009). In many of these applications the dynamics (13) governing the system are *not known* explicitly. Instead, in data-driven control Dörfler (2023a;b); Markovsky et al. (2023) one assumes that (13) is given *implicitly* by a finite number of measured input-output trajectories

$$\mathcal{D}_{\text{train}} = \{(\boldsymbol{u}_i, \boldsymbol{y}_i) \mid i = 1, \dots, N_{\text{data}}\} \subseteq \mathcal{M}_{\text{IO}},$$

where the input u is persistently exciting Willems et al. (2005), e.g. given by (white) noise Ljung (1999). In particular the so-called *system behavior* manifold  $\mathcal{M}_{IO}$  Willems & Polderman (1997) is



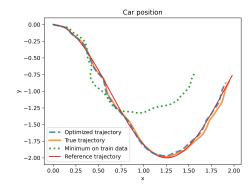


Figure 2: Denoising Riemannian gradient descent: Angle of the first pendulum (left) and unicycle car position (right) with the optimized output trajectory  $y^*$  (blue, dashed), the true system trajectory  $y^{\text{true}}$  (orange), the initial trajectory  $y_0$  (green, dotted) and the reference trajectory y (red)

given by samples from a distribution  $\mu$  on its in- and outputs and fits precisely into our framework of data-driven Riemannian optimization (1). We test our proposed denoising Riemannian gradient descent on two classical systems from the control domain: The discretized double pendulum system and the unicycle car model LaValle (2006) (see Appendix G.2.1 detailed information on the systems and the particular choice of r, R and Q in (15)), each on a horizon of  $N_h=100$ . To apply our proposed methods, we train a diffusion model (see Appendix G.2.2 for implementation details) on the measured trajectories  $\mathcal{D}_{\text{train}}$  and solve (14) via the denoising Riemannian gradient descent to obtain a solution  $(u^*, y^*)$ . As initial values we take the trajectories from the training set that minimize the objective cost f, i.e.  $(u_0, y_0) = \arg\min_{(u, y) \in \mathcal{D}_{\text{train}}} f(u, y)$ . Note that, as seen in Section 5, in general the final iterate will not exactly lie on the input-output manifold, i.e.  $(u^*, y^*) \notin \mathcal{M}_{\text{IO}}$ . To account for this deviation we back-test our generated input trajectory  $u^*$  by implementing it on the true underlying system (13) to obtain the real output  $y^{\text{true}}$ .

**Results and discussion:** In Figure 4 (in Appendix G.2.3) we depict the evolution of the objective value w.r.t. iteration count and in Figure 2 we depict the final optimizing trajectories  $y^*$  and  $y^{\rm true}$ . We can observe that the error  $\|y^*-y^{\rm true}\|$  is small, which shows that  $(u^*,y^*)$  is close to the true system behavior  $\mathcal{M}_{\rm IO}$ . Moreover, we can see that the trajectory  $y^{\rm true}$  tracks the corresponding reference r much better than the train set minimum  $y_0$ , which shows a generalization capability of our diffusion model. We from Figure 4 (left) that the current objective can depart from the true objective significantly. This is due to the iterates deviating from  $\mathcal{M}_{\rm IO}$ . In this example the algorithm (DRGD) recovers and we have found it to be robust w.r.t. moderate deviations from the manifold. Note that we have set a iteration budget of  $N_{\rm iter}=3000$  and  $N_{\rm iter}=2500$ , respectively, while the objective is still decreasing. Accelerating the convergence of DRGD is left for future work.

## 7 CONCLUSION AND FUTURE WORK

In this paper, we show that the denoising score and its Jacobian allow to perform manifold operations such as the closest-point and tangent space projection without the explicit knowledge of the manifold. We then propose a landing flow for the corresponding manifold-constrained Riemannian optimization problem and show that its limit points approximate critical points of the original problem. Moreover, we investigate the approximate version of the Riemannian gradient descent and provide a bound on the average-gradient-norm of its iterates, which converges to zero as the manifold operations become more exact. We apply both algorithms on known manifolds and to finite-horizon reference tracking in the domain of data-driven control. Future work will consist of deriving error bounds for this flow when the denoising score is trained with a non-zero  $L^2$ -error as well as the study of more sophisticated classical Riemannian optimization algorithms such as Newton and trust region methods when using the approximate manifold operations with the trained score to accelerate convergence.

# REFERENCES

- Theagenis J Abatzoglou. The minimum norm projection on c2-manifolds in rn. *Transactions of the American Mathematical Society*, 243:115–122, 1978.
- Pierre Ablin and Gabriel Peyré. Fast and accurate optimization on the orthogonal manifold without retraction. In *International Conference on Artificial Intelligence and Statistics*, pp. 5636–5657. PMLR, 2022.
- P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix manifolds*.
   Princeton University Press, 2008.
  - Felipe Alvarez, Jerome Bolte, and Olivier Brahic. Hessian riemannian gradient flows in convex programming. *SIAM journal on control and optimization*, 43(2):477–501, 2004.
  - Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows: in metric spaces and in the space of probability measures.* Springer, 2005.
  - Nitin Bansal, Xiaohan Chen, and Zhangyang Wang. Can we gain more from orthogonality regularizations in training deep networks? *Advances in Neural Information Processing Systems*, 31, 2018.
  - Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6):1373–1396, 2003.
  - Mikhail Belkin and Partha Niyogi. Towards a theoretical foundation for laplacian-based manifold methods. In *International conference on computational learning theory*, pp. 486–500. Springer, 2005.
  - Nicolas Boumal. *An introduction to optimization on smooth manifolds*. Cambridge University Press, 2023.
  - Nicolas Boumal, Pierre-Antoine Absil, and Coralia Cartis. Global rates of convergence for nonconvex optimization on manifolds. *IMA Journal of Numerical Analysis*, 39(1):1–33, 2019.
  - Paul Breiding and Nick Vannieuwenhoven. The condition number of riemannian approximation problems. *SIAM Journal on Optimization*, 31(1):1049–1077, 2021.
  - Lawrence Cayton et al. Algorithms for manifold learning. *Univ. of California at San Diego Tech. Rep*, 12(1-17):1, 2005.
  - Long Chen, Emre Oezkaya, Jan Rottmayer, Nicolas R Gauger, Zebang Shen, and Yinyu Ye. Adjoint-based aerodynamic shape optimization with a manifold constraint learned by diffusion models. *arXiv preprint arXiv:2507.23443*, 2025.
  - Hyungjin Chung, Jeongsol Kim, Michael T Mccann, Marc L Klasky, and Jong Chul Ye. Diffusion posterior sampling for general noisy inverse problems. *arXiv preprint arXiv:2209.14687*, 2022a.
  - Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Improving diffusion models for inverse problems using manifold constraints. *Advances in Neural Information Processing Systems*, 35:25683–25696, 2022b.
  - Mitchell R Cohen, Khairi Abdulrahim, and James Richard Forbes. Finite-horizon lqr control of quadrotors on se\_2(3). *IEEE Robotics and Automation Letters*, 5(4):5748–5755, 2020.
  - Ronald R Coifman and Stéphane Lafon. Diffusion maps. *Applied and computational harmonic analysis*, 21(1):5–30, 2006.
  - Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
  - Vincent Divol. Measure estimation on manifolds: an optimal transport approach. *Probability Theory and Related Fields*, 183(1):581–647, 2022.

- Carles Domingo-Enrich, Michal Drozdzal, Brian Karrer, and Ricky TQ Chen. Adjoint matching:
   Fine-tuning flow and diffusion generative models with memoryless stochastic optimal control.
   arXiv preprint arXiv:2409.08861, 2024.
  - Florian Dörfler. Data-driven control: Part one of two: A special issue sampling from a vast and dynamic landscape. *IEEE Control Systems Magazine*, 43(5):24–27, 2023a.
  - Florian Dörfler. Data-driven control: Part two of two: Hot take: Why not go with models? *IEEE Control Systems Magazine*, 43(6):27–31, 2023b.
  - Thai Duong, Abdullah Altawaitan, Jason Stanley, and Nikolay Atanasov. Port-hamiltonian neural ode networks on lie groups for robot dynamics learning and control. *IEEE Transactions on Robotics*, 2024.
  - Bradley Efron. Tweedie's formula and selection bias. *Journal of the American Statistical Association*, 106(496):1602–1614, 2011.
  - Bálint Farkas and Sven-Ake Wegner. Variations on barbălat's lemma. *The American Mathematical Monthly*, 123(8):825–830, 2016.
  - Charles Fefferman, Sanjoy Mitter, and Hariharan Narayanan. Testing the manifold hypothesis. *Journal of the American Mathematical Society*, 29(4):983–1049, 2016.
  - Alexandros Graikos, Nikolay Malkin, Nebojsa Jojic, and Dimitris Samaras. Diffusion models as plug-and-play priors. *Advances in Neural Information Processing Systems*, 35:14715–14728, 2022.
  - Enkelejd Hashorva, Dmitry Korshunov, and Vladimir I Piterbarg. On laplace asymptotic method, with application to random chaos. 2015.
  - Uwe Helmke and John B Moore. *Optimization and dynamical systems*. Springer Science & Business Media, 2012.
  - Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint* arXiv:2207.12598, 2022.
    - Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
    - Jiang Hu, Xin Liu, Zai-Wen Wen, and Ya-Xiang Yuan. A brief introduction to manifold optimization. *Journal of the Operations Research Society of China*, 8(2):199–248, 2020.
    - Chii-Ruey Hwang. Laplace's method revisited: weak convergence of probability measures. *The Annals of Probability*, pp. 1177–1182, 1980.
    - Tadeusz Inglot and Piotr Majerski. Simple upper and lower bounds for the multivariate laplace approximation. *Journal of Approximation Theory*, 186:1–11, 2014.
    - Amin G Jaffer and Someshwar C Gupta. On relations between detection and estimation of discrete time processes. *Information and Control*, 20(1):46–54, 1972.
    - Hamid Kamkari, Brendan Ross, Rasa Hosseinzadeh, Jesse Cresswell, and Gabriel Loaiza-Ganem. A geometric view of data complexity: Efficient local intrinsic dimension estimation with diffusion models. *Advances in Neural Information Processing Systems*, 37:38307–38354, 2024.
- Hassan K Khalil and Jessy W Grizzle. *Nonlinear systems*, volume 3. Prentice hall Upper Saddle River, NJ, 2002.
  - Konstantin Kostić. Inventory control as a discrete system control for the fixed-order quantity system. *Applied Mathematical Modelling*, 33(11):4201–4214, 2009.
    - Tomasz M Lapinski. Multivariate laplace's approximation with estimated error and application to limit theorems. *Journal of Approximation Theory*, 248:105305, 2019.

- Rémi Laumont, Valentin De Bortoli, Andrés Almansa, Julie Delon, Alain Durmus, and Marcelo Pereyra. Bayesian imaging using plug & play priors: when langevin meets tweedie. *SIAM Journal on Imaging Sciences*, 15(2):701–737, 2022.
- Steven M LaValle. *Planning algorithms*. Cambridge university press, 2006.
- Kyowoon Lee and Jaesik Choi. Local manifold approximation and projection for manifold-aware diffusion planning. *arXiv* preprint arXiv:2506.00867, 2025.
- Tong Lin and Hongbin Zha. Riemannian manifold learning. *IEEE transactions on pattern analysis and machine intelligence*, 30(5):796–809, 2008.
- Lennart Ljung. System identification (2nd ed.): theory for the user. Prentice Hall PTR, USA, 1999. ISBN 0136566952.
- Gabriel Loaiza-Ganem, Brendan Leigh Ross, Rasa Hosseinzadeh, Anthony L Caterini, and Jesse C Cresswell. Deep generative models through the lens of the manifold hypothesis: A survey and new connections. *Transactions on Machine Learning Research*, 2024.
- Piotr Majerski. Simple error bounds for the multivariate laplace approximation under weak local assumptions. *arXiv preprint arXiv:1511.00302*, 2015.
- Ivan Markovsky, Linbin Huang, and Florian Dörfler. Data-driven control based on the behavioral approach: From theory to applications in power systems. *IEEE Control Systems Magazine*, 43 (5):28–68, 2023.
- Marina Meilă and Hanyu Zhang. Manifold learning: What, how, and why. *Annual Review of Statistics and Its Application*, 11(1):393–417, 2024.
- Yasunori Nishimori. Learning algorithm for independent component analysis by geodesic flows on orthogonal group. In *IJCNN'99*. *International Joint Conference on Neural Networks*. *Proceedings* (*Cat. No. 99CH36339*), volume 2, pp. 933–938. IEEE, 1999.
- Kazusato Oko, Shunta Akiyama, and Taiji Suzuki. Diffusion models are minimax optimal distribution estimators. In *International Conference on Machine Learning*, pp. 26517–26582. PMLR, 2023.
- Neal Patwari and Alfred O Hero. Manifold learning algorithms for localization in wireless sensor networks. In 2004 IEEE international conference on acoustics, speech, and signal processing, volume 3, pp. iii–857. IEEE, 2004.
- Scott Pesme, Giacomo Meanti, Michael Arbel, and Julien Mairal. Map estimation with denoisers: Convergence rates and guarantees. *arXiv preprint arXiv:2507.15397*, 2025.
- Phil Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. The intrinsic dimension of images and its impact on learning. In *International Conference on Learning Representations*.
- Herbert E Robbins. An empirical bayes approach to statistics. In *Breakthroughs in Statistics: Foundations and basic theory*, pp. 388–394. Springer, 1992.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241. Springer, 2015.
- Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear embedding. *science*, 290(5500):2323–2326, 2000.
- Hiroyuki Sato. Riemannian optimization and its applications, volume 670. Springer, 2021.
- Sholom Schechtman, Daniil Tiapkin, Michael Muehlebach, and Eric Moulines. Orthogonal directions constrained gradient method: from non-linear equality constraints to stiefel manifold. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 1228–1258. PMLR, 2023.

- Gianluca Serale, Massimo Fiorentini, Alfonso Capozzoli, Daniele Bernardini, and Alberto Bemporad. Model predictive control (mpc) for enhancing building and hvac system energy efficiency: Problem formulation, applications and opportunities. *Energies*, 11(3):631, 2018.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv* preprint *arXiv*:2011.13456, 2020.
- Jan Pawel Stanczuk, Georgios Batzolis, Teo Deveney, and Carola-Bibiane Schönlieb. Diffusion models encode the intrinsic dimension of data manifolds. In *Forty-first International Conference on Machine Learning*, 2024.
- Rong Tang and Yun Yang. Adaptivity of diffusion models to manifold structures. In *International Conference on Artificial Intelligence and Statistics*, pp. 1648–1656. PMLR, 2024.
- Wenpin Tang and Hanyang Zhao. Score-based diffusion models via stochastic differential equations. *Statistic Surveys*, 19:28–64, 2025.
- Joshua B Tenenbaum, Vin de Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
- Bart Vandereycken. Low-rank matrix completion by riemannian optimization. *SIAM Journal on Optimization*, 23(2):1214–1236, 2013.
- Singanallur V Venkatakrishnan, Charles A Bouman, and Brendt Wohlberg. Plug-and-play priors for model based reconstruction. In *2013 IEEE global conference on signal and information processing*, pp. 945–948. IEEE, 2013.
- Enrico Ventura, Beatrice Achilli, Gianluigi Silvestri, Carlo Lucibello, and Luca Ambrogioni. Manifolds, random matrices and spectral gaps: The geometric phases of generative diffusion. *arXiv* preprint arXiv:2410.05898, 2024.
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.
- Trieu Minh Vu, Reza Moezzi, Jindrich Cyrus, and Jaroslav Hlava. Model predictive control for autonomous driving vehicles. *Electronics*, 10(21):2593, 2021.
- Jan C Willems and Jan W Polderman. *Introduction to mathematical systems theory: a behavioral approach*, volume 26. Springer Science & Business Media, 1997.
- Jan C Willems, Paolo Rapisarda, Ivan Markovsky, and Bart LM De Moor. A note on persistency of excitation. *Systems & Control Letters*, 54(4):325–329, 2005.
- Yue Xie and Stephen J Wright. Complexity of proximal augmented lagrangian for nonconvex optimization with nonlinear equality constraints. *Journal of Scientific Computing*, 86(3):38, 2021.
- Boda Zheng, Abhijith Moni, Weigang Yao, and Min Xu. Manifold learning for aerodynamic shape design optimization. *Aerospace*, 12(3):258, 2025.
- Yufan Zhou, Changyou Chen, and Jinhui Xu. Learning manifold implicitly via explicit heat-kernel learning. *Advances in Neural Information Processing Systems*, 33:477–487, 2020.

## A MORE ON MANIFOLDS AND DISTANCE FUNCTIONS

In addition to the notation introducted in Section 2.1, we note that  $\tau_{\mathcal{M}}$  is also the largest  $\tau \geq 0$  such that the map  $\{(p,v) \in \mathbb{N} \,\mathcal{M} \mid \|v\| < \tau\} \to \mathbb{R}^d : (p,v) \mapsto p+v$  is a diffeomorphism. By a tubular neighborhood of radius  $\tau \in (0,\tau_{\mathcal{M}}]$  we mean a set of the form  $\mathcal{T}(\tau) = \{p+v \mid p \in \mathcal{M}, v \in \mathbb{N}_p \,\mathcal{M}, \|v\| < \tau\}$ . Moreover, for  $x \in \mathbb{R}^d$  let  $\mathrm{dist}_{\mathcal{M}}(x) = \inf_{p \in \mathcal{M}} \|x-p\|$  denote the distance function so that  $\mathrm{d}(x) = \frac{1}{2} \, \mathrm{dist}_{\mathcal{M}}(x)^2$ . The second fundamental form of  $\mathcal{M}$  at a point  $p \in \mathcal{M}$  will be denoted by  $\mathrm{II}_p$  and is a symmetric bilinear map  $\mathrm{II}_p : \mathrm{T}_p \,\mathcal{M} \times \mathrm{T}_p \,\mathcal{M} \to \mathrm{N}_p \,\mathcal{M}$  intrinsic to the manifold  $\mathcal{M}$ . Fixing some  $u \in \mathrm{N}_p \,\mathcal{M}$  we also define the directed second fundamental form

  $\mathbb{I}_p^u: \mathrm{T}_p \, \mathcal{M} \times \mathrm{T}_p \, \mathcal{M} \to \mathbb{R}: (v,w) \mapsto \langle \mathbb{I}_p(v,w), u \rangle_{\mathrm{N}_p \, \mathcal{M}}$ . The Weingarten map  $S_p^u$  at a point  $p \in \mathcal{M}$  in the direction  $u \in \mathrm{N}_p \, \mathcal{M}$  is defined as the unique self-adjoint linear operator  $S_p^u: \mathrm{T}_p \, \mathcal{M} \to \mathrm{T}_p \, \mathcal{M}$  such that  $\langle w, S_p^u(v) \rangle_{\mathrm{T}_p \, \mathcal{M}} = \mathbb{I}_p^u(v,w)$  for all  $v,w \in \mathrm{T}_p \, \mathcal{M}$ . A useful operator that has been studied in Abatzoglou (1978); Breiding & Vannieuwenhoven (2021); Alvarez et al. (2004) is

$$H_x = I_{T_{\pi(x)}\mathcal{M}} + S_{\pi(x)}^{\pi(x)-x} : T_{\pi(x)}\mathcal{M} \to T_{\pi(x)}\mathcal{M}.$$
 (16)

In Breiding & Vannieuwenhoven (2021) it has been shown (see Lemma 8) that  $H_x$  is invertible on  $\mathcal{T}$ . Now we have the following useful identities that hold for  $x \in \mathcal{T}$ 

$$\begin{split} \pi'(x) &= \mathrm{I}_{\mathrm{T}_{\pi(x)} \, \mathcal{M}} \, H_x^{-1} \, \mathrm{P}_{\mathrm{T}_{\pi(x)} \, \mathcal{M}} \\ \nabla \, \mathrm{d}(x) &= x - \pi(x) \\ \nabla^2 \, \mathrm{d}(x) &= I - \pi'(x) = I - \mathrm{I}_{\mathrm{T}_{\pi(x)} \, \mathcal{M}} \, H_x^{-1} \, \mathrm{P}_{\mathrm{T}_{\pi(x)} \, \mathcal{M}} \; . \end{split}$$

For any  $p \in \mathcal{M}$  the map  $\operatorname{pr}_p : \mathcal{M} \to \operatorname{T}_p \mathcal{M} : q \mapsto \operatorname{P}_{\operatorname{T}_p \mathcal{M}}(q-p)$  is a local diffeomorphism at p with inverse  $\psi_p$  defined on  $B_{\tau_{\mathcal{M}}/4}^{\operatorname{T}_p \mathcal{M}}(0) := B_{\tau_{\mathcal{M}}/4}(0) \cap \operatorname{T}_p \mathcal{M}$ . Following Divol (2022) we define  $\mathcal{M}_k(\tau, M)$  as the set of all  $C^k$ -manifolds  $\mathcal{M}$  as above such that  $\tau_{\mathcal{M}} > \tau$  and  $\sup_{p \in \mathcal{M}} \|\psi_p\|_{C^k} \leq M$ . For the class  $\mathcal{M}_k(\tau, M)$ , a manifold  $\mathcal{M} \in \mathcal{M}_k(\tau, M)$  and  $p \in \mathcal{M}$ , we denote by  $\psi_p$  always the inverse of the orthogonal projection  $\operatorname{pr}_p$  (also called  $\operatorname{Monge}$  or  $\operatorname{graph\ chart}$ ) restricted to the particular neighborhood  $B_{\min\{\tau_{\mathcal{M}}, M\}/4}^{\operatorname{T}_p \mathcal{M}}(0)$ , which will be (isometrically) identified with the ball  $B_{\min\{\tau_{\mathcal{M}}, M\}/4}(0) \subseteq \mathbb{R}^k$ . We make frequent use of the following useful result from Divol (2022).

**Lemma 7** (Lemma A.1 in Divol (2022)). Suppose  $\mathcal{M} \in \mathcal{M}_k(\tau, M)$  and  $p \in \mathcal{M}$ . Then  $\psi_p : B_{\min\{\tau_{\mathcal{M}}, M\}/4}^{\mathrm{T}_p \mathcal{M}}(0) \to \mathcal{M}$  is well-defined,  $C^k$ -smooth and the following holds:

- (i) For all  $r \leq \min\{\tau_{\mathcal{M}}, M\}/4$  it holds that  $B_r(p) \cap \mathcal{M} \subseteq \psi_p(B_r^{\mathrm{T}_p \mathcal{M}}(0)) \subseteq B_{8r/7}(p) \cap \mathcal{M}$ . For  $z \in B_{\min\{\tau_{\mathcal{M}}, M\}/4}^{\mathrm{T}_p \mathcal{M}}(0)$  it holds that  $\|z\| \leq \|\psi_p(z) - p\| \leq 8\|z\|/7$ .
- (ii) There exists a map  $W_p: B^{\mathrm{T}_p\mathcal{M}}_{\min\{\tau_{\mathcal{M}},M\}/4}(0) \to \mathrm{N}_p\mathcal{M}$  with  $W_p'(0) = 0$  and such that  $\psi_p(z) = p + z + W_p(z)$  and  $\|W_p(z)\| \le M\|z\|^2$  for all  $z \in B^{\mathrm{T}_p\mathcal{M}}_{\min\{\tau_{\mathcal{M}},M\}/4}(0)$ .
- (iii) For  $G_{\psi_p}: B_{\min\{\tau_{\mathcal{M}}, M\}/4}^{\mathrm{T}_p \mathcal{M}}(0) \to \mathbb{R}: z \mapsto \sqrt{\det \psi_p'(z)^\top \psi_p'(z)}$  it holds that  $G_{\psi_p}(0) = 1$  and  $\nabla G_{\psi_p}(0) = 0$ .

Note that for the graph chart  $\psi_n$  we always have

$$\psi'_{p}(0) = \mathbf{I}_{\mathbf{T}_{p}\mathcal{M}}, \quad \psi''_{p}(0)[\cdot, \cdot] = W''_{p}(0)[\cdot, \cdot] = \mathbf{I}_{p}(\cdot, \cdot),$$
 (17)

and hence  $\|\psi_p'(0)\| \le 1$  and  $\|\psi_p''(0)\| = \|\mathbf{II}_p\| \le 1/\tau_{\mathcal{M}}$ .

## A.1 Properties of (16)

We study the invertibility and boundedness of the operator (16). For this purpose, let us recall first the definition of the (normalized) curvature radius of  $\mathcal{M}$  at p in the direction of  $u \in N_p \mathcal{M}$ :

$$\frac{1}{\rho(p,u)} = \max_{\substack{v \in \mathcal{T}_p \, \mathcal{M} \\ \mathbb{I}_u^u(v,v) > 0}} \frac{\mathbb{I}_p^{u/\|u\|}(v,v)}{\|v\|^2} = \max(\text{eig}(S_p^{u/\|u\|}) \cup \{0\}) \,.$$

If  $S_p^u$  has only non-positive eigenvalues, then  $\mathcal{M}$  is curved away from the unit vector u and thus the curvature radius is infinite. Moreover, we define the (normalized) curvature of  $\mathcal{M}$  to be

$$\kappa_p^{\mathcal{M}}(u) = \max \left| \operatorname{eig}(S_p^{u/\|u\|}) \right| \text{ for } u \in \mathcal{N}_p \mathcal{M},$$

and the maximal curvature by

$$\kappa_{\mathcal{M}} = \max_{(p,u)\in\mathbb{N}} \kappa_p^{\mathcal{M}}(u). \tag{18}$$

We have the following

**Lemma 8.** The operator (16) is invertible on  $\mathcal{T}(\tau_{\mathcal{M}})$ . Moreover, (16) satisfies <sup>1</sup>

$$||P_0(x)|| = ||H_x^{-1}|| = \left(1 - \frac{||\pi(x) - x||}{\rho(\pi(x), x - \pi(x))}\right)^{-1} \le \frac{1}{1 - ||x - \pi(x)|| \kappa_{\mathcal{M}}} \le \frac{1}{1 - ||x - \pi(x)|| / \tau_{\mathcal{M}}},$$

and if  $\psi$  is a local parametrization of  $\mathcal{M}$  with  $\psi(0) = p$ , then

$$P_0(x) = \psi'(0) \left( \psi'(0)^{\top} \psi'(0) + \sum_{i=1}^d (p-x)_i \nabla^2 \psi_i(0) \right)^{-1} \psi'(0)^{\top}.$$

In particular, for any  $\tau \in (0, \tau_{\mathcal{M}})$ , it holds

$$\sup_{x\in\mathcal{T}(\tau)} \lVert H_x^{-1}\rVert < \infty .$$

*Proof.* In (Breiding & Vannieuwenhoven, 2021, Lemma A.2) the following condition has been established: Let  $\mathcal{S} = \{(a,p) \in \mathbb{R}^n \times \mathcal{M} \mid a-p \in \mathrm{N}_p \mathcal{M}\}$ . Then  $\mathcal{S}$  is diffeomorphic to the normal bundle  $\mathrm{N}\,\mathcal{M}$  via the diffeomorphism  $\Phi: \mathrm{N}\,\mathcal{M} \to \mathcal{S}: (v,p) \mapsto (p+\mathrm{I}_{\mathrm{N}_p}\,\mathcal{M}(v),p)$ . Consider the operator  $\mathrm{\Pi}: \mathcal{S} \to \mathbb{R}^n: (a,p) \mapsto a$  and the domain where its differential is invertible  $\mathcal{W} = \{(a,p) \in \mathcal{S} \mid \Pi'(a,p): T_{(a,p)}\,\mathcal{S} \to \mathbb{R}^n \text{ invertible}\}$ . Then  $H_x$  is invertible iff  $(x,\pi(x)) \in \mathcal{W}$ . But  $\mathrm{\Pi} \circ \Phi: \mathrm{N}\,\mathcal{M} \mapsto \mathbb{R}^n: (v,p) \mapsto p+\mathrm{I}_{\mathrm{N}_p}\,\mathcal{M}(v)$  being a diffeomorphism (and thus having an inveritble differential) is precisely the condition in the definition of the tubular neighborhood  $\mathcal{T}$ . The formulas for  $P_0(x)$  in local coordinates as well as  $\|P_0(x)\|$  are given in (Abatzoglou, 1978, Theorem 4.1, Corollary 4.1). To see that  $P_0(x)$  is bounded on  $\mathcal{T}(\tau)$  for any  $\tau \in (0,\tau_{\mathcal{M}})$  it sufficies to note that in the tubular neighborhood  $\mathcal{T} = \mathcal{T}(\tau_{\mathcal{M}})$  we always have  $\|\pi(x) - x\| < \rho(\pi(x), x - \pi(x))$  and that  $\overline{\mathcal{T}(\tau)}$  is a compact subset thereof. The second inequality follows from  $1/\rho(p,u) \leq \kappa_{\mathcal{M}}$ .

Now let us derive bounds for the quantity  $||P_0(\pi(x)) - P_0(x)||$  when  $x \in \mathcal{T}(\tau_{\mathcal{M}})$ .

**Lemma 9.** If  $x \in \mathcal{T}(\tau_{\mathcal{M}})$ , then

$$||P_0(\pi(x)) - P_0(x)|| \le \kappa_{\pi(x)}^{\mathcal{M}}(x - \pi(x)) \left( 1 - \frac{||\pi(x) - x||}{\rho(\pi(x), x - \pi(x))} \right)^{-1} ||x - \pi(x)||,$$

$$\le \frac{||x - \pi(x)|| \kappa_{\mathcal{M}}}{1 - ||x - \pi(x)|| \kappa_{\mathcal{M}}} \le \frac{||x - \pi(x)|| / \tau_{\mathcal{M}}}{1 - ||x - \pi(x)|| / \tau_{\mathcal{M}}}.$$

*Proof.* We clearly have for  $x \in \mathcal{T}$ 

$$P_0(\pi(x)) - P_0(x) = I_{T_{\pi(x)} \mathcal{M}}(I - H_x^{-1}) P_{T_{\pi(x)} \mathcal{M}}.$$

Moreover,  $(I - H_x^{-1}) : T_{\pi(x)} \mathcal{M} \to T_{\pi(x)} \mathcal{M}$  is symmetric with eigenvalues

$$\operatorname{eig}(I - H_x^{-1}) = \left\{ \frac{\zeta}{1 + \zeta} \mid \zeta \in \operatorname{eig}(S_{\pi(x)}^{\pi(x) - x}) \right\} = \left\{ \frac{-\|\pi(x) - x\|\zeta}{1 - \|\pi(x) - x\|\zeta} \mid \zeta \in \operatorname{eig}(S_{\pi(x)}^u) \right\},$$

where  $u = \frac{x - \pi(x)}{\|x - \pi(x)\|}$ . Thus

$$||P_0(\pi(x)) - P_0(x)|| \le \max_{\zeta \in eig(S_{\pi(x)}^u)} \left| \frac{\zeta}{1 - ||x - \pi(x)||\zeta} \right| ||x - \pi(x)||$$

$$\le ||S_{\pi(x)}^u|| \left( 1 - \frac{||\pi(x) - x||}{\rho(\pi(x), x - \pi(x))} \right)^{-1} ||x - \pi(x)||,$$

which shows the first inequality. The second and third inequalities follow from  $1/\rho(p,u) \le \kappa_{\mathcal{M}} \le 1/\tau_{\mathcal{M}}$ .

The next lemma establishes a bound on the Lipschitz-constant of  $P_0$  of some  $\mathcal{M} \in \mathcal{M}_k(\tau, M)$  in terms of M and  $\tau_{\mathcal{M}}$ .

<sup>&</sup>lt;sup>1</sup>In (Breiding & Vannieuwenhoven, 2021, Theorem 4.3) the quantity  $||H_x^{-1}||$  has been shown to equal the condition number of a certain critical point problem associated with  $\mathcal{M}$ .

**Lemma 10.** If  $\mathcal{M} \in \mathcal{M}_k(\tau, M)$ , then

$$\sup_{x \in \mathcal{T}(\tau)} \lVert P_0'(x) \rVert \leq (\frac{1}{1 - \tau/\tau_{\mathcal{M}}})^2 \left( (\frac{3}{\tau_{\mathcal{M}}} + \tau M) (\frac{1}{1 - \tau/\tau_{\mathcal{M}}})^2 + \frac{2}{\tau_{\mathcal{M}}} \right)$$

*Proof.* Let  $x \in \mathcal{T}(\tau)$  where  $\tau \in (0, \tau_{\mathcal{M}})$ . First let us relate the quantity  $\pi(x)$  and a fixed chart  $\psi : \mathcal{V} \to \mathcal{U}$  with  $\pi(x) \in \mathcal{U}$  for all  $x \in \mathcal{W}$  for some (small enough) open set  $\mathcal{W} \subseteq \mathbb{R}^d$ . The map  $h_{\psi}(z) = \frac{1}{2} \|x - \psi(z)\|^2$  attains its minimum in some  $z_{\psi}(x) \in \mathcal{V}$  and if  $F^{\psi}(z, x) = h'_{\psi}(z) = \psi'(z)(\psi(z) - x)$ , then  $F^{\psi}(z_{\psi}(x), x) = 0$ . Moreover, we have  $\pi(x) = \psi(z_{\psi}(x))$  and hence for  $w \in \mathbb{R}^d$ 

$$P_0'(x)[w,w] = \pi''(x)[w,w] = \psi''(z_{\psi}(x))[z_{\psi}'(x)[w], z_{\psi}'(x)[w]] + \psi'(z_{\psi}(x))[z_{\psi}''(x)[w,w]],$$

and hence

$$||P'_0(x)|| \le ||\psi''(z_{\psi}(x))||||z'_{\psi}(x)||^2 + ||\psi'(z_{\psi}(x))||||z''_{\psi}(x)||$$

By the implicit function theorem we can bound the derivatives of  $z_{\psi}$  in terms of derivatives of  $F^{\psi}$ . Indeed, we have (here again  $w \in \mathbb{R}^d$  and  $v \in \mathbb{R}^k$  are place-holder vectors to express the differentials)

$$0 = F_z^{\psi}(z_{\psi}(x), x)[v, z'_{\psi}(x)[w]] + F_x^{\psi}(z_{\psi}(x), x)[v, w],$$

$$0 = F_z^{\psi}(z_{\psi}(x), x)[v, z''_{\psi}(x)[w, w]] + F_{zz}^{\psi}(z_{\psi}(x), x)[v, z'_{\psi}(x)[w], z'_{\psi}(x)[w]]$$

$$+ 2F_{xz}^{\psi}(z_{\psi}(x), x)[v, z_{\psi}(x)[w], w] + F_{xx}^{\psi}(z_{\psi}(x), x)[v, w, w],$$
(19)

with

$$\begin{split} F_z^{\psi}(z,x)[v,v] &= \langle \psi'(z)[v], \psi'(z)[v] \rangle + \langle \psi(z) - x, \psi''(z)[v,v] \rangle \;, \\ F_x^{\psi}(z,x)[w,v] &= \langle \psi'(z)[v], w \rangle \;, \\ F_{zz}^{\psi}(z,x)[v,v,v] &= 3 \left< \psi''(z)[v,v], \psi'(z)[v] \right> + \left< \psi(z) - x, \psi'''(z)[v,v,v] \right> \;, \\ F_{xx}^{\psi}(z,x)[w,v,v] &= \left< \psi''(z)[v,v], w \right> \;, \\ F_{xx}^{\psi}(z,x)[w,v,v] &= 0 \;. \end{split}$$

Now we set  $\psi=\psi_p$  for the graph chart defined on  $\mathcal{V}=B^{\mathrm{T}_p\,\mathcal{M}}_{\min\{\tau,M\}/4}(0)$ , where  $p=\pi(x)$ . In this case  $z=z_\psi(x)=0$  and  $\|F_z^\psi(0,x)^{-1}\|=\|P_0(x)\|\leq (1-\tau/\tau_\mathcal{M})^{-1}$  by Lemma 8. Solving for  $z_\psi'(x)$  and  $z_\psi''(x)$  in (19), using (17) with  $\|\mathrm{II}_p\|\leq 1/\tau_\mathcal{M}$  and taking norms yields then

$$||z'_{\psi}(x)|| \leq \frac{1}{1 - \tau/\tau_{\mathcal{M}}},$$

$$||z''_{\psi}(x)|| \leq \frac{1}{1 - \tau/\tau_{\mathcal{M}}} \left( (\frac{3}{\tau_{\mathcal{M}}} + \tau M) ||z'_{\psi}(x)||^{3} + \frac{1}{\tau_{\mathcal{M}}} ||z'_{\psi}(x)|| \right).$$

Plugging this back into the upper bound of  $||P_0'(x)||$  and using once again (17) finishes the proof

The next lemma establishes a lower bound on the distance between a point x in a tubular neighborhood and any other point that is sufficiently away from  $p = \pi(x)$  uniformly in x.

**Lemma 11.** Let  $\tau \in (0, \tau_{\mathcal{M}})$  and let  $\eta > 0$ . Then

$$\inf_{x \in \mathcal{T}(\tau)} \frac{1}{2} \operatorname{dist}_{\mathcal{M} \setminus B_{\eta}(\pi(x))}(x)^{2} - \frac{1}{2} \operatorname{dist}_{\mathcal{M}}(x)^{2} \ge \frac{1}{2} \frac{\tau_{\mathcal{M}} - \tau}{\tau_{\mathcal{M}} + \tau} \eta^{2} > 0.$$

*Proof.* Let  $\tau' = \frac{\tau + \tau_{\mathcal{M}}}{2} \in (0, \tau_{\mathcal{M}})$ . Take  $x \in \mathcal{T}(\tau)$  and set  $p = \pi(x)$  and  $y = p + \tau' \widehat{n}$  with  $\widehat{n} = \frac{x - p}{\|x - p\|}$ . Then  $y \in \mathcal{T}(\tau_{\mathcal{M}})$  with  $\mathrm{dist}_{\mathcal{M}}(y) = \tau'$  and thus  $\overline{B}_{\tau'}(y) \cap \mathcal{M} = \{p\}$ . Then, since  $\mathcal{M} \setminus B_{\eta}(p) \subseteq \mathbb{R}^d \setminus (\overline{B}_{\tau'}(y) \cup B_{\eta}(p))$ , it is sufficient to show

$$\inf_{q \in \mathbb{R}^d \setminus (\overline{B}_{\tau'}(y) \cup B_{\eta}(p))} \|x - q\|^2 - \|x - p\|^2 \ge \frac{\tau_{\mathcal{M}} - \tau}{\tau_{\mathcal{M}} + \tau} \eta^2.$$

To see this, note that  $q \in \mathbb{R}^d \setminus (\overline{B}_{\tau'}(y) \cup B_{\eta}(p))$  implies  $\|y-q\| \ge \tau' = \|y-p\|$  and  $\|q-p\| \ge \eta$ , as well as  $2 \langle p-q, y-q \rangle \ge \|p-q\|^2$ . Then, abbreviating  $t = \frac{\|x-p\|}{\|y-p\|} \in (0,1)$ , we have x = ty + (1-t)p = p + t(y-p) and

$$\begin{split} \|x-q\|^2 - \|x-p\|^2 &= \|ty + (1-t)p - q\|^2 - t^2 \|y-p\|^2 \\ &\geq \|t(y-q) + (1-t)(p-q)\|^2 - t^2 \|y-q\|^2 \\ &= (1-t)^2 \|p-q\|^2 + 2t(1-t) \left\langle p-q, y-q \right\rangle \\ &\geq (1-t) \|p-q\|^2 \\ &\geq \frac{\tau_{\mathcal{M}} - \tau}{\tau_{\mathcal{M}} + \tau} \eta^2 \,. \end{split}$$

#### A.2 DENSITIES ON MANIFOLDS

Given a manifold  $\mathcal{M}$  as in Section A and any chart  $\psi: \mathcal{V} \to \mathcal{U} \subseteq \mathcal{M}$ , the volume measure  $\operatorname{Vol}_{\mathcal{M}}$  is uniquely defined by

$$\operatorname{Vol}_{\mathcal{M}}(E) = \int_{\psi^{-1}(E)} G_{\psi}(z) \, \mathrm{d}z \ \text{ for } E \subseteq \mathcal{U} \text{ Borel measurable}.$$

Here  $G_{\psi}(z) = \sqrt{\det \psi'(z)^{\top} \psi'(z)}$  and  $\operatorname{Vol}_{\mathcal{M}}$  is independent of the chart  $\psi$ . If  $\mu \in \mathcal{P}(\mathbb{R}^d)$  is absolutely continuous w.r.t.  $\operatorname{Vol}_{\mathcal{M}}$ , then for the density  $\mu(y) := \frac{\operatorname{d} \mu}{\operatorname{d} \operatorname{Vol}_{\mathcal{M}}}(y)$  on  $\mathcal{M}$  we have the local representations

$$\mu(\psi(z)) = \frac{\mathrm{d}\,\lambda}{\mathrm{d}(\psi^{-1} \# \operatorname{Vol}_{\mathcal{M}})}(z)\,, \quad \lambda(z) = G_{\psi}(z) \mu(\psi(z))$$

with  $\lambda=\psi^{-1}\#\mu$  the pullback under the chart  $\psi$  with density  $\lambda(z)=\frac{\mathrm{d}\,\lambda}{\mathrm{d}\,m_{\mathcal{V}}}(z)$  w.r.t. the Lebesgue measure  $m_{\mathcal{V}}$  on  $\mathcal{V}$ . In particular when  $\mathcal{M}\in\mathcal{M}_k(\tau,M)$  with graph chart  $\psi_p$ , we have

$$\lambda(0) = \mu(p), \quad \nabla \lambda(0) = \operatorname{grad}_{\mathcal{M}} \mu(p) \in T_p \,\mathcal{M} \cong \mathbb{R}^k.$$
 (20)

## B More on the Stein score function

In this section we derive the representations (4). For  $\sigma \geq 0$  let  $\mathcal{N}(0, \sigma^2 I_d)$  denote the Gaussian distribution with mean zero and variance  $\sigma^2$  and for  $\sigma > 0$  the Gaussian (heat) kernel in the ambient space  $\mathbb{R}^d$  by

$$\varphi_{\sigma}(x) = \frac{1}{Z_{\sigma}} e^{-\|x\|^2/2\sigma^2}, \quad Z_{\sigma} = \frac{1}{(2\pi\sigma^2)^{d/2}}.$$

Then  $p_{\sigma} = \mathcal{N}(0, \sigma^2 I) * \mu$  has for  $\sigma > 0$  a fully supported  $C^{\infty}$ -density, denoted by  $p_{\sigma}$  as well, given by

$$p_{\sigma}(x) = \int_{\mathcal{M}} \varphi_{\sigma}(x - y) \, \mathrm{d}\mu(y), \quad x \in \mathbb{R}^d.$$

Let  $\nu_{x,\sigma}$  be the posterior of observing x under  $p_{\sigma}$  with prior  $\mu$ , i.e.

$$\nu_{x,\sigma}(E) = \frac{1}{p_\sigma(x)} \int_E \varphi_\sigma(x-y) \,\mathrm{d}\mu(y)\,, \quad E \subseteq \mathbb{R}^d \text{ Borel }.$$

We have the following representation

**Lemma 12.** For each  $\sigma > 0$  representations (4) hold.

Proof. Clearly

$$\nabla p_{\sigma}(x) = -\frac{1}{\sigma^2} \int_{\mathcal{M}} (x - y) \varphi_{\sigma}(x - y) \, \mathrm{d}\mu(y) \,, \quad x \in \mathbb{R}^d \,,$$

$$\nabla^2 p_{\sigma}(x) = -\frac{1}{\sigma^2} \left( p_{\sigma}(x) I - \frac{1}{\sigma^2} \int_{\mathcal{M}} (x - y) (x - y)^{\top} \varphi_{\sigma}(x - y) \, \mathrm{d}\mu(y) \right) \,, \quad x \in \mathbb{R}^d \,.$$

and hence

$$\nabla \log p_{\sigma}(x) = \frac{\nabla p_{\sigma}(x)}{p_{\sigma}(x)} = -\frac{1}{\sigma^2} (x - \mathbb{E}\nu_{x,\sigma}),$$

which yields the first representation in (4). Further

$$\frac{\nabla^2 p_{\sigma}(x)}{p_{\sigma}(x)} = -\frac{1}{\sigma^2} \left( I - \frac{1}{\sigma^2} \int_{\mathcal{M}} (x - y)(x - y)^{\top} d\nu_{x,\sigma}(y) \right)$$

and hence

$$\nabla^{2} \log p_{\sigma}(x) = \frac{\nabla^{2} p_{\sigma}(x)}{p_{\sigma}(x)} - \frac{\nabla p_{\sigma}(x) \nabla p_{\sigma}(x)^{\top}}{p_{\sigma}(x)^{2}}$$

$$= -\frac{1}{\sigma^{2}} \left( I - \frac{1}{\sigma^{2}} \int_{\mathcal{M}} (x - y)(x - y)^{\top} d\nu_{x,\sigma}(y) \right)$$

$$-\frac{1}{\sigma^{4}} (x - \mathbb{E}\nu_{x,\sigma})(x - \mathbb{E}\nu_{x,\sigma})^{\top}$$

$$= -\frac{1}{\sigma^{2}} I + \frac{1}{\sigma^{4}} \left( \int_{\mathcal{M}} y y^{\top} d\nu_{x,\sigma}(y) - (\mathbb{E}\nu_{x,\sigma})(\mathbb{E}\nu_{x,\sigma})^{\top} \right)$$

$$= -\frac{1}{\sigma^{2}} I + \frac{1}{\sigma^{4}} \operatorname{Cov}(\nu_{x,\sigma}),$$

which yields the second representation in (4).

## C VARIANCE-EXPLODING DIFFUSION MODELS

In the variance-exploding scheme of score-based diffusion models Song et al. (2020); Tang & Zhao (2025), one considers the following stochastic differential equation on the finite interval [0, T]:

$$dX_t = \sqrt{2t} dW_t, \quad t \in [0, T], \quad X_0 \sim \mu. \tag{21}$$

The solution  $X_t$  of this SDE is distributed according to  $X_t \sim p_{\sigma(t)}$ , where  $\sigma^2(t) = t^2$ . To sample from  $\mu$  one exploits the fact that the reverse SDE

$$\mathrm{d}\,\overline{X}_t = 2(T-t)\nabla\log p_{\sigma(T-t)}(\overline{X}_t)\,\mathrm{d}\,t + \sqrt{2(T-t)}\,\mathrm{d}\,W_t\,,\quad t\in[0,T]\,,\quad \overline{X}_0\sim p_{\sigma(T)}\,,\quad (22)$$

satisfies  $\overline{X}_t \sim p_{\sigma(T-t)}$  and in particular  $\overline{X}_0 \sim p_0 = \mu$ . Here the initial distribution is approximated by  $p_{\sigma(T)} \approx \mathcal{N}(0, \sigma(T)^2 I)$  and unknown score  $\nabla \log p_{\sigma(T-t)}$  is learned by minimizing the conditional score matching loss defined as

$$L_{\text{CSM}}(s) = \mathbb{E}_{t \sim \text{Unif}[0,T]} \mathbb{E}_{x_0 \sim \mu} \mathbb{E}_{x \sim p_{\sigma(t)}(\cdot \mid x_0)} \sigma(t)^2 \|s_{\sigma(t)}(x) - \nabla \log p_{\sigma(t)}(x \mid x_0)\|^2,$$
 (23)

and which admits the unique minimizer  $s_{\sigma}(x) = \nabla \log p_{\sigma}(x)$ . The loss  $L_{\text{CSM}}$  can be evaluated because  $\nabla \log p_{\sigma}(x \mid x_0) = -\frac{1}{\sigma^2}(x - x_0)$  is known explicitly.

## D PROOF OF THEOREM 1

#### D.1 NONASYMPTOTIC LAPLACE METHOD

In this section we derive some non-asymptotic error estimates for the Laplace method Hashorva et al. (2015); Hwang (1980), which concerns itself with the asymptotic of integrals of the form

$$\int_{\mathcal{M}} f(z)e^{-\frac{1}{\sigma^2}h(z)} \,\mathrm{d}z \quad (\sigma \to 0)$$

for some functions  $f,h:\mathcal{V}\to\mathbb{R}$  from an open set  $\mathcal{V}\subseteq\mathbb{R}^k$ , where h is non-negative and attains a unique minimum in, say,  $z=0\in\mathcal{V}$ . While there have been results providing such an error estimation to the first order expansion Inglot & Majerski (2014); Majerski (2015); Lapinski (2019), we need an error estimate up to the second order expansion, as provided in Theorem 14. For the sake of completeness we also include the statement and proof for the first order expansion in Theorem 13.

Before we state the results, we need to introduce some notation and conventions. For any  $f \in C^0(\mathcal{V})$  and  $h \in C^2(\mathcal{V})$ , respectively we write

$$\zeta_f(z) = f(z) - f(0) \stackrel{z \to 0}{=} o(1)$$

$$\chi_h(z) = h(z) - h(0) - \nabla h(0)^{\top} z - \frac{1}{2} z^{\top} \nabla^2 h(0) z \stackrel{z \to 0}{=} o(||z||^2)$$

for its first and second order remainder terms and abbreviate

$$\zeta_f^\sigma(w)=\zeta_f(\sigma w) \text{ and } \chi_h^\sigma(w)=\frac{1}{\sigma^2}\chi_h(\sigma w) \text{ for } \sigma>0 \,.$$

Let us note the following: If additionally  $f \in C^1(\mathcal{V})$  and  $g \in C^3(\mathcal{V})$ , then  $\zeta_f(z) = O(\|z\|)$  and  $\chi_g(z) = O(\|z\|^3)$  for  $z \to 0$ , i.e. there exist some  $\eta > 0$  and C > 0 such that  $B_{\eta}(0) \subseteq \mathcal{V}$  and

$$|\zeta_f(z)| \le C||z||$$
 and  $|\chi_h(z)| \le C||z||^3$  for all  $z \in B_\eta(0)$ . (24)

In particular for  $\sigma > 0$  and  $\beta : [0, \infty) \to [0, \infty)$  it holds

$$\sup_{\|w\| \le \min\{\eta/\sigma, \beta(\sigma)\}} |\zeta_f^{\sigma}(w)| \le C\sigma\beta(\sigma) \text{ and } \sup_{\|w\| \le \min\{\eta/\sigma, \beta(\sigma)\}} |\chi_g^{\sigma}(w)| \le C\sigma\beta(\sigma)^3. \tag{25}$$

For a non-negative function h with a global minimum at z=0 we define

$$\gamma_h(\delta) = \inf_{\substack{z \in \mathcal{V} \\ \|z\| \ge \delta}} h(z) - h(0).$$

We say that h satisfies a local quadratic growth condition in  $B_{\eta}(0) \subseteq \mathcal{V}$  if there exists a c > 0 with

$$h(z) - h(0) \ge c||z||^2 \text{ for all } z \in B_n(0),$$
 (26)

and that h admits minimum separation outside of  $B_n(0)$  in V if there exists a  $\Delta > 0$  with

$$h(z) - h(0) \ge \Delta > 0 \text{ for all } z \in \mathcal{V} \setminus B_{\eta}(0).$$
 (27)

Clearly, if h satisfies (24), then h also satisfies (26) (for a potentially smaller  $\eta$ ). On the other hand, the global assumption (27) cannot be inferred from properties of h around z=0 alone. Together, (26) and (27) imply that  $\gamma_h$  can be bounded below by

$$\gamma_h(\delta) \ge \min\{c\delta^2, \Delta\},$$
(28)

which is crucial to obtain a quantitative bound on the convergence of the Laplace method. While we state the following results in this full generality, in our application we have  $B_{\eta}(0) = \mathcal{V}$  and thus (27) is not needed, with (24) and (26) holding globally on  $\mathcal{V}$  and implying that

$$\gamma_h(\delta) \ge c\delta^2 \,. \tag{29}$$

**Theorem 13** (First Order Laplace method). Let  $h \in C^3(\mathcal{V})$  with  $\nabla h(0) = 0$  and  $f \in C^1(\mathcal{V})$  for some open  $B_r(0) \subseteq \mathcal{V} \subseteq \mathbb{R}^k$ . Let  $\eta \in (0,r]$  and  $C,c,\Delta>0$  be such that (24), (26) and (27) hold. Then for all  $\sigma \in (0,\overline{\sigma})$  with  $\overline{\sigma} = \min\{\eta^2, (\log(2)/(2C))^2\}$  with it holds

$$\left| e^{\frac{1}{\sigma^2}h(0)} \frac{\sqrt{\det \Sigma}}{Z_{\sigma}} \int_{\mathcal{V}} f(z) e^{-\frac{1}{\sigma^2}h(z)} dz - f(0) \right| \le E_1(\sigma; f, h, r)$$
(30)

with  $\Sigma = \nabla^2 h(0)$  and  $E_1(\sigma; f, h, r) = O(\sigma |\log \sigma|^3)$  for  $\sigma \to 0$  given in (31).

*Proof.* Without loss of generality we can assume h(0)=0. For brevity we also write  $Z^\Sigma_\sigma=\frac{Z_\sigma}{\sqrt{\det\Sigma}}$  for the normalizing constant. Let  $\alpha:[0,\infty)\to[0,\infty)$  be any function and denote  $B(\sigma)=B_{\alpha(\sigma)}(0)$ . We can split

$$\frac{1}{Z^\Sigma_\sigma} \int_{\mathcal{V}} f(z) e^{-\frac{1}{\sigma^2} h(z)} \, \mathrm{d}z = \frac{1}{Z^\Sigma_\sigma} \int_{\mathcal{V} \backslash B(\sigma)} f(z) e^{-\frac{1}{\sigma^2} h(z)} \, \mathrm{d}z + \frac{1}{Z^\Sigma_\sigma} \int_{B(\sigma)} f(z) e^{-\frac{1}{\sigma^2} h(z)} \, \mathrm{d}z \,.$$

We will craft  $\alpha:[0,\infty)\to[0,\infty)$  in such a way that the first integral vanishes and the second converges to f(0) for  $\sigma\to0$ . For the first integral we have

$$\left| \frac{1}{Z_{\sigma}^{\Sigma}} \int_{\mathcal{V} \setminus B(\sigma)} f(z) e^{-\frac{1}{\sigma^2} h(z)} \, \mathrm{d}z \right| \leq \frac{1}{Z_{\sigma}^{\Sigma}} e^{-\frac{1}{\sigma^2} \gamma_h(\alpha(\sigma))} \|f\|_{L^1(\mathcal{V})} \, .$$

For the second integral let us write  $f(z) = f(0) + \zeta_f(z)$  and split

$$\begin{split} \frac{1}{Z_{\sigma}^{\Sigma}} \int_{B(\sigma)} f(z) e^{-\frac{1}{\sigma^2}h(z)} \, \mathrm{d}z &= \frac{f(0)}{Z_{\sigma}^{\Sigma}} \int_{B(\sigma)} e^{-\frac{1}{\sigma^2}h(z)} \, \mathrm{d}z + \frac{1}{Z_{\sigma}^{\Sigma}} \int_{B(\sigma)} \zeta_f(z) e^{-\frac{1}{\sigma^2}h(z)} \, \mathrm{d}z \\ &=: I(\sigma) + J(\sigma) \, . \end{split}$$

First let us estimate the difference between  $I(\sigma)$  and f(0). We have, using the substitution  $z = \sigma w$  and abbreviation  $B_{\mathcal{V}}(\sigma) = (\mathcal{V} \cap B(\sigma))/\sigma$ ,

$$\begin{split} |I(\sigma) - f(0)| \\ & \leq f(0) \left| \frac{1}{Z_1^{\Sigma}} \int_{B_{\mathcal{V}}(\sigma)} e^{-\frac{1}{\sigma^2} h(\sigma w)} \, \mathrm{d}w - 1 \right| \\ & \leq \frac{f(0)}{Z_1^{\Sigma}} \left( \left| \int_{B_{\mathcal{V}}(\sigma)} e^{-\frac{1}{2} w^{\top} \Sigma w} (e^{-\chi_h^{\sigma}(w)} - 1) \, \mathrm{d}w \right| + \left| \int_{\mathbb{R}^k \backslash B_{\mathcal{V}}(\sigma)} e^{-\frac{1}{2} w^{\top} \Sigma w} \, \mathrm{d}w \right| \right) \\ & =: \frac{f(0)}{Z_1^{\Sigma}} (|I_1(\sigma)| + |I_2(\sigma)|) \,. \end{split}$$

Estimating  $I_1(\sigma)$  we obtain

$$|I_1(\sigma)| \le Z_1^{\Sigma} \sup_{w \in B_{\mathcal{V}}(\sigma)} \left| e^{-\chi_h^{\sigma}(w)} - 1 \right|.$$

To estimate  $I_2(\sigma)$ , let us note that  $\Sigma^{1/2}(\mathbb{R}^k \backslash B_R(0)) \subseteq \mathbb{R}^k \backslash B_{\lambda_{\min}(\Sigma)R}(0)$ . Then, via the substitution  $u = \Sigma^{1/2}w$ , we obtain

$$|I_{2}(\sigma)| \leq \left| \int_{\mathbb{R}^{k} \setminus B_{\mathcal{V}}(\sigma)} e^{-\frac{1}{2}w^{\top} \Sigma w} \, \mathrm{d}w \right| \leq \left| \int_{\mathbb{R}^{k} \setminus B_{\min\{r,\alpha(\sigma)\}/\sigma}(0)} e^{-\frac{1}{2}w^{\top} \Sigma w} \, \mathrm{d}w \right|$$

$$\leq Z_{1}^{\Sigma} \frac{1}{(2\pi)^{k/2}} \left| \int_{\Sigma^{1/2}(\mathbb{R}^{k} \setminus B_{\min\{r,\alpha(\sigma)\}/\sigma}(0))} e^{-\frac{1}{2}\|u\|^{2}} \, \mathrm{d}u \right| \leq Z_{1}^{\Sigma} G_{0} \left( \lambda_{\min}(\Sigma) \frac{\min\{r,\alpha(\sigma)\}}{\sigma} \right) ,$$

with the Gaussian tail

$$\mathrm{G}_0(R) = \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n \backslash B_R(0)} e^{-\frac{1}{2} \|u\|^2} \, \mathrm{d} u \,.$$

Next we estimate  $J(\sigma)$  via

$$|J(\sigma)| \leq \frac{1}{Z_1^{\Sigma}} \int_{B_{\mathcal{V}}(\sigma)} \zeta_f(\sigma w) e^{-\frac{1}{2}w^{\top} \Sigma w} e^{-\chi_h^{\sigma}(w)} dw$$

$$\leq \left( \sup_{w \in B_{\mathcal{V}}(\sigma)} |\zeta_f(\sigma w)| \right) \left( \sup_{w \in B_{\mathcal{V}}(\sigma)} |e^{-\chi_h^{\sigma}(w)} - 1| + 1 \right).$$

In all, we obtain the following error estimate

$$\begin{split} & \left| e^{\frac{1}{\sigma^{2}}h(0)} \frac{1}{Z_{\sigma}^{\Sigma}} \int_{\mathcal{V}} f(z) e^{-\frac{1}{\sigma^{2}}h(z)} \, \mathrm{d}z - f(0) \right| \\ & \leq \frac{1}{Z_{\sigma}^{\Sigma}} e^{-\frac{1}{\sigma^{2}}\gamma(\alpha(\sigma))} \|f\|_{L^{1}(\mathcal{V})} + \frac{f(0)}{Z_{1}^{\Sigma}} (|I_{1}(\sigma)| + |I_{2}(\sigma)|) + |J(\sigma)| \\ & \leq \frac{1}{Z_{\sigma}^{\Sigma}} e^{-\frac{1}{\sigma^{2}}\gamma(\alpha(\sigma))} \|f\|_{L^{1}(\mathcal{V})} \\ & + f(0) \left( \sup_{w \in B_{\mathcal{V}}(\sigma)} \left| e^{-\chi_{h}^{\sigma}(w)} - 1 \right| + G_{0} \left( \lambda_{\min}(\Sigma) \frac{\min\{r, \alpha(\sigma)\}}{\sigma} \right) \right) \\ & + \left( \sup_{w \in B_{\mathcal{V}}(\sigma)} |\zeta_{f}(\sigma w)| \right) \left( \sup_{w \in B_{\mathcal{V}}(\sigma)} |e^{-\chi_{h}^{\sigma}(w)} - 1| + 1 \right) . \end{split}$$

Thus we observe that we need to select  $\alpha$  in such a way that for  $\sigma \to 0$ 

(i) 
$$(Z_{\sigma}^{\Sigma})^{-1} \exp(-\gamma_h(\alpha(\sigma))/\sigma^2) \to 0$$

(ii) 
$$G_0\left(\lambda_{\min}(\Sigma)\frac{\min\{r,\alpha(\sigma)\}}{\sigma}\right) \to 0$$

(iii) 
$$\sup_{w \in B_{\mathcal{V}}(\sigma)} |\chi_h^{\sigma}(w)| \to 0$$

(iv) 
$$\sup_{w \in B_{\mathcal{V}}(\sigma)} |\zeta_f(\sigma w)| \to 0$$

Let us make the Ansatz  $\alpha(\sigma) = \sigma\beta(\sigma)$  for some  $\beta: (0,\infty) \to (0,\infty)$  with  $\beta(\sigma) \to \infty$  for  $\sigma \to 0$ . Taking  $\beta(\sigma) = |\log(\sigma)|$ , noting that  $|e^{-x} - 1| \le 2|x|$  for  $|x| \le \log 2$  as well as using (28), (25) and the fact that  $\sigma|\log(\sigma)|$ ,  $\sigma|\log(\sigma)|^3 \le \sigma^{1/2}$  when  $\sigma \in (0,1)$ , yields the following estimates

$$\begin{split} \exp(-\frac{1}{\sigma^2}\gamma(\alpha(\sigma))) &\leq \exp(-\min\{c\log(\sigma)^2, \Delta\sigma^{-2}\})\,, \\ G_0\left(\lambda_{\min}(\Sigma)\frac{\min\{r, \alpha(\sigma)\}}{\sigma}\right) &\leq G_0\left(\lambda_{\min}(\Sigma)\min\{r\sigma^{-1}, |\log(\sigma)|\}\right)\,, \\ \sup_{w \in B_{\mathcal{V}}(\sigma)} \left|e^{-\chi_h^{\sigma}(w)} - 1\right| &\leq 2\sup_{w \in B_{\mathcal{V}}(\sigma)} |\chi_h^{\sigma}(w)| \leq 2C\sigma\beta(\sigma)^3 = 2C\sigma|\log(\sigma)|^3\,, \\ \sup_{w \in B_{\mathcal{V}}(\sigma)} |\zeta_f(\sigma w)| &\leq C\sigma\beta(\sigma) = C\sigma|\log(\sigma)|\,, \end{split}$$

when  $0 < \sigma \le \overline{\sigma} := \min\{\eta^2, (\log 2/(2C))^2\}$ . Plugging all these estimates back yields (30) with

$$E(\sigma) = \frac{1}{Z_{\sigma}^{\Sigma}} \exp(-\min\{c \log(\sigma)^{2}, \Delta \sigma^{-2}\}) \|f\|_{L^{1}(\mathcal{V})}$$

$$+ f(0) \left(2C\sigma |\log(\sigma)|^{3} + G_{0}(\lambda_{\min}(\Sigma) \min\{r\sigma^{-1}, |\log(\sigma)|\})\right)$$

$$+ C\sigma |\log(\sigma)|(1 + 2C\sigma |\log(\sigma)|^{3})$$

$$\stackrel{\sigma \to 0}{=} O(\sigma \log(\sigma)^{3}) = \widetilde{O}(\sigma).$$
(31)

The next result provides an nonasymptotic estimate of the Laplace method for the second order expansion. We will need a version that allows for a non-zero gradient of f, as long as it is contained in the subspace  $\{(q,-q)\mid q\in\mathbb{R}^n\}\subseteq\mathbb{R}^n\times\mathbb{R}^n$ .

**Theorem 14** (Second Order). Let  $f, h \in C^3(\mathcal{V} \times \mathcal{V})$  be such that f(0) = 0 and  $\nabla h(0) = 0$  for some open  $B_r(0) \subseteq \mathcal{V} \subseteq \mathbb{R}^n$  for some r > 0. Further, suppose that h is additively separable as  $h(z) = \overline{h}(z_1) + \overline{h}(z_2)$  and  $\nabla f(0) = \binom{q}{-q}$  for some  $q \in \mathbb{R}^n$ . Additionally, let  $\eta \in (0, r]$  and  $C, c, \Delta > 0$  be such that (24), (26) and (27) hold, where this time (24) applies also for h = f. Then for  $\sigma \in (0, \overline{\sigma})$  with  $\overline{\sigma} = \min\{\eta^2, (\log(2)/(2C))^2\}$  it holds

$$\left| e^{\frac{1}{\sigma^2}h(0)} \frac{\sqrt{\det \Sigma}}{Z_{\sigma}} \int_{\mathcal{V} \times \mathcal{V}} f(z) e^{-\frac{1}{\sigma^2}h(z)} dz - \sigma^2 A(f, h) \right| \le E_2(\sigma; f, h, r), \tag{32}$$

where  $\Sigma = \nabla^2 h(0)$  with  $E_2(\sigma; f, h, r) = O(\sigma^3 |\log(\sigma)|^3)$  for  $\sigma \to 0$ . Here the second order coefficient  $A_2(f, h)$  is given by

$$A_2(f,h) = \frac{1}{2} \operatorname{tr}(\Sigma^{-1} \nabla^2 f(0)).$$
 (33)

and  $E_2(\sigma; f, h, r)$  is given in (34).

*Proof.* The proof is similar to the proof of Theorem 13 with the modification for the product space  $\mathbb{R}^k = \mathbb{R}^n \times \mathbb{R}^n$ . Again we can assume h(0) = 0. Let  $\alpha : [0, \infty) \to [0, \infty)$  be any function and denote  $\widehat{B}(\sigma) = B_{\alpha(\sigma)}(0) \times B_{\alpha(\sigma)}(0) \subseteq \mathbb{R}^n \times \mathbb{R}^n$  and  $\widehat{\mathcal{V}} = \mathcal{V} \times \mathcal{V}$ . Write

$$\begin{split} &\frac{1}{\sigma^2} \frac{1}{Z_{\sigma}^{\Sigma}} \int_{\widehat{\mathcal{V}}} f(z) e^{-\frac{1}{\sigma^2} h(z)} \, \mathrm{d}z \\ &= \frac{1}{\sigma^2} \frac{1}{Z_{\sigma}^{\Sigma}} \int_{\widehat{\mathcal{V}} \setminus \widehat{B}(\sigma)} f(z) e^{-\frac{1}{\sigma^2} h(z)} \, \mathrm{d}z + \frac{1}{\sigma^2} \frac{1}{Z_{\sigma}^{\Sigma}} \int_{\widehat{\mathcal{V}} \cap \widehat{B}(\sigma)} f(z) e^{-\frac{1}{\sigma^2} h(z)} \, \mathrm{d}z \,. \end{split}$$

We will carefully craft  $\alpha$  such that  $\lim_{\sigma\to 0}\alpha(\sigma)=0$  and the first and second integral converge to 0 and  $\frac{1}{2}\operatorname{tr}(\Sigma^{-1}\nabla^2 f(0))$ , respectively, with explicit error bounds. For the first one we have

$$\left|\frac{1}{\sigma^2}\frac{1}{Z_\sigma^\Sigma}\int_{\widehat{\mathcal{V}}\backslash\widehat{B}(\sigma)}f(z)e^{-\frac{1}{\sigma^2}h(z)}\,\mathrm{d}z\right|\leq \frac{e^{-\frac{1}{\sigma^2}\gamma_h(\alpha(\sigma))}}{Z_\sigma^\Sigma\sigma^2}\|f\|_{L^1(\mathcal{V}\times\mathcal{V})}\,.$$

Now consider the second integral. By a substitution  $z = \sigma w$  and abbreviating  $\widetilde{B}_{\widehat{\mathcal{V}}}(\sigma) = (\widehat{\mathcal{V}} \cap \widehat{B}(\sigma))/\sigma = \widetilde{B}_{\mathcal{V}}(\sigma) \times \widetilde{B}_{\mathcal{V}}(\sigma)$  with  $\widetilde{B}_{\mathcal{V}}(\sigma) = (\mathcal{V} \cap B_{\alpha(\sigma)}(0))/\sigma$  as well as  $\chi_g^{\sigma}(w) = \frac{1}{\sigma^2}\chi_g(\sigma w)$  for  $g \in \{f, h\}$ , we have

$$\begin{split} &\frac{1}{\sigma^2} \frac{1}{Z_{\sigma}^{\Sigma}} \int_{\widehat{\mathcal{V}} \cap \widehat{B}(\sigma)} f(z) e^{-\frac{1}{\sigma^2} h(z)} \, \mathrm{d}z \\ &= \frac{1}{Z_1^{\Sigma}} \int_{\widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} \left( \frac{1}{\sigma} \nabla f(0)^\top w + \frac{1}{2} w^\top \nabla^2 f(0) w + \chi_f^{\sigma}(w) \right) e^{-\frac{1}{\sigma^2} h(\sigma w)} \, \mathrm{d}w \\ &=: H(\sigma) + I(\sigma) + J(\sigma) \, . \end{split}$$

Note that  $H(\sigma)$  vanishes due to separability of h and the particular form of  $\nabla f(0)$ :

$$\begin{split} \sigma Z_1^{\Sigma} H(\sigma) &= \int_{\widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} \nabla f(0)^{\top} w \cdot e^{-\frac{1}{\sigma^2} h(\sigma w)} \, \mathrm{d}w \\ &= \int_{\widetilde{B}_{\mathcal{V}}(\sigma)} \int_{\widetilde{B}_{\mathcal{V}}(\sigma)} (q^{\top} w_1 - q^{\top} w_2) e^{-\frac{1}{\sigma^2} \overline{h}(\sigma w_1)} e^{-\frac{1}{\sigma^2} \overline{h}(\sigma w_2)} \, \mathrm{d}w_1 \, \mathrm{d}w_2 \\ &= 0 \,. \end{split}$$

Next, let us consider the difference between  $I(\sigma)$  and  $\frac{1}{2}\operatorname{tr}(\Sigma^{-1}\nabla^2 f(0))$  given by

$$\left| \frac{1}{Z_1^{\Sigma}} \int_{\widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} \frac{1}{2} w^{\top} \nabla^2 f(0) w \cdot e^{-\frac{1}{2} w^{\top} \Sigma w} e^{-\chi_h^{\sigma}(w)} dw - \frac{1}{2} \operatorname{tr}(\Sigma^{-1} \nabla^2 f(0)) \right| \\
\leq \left| \frac{1}{Z_1^{\Sigma}} \int_{\widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} \frac{1}{2} w^{\top} \nabla^2 f(0) w \cdot e^{-\frac{1}{2} w^{\top} \Sigma w} (e^{-\chi_h^{\sigma}(w)} - 1) dw \right| \\
+ \left| \frac{1}{Z_1^{\Sigma}} \int_{\mathbb{R}^{2n} \setminus \widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} \frac{1}{2} w^{\top} \nabla^2 f(0) w \cdot e^{-\frac{1}{2} w^{\top} \Sigma w} dw \right| \\
=: |I_1(\sigma)| + |I_2(\sigma)|.$$

Then

$$\begin{split} |I_1(\sigma)| & \leq \left(\sup_{w \in \widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} |e^{-\chi_h^{\sigma}(w)} - 1|\right) \frac{1}{Z_1^{\Sigma}} \int_{\mathbb{R}^{2n}} \frac{1}{2} \left|w^{\top} \nabla^2 f(0) w\right| \cdot e^{-\frac{1}{2} w^{\top} \Sigma w} \, \mathrm{d}w \\ & \leq \frac{n}{2} \|\Sigma^{-1/2} \nabla^2 f(0) \Sigma^{-1/2} \| \left(\sup_{w \in \widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} |e^{-\chi_h^{\sigma}(w)} - 1|\right) \end{split}$$

Furthermore, by a variable substitution  $u = \Sigma^{1/2}w$  we obtain

$$|I_{2}(\sigma)| \leq \left| \frac{1}{Z_{1}} \int_{\mathbb{R}^{2n} \setminus \Sigma^{1/2} \widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} \frac{1}{2} u^{\top} \Sigma^{-1/2} \nabla^{2} f(0) \Sigma^{-1/2} u \cdot e^{-\frac{1}{2} \|u\|^{2}} du \right|$$

$$\leq \frac{1}{2} L(\sigma) \|\Sigma^{-1/2} \nabla^{2} f(0) \Sigma^{-1/2} \|,$$

where

$$L(\sigma) = \frac{1}{Z_1} \int_{\mathbb{R}^k \setminus \overline{\Sigma}^{1/2} \widetilde{B}_{\mathcal{V}}(\sigma)} \|u\|^2 e^{-\frac{1}{2} \|u\|^2} du.$$

Noting that for  $\overline{\Sigma} = \nabla^2 h(0)$  it holds  $\Sigma^{1/2} \widetilde{B}_{\widehat{\mathcal{V}}}(\sigma) = \overline{\Sigma}^{1/2} \widetilde{B}_{\mathcal{V}}(\sigma) \times \overline{\Sigma}^{1/2} \widetilde{B}_{\mathcal{V}}(\sigma)$  and

$$\overline{\Sigma}^{1/2}\widetilde{B}_{\mathcal{V}}(\sigma) = \overline{\Sigma}^{1/2}(\mathcal{V} \cap B_{\alpha(\sigma)}(0))/\sigma \supseteq B_{\lambda_{\min}(\Sigma)\min\{r,\alpha(\sigma)\}/\sigma}(0)$$

and hence

$$L(\sigma) \le 2 G_2 \left( \lambda_{\min}(\Sigma) \frac{\min\{r, \alpha(\sigma)\}}{\sigma} \right)$$
,

with the Gaussian second moment tail

$$G_2(R) := \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n \setminus B_R(0)} ||u||^2 e^{-\frac{1}{2}||u||^2} du.$$

Now consider the final expression

$$\begin{split} |J(\sigma)| &= \left| \frac{1}{Z_1^{\Sigma}} \int_{\widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} \chi_f^{\sigma}(w) e^{-h^{\sigma}(w)} \, \mathrm{d}w \right| \\ &\leq \left( \sup_{w \in \widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} |\chi_f^{\sigma}(w)| \right) \frac{1}{Z_1^{\Sigma}} \int_{\widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} e^{-h^{\sigma}(w)} \, \mathrm{d}w \end{split}$$

Let us further estimate

$$\begin{split} \frac{1}{Z_1^{\Sigma}} \int_{\widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} e^{-h^{\sigma}(w)} \, \mathrm{d}w &= \frac{1}{Z_1^{\Sigma}} \int_{\widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} e^{-\frac{1}{2}w^{\top} \Sigma w} e^{-\chi_h^{\sigma}(w)} \, \mathrm{d}w \\ &\leq \left( \sup_{w \in \widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} |e^{-\chi_h^{\sigma}(w)} - 1| \right) + 1 \, . \end{split}$$

Hence we see that we need to pick  $\alpha$  in such a way that for  $\sigma \to 0$ 

(i) 
$$(Z_{\sigma}^{\Sigma})^{-1}\sigma^{-2}\exp(-\gamma_h(\alpha(\sigma))/\sigma^2) \to 0.$$

(ii) 
$$G_2\left(\lambda_{\min}(\Sigma)\frac{\min\{r,\alpha(\sigma)\}}{\sigma}\right) \to 0$$

(iii) 
$$\sup_{w \in \widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} |\chi_h^{\sigma}(w)| \to 0$$

(iv) 
$$\sup_{w \in \widetilde{B}_{\widehat{\mathcal{V}}}(\sigma)} |\chi_f^{\sigma}(w)| \to 0$$

Again take  $\alpha(\sigma) = \sigma\beta(\sigma)$  for  $\beta(\sigma) = |\log(\sigma)|$ . This yields

$$\exp(-\frac{1}{\sigma^{2}}\gamma_{h}(\alpha(\sigma))) \leq \exp(-\min\{c\log(\sigma)^{2}, \Delta\sigma^{-2}\}),$$

$$G_{2}\left(\lambda_{\min}(\Sigma)\frac{\min\{r, \alpha(\sigma)\}}{\sigma}\right) \leq G_{2}\left(\lambda_{\min}(\Sigma)\min\{r\sigma^{-1}, |\log(\sigma)|\}\right),$$

$$\sup_{w \in B_{\mathcal{V}}(\sigma)} \left|e^{-\chi_{h}^{\sigma}(w)} - 1\right| \leq 2\sup_{w \in B_{\mathcal{V}}(\sigma)} |\chi_{h}^{\sigma}(w)| \leq 2C\sigma\beta(\sigma)^{3} = 2C|\log(\sigma)|^{3},$$

$$\sup_{w \in B_{\mathcal{V}}(\sigma)} \left|\chi_{f}^{\sigma}(w)\right| \leq C\sigma\beta(\sigma)^{3} = C|\log(\sigma)|^{3},$$

when  $0 < \sigma \le \overline{\sigma} := \min\{\eta^2, (\log(2)/(2C))^2\}$ . Combining all these estimates one obtains, as in the proof of Theorem 13, that (32) holds with

$$\frac{1}{\sigma^2} E(\sigma) = \frac{\sqrt{\det \Sigma}}{Z_{\sigma} \sigma^2} \exp(-\min\{c \log(\sigma)^2, \Delta \sigma^{-2}\}) \|f\|_{L^1(\mathcal{V} \times \mathcal{V})} 
+ \|\Sigma^{-1/2} \nabla^2 f(0) \Sigma^{-1/2} \| \left( nC\sigma |\log(\sigma)|^3 + G_2 \left( \lambda_{\min}(\Sigma) \min\{r\sigma^{-1}, |\log(\sigma)|^3 \right) \right) 
+ C\sigma |\log(\sigma)|^3 (1 + 2C\sigma |\log(\sigma)|^3) 
\stackrel{\sigma \to 0}{=} O(\sigma |\log(\sigma)|^3) = \widetilde{O}(\sigma).$$
(34)

#### D.2 LOCAL VERSION OF THEOREM 1

 First we prove a local version of Theorem 1.

**Theorem 15.** Assume that there exists an subset  $\mathcal{U} \subseteq \mathcal{M}$  that is  $C^3$ -diffeomorphic to an open subset of  $\mathcal{V} \subseteq \mathbb{R}^k$  with  $B_r(0) \subseteq \mathcal{V}$  via  $\psi : \mathcal{V} \to \mathcal{U}$  for some  $0 \le k \le d$  and that  $\mu \ll \operatorname{Vol}_{\mathcal{U}}$ , where  $\operatorname{Vol}_{\mathcal{U}}$  is the volume measure on  $\mathcal{U}$ . Moreover, assume that  $\mu(\cdot) = \frac{\mathrm{d}\,\mu}{\mathrm{d}\,\mathrm{Vol}_{\mathcal{U}}} \in C^3(\mathcal{U})$ . Let  $x \in \mathbb{R}^d$  be any point with  $\psi(0) =: p = \arg\min_{p \in \mathcal{M}} \|x-p\|$  and  $\delta = \inf_{y \in \mathcal{M} \setminus \mathcal{U}} \|x-y\|^2 - \|x-p\|^2 > 0$ , and such that

$$\Sigma = \psi'(0)^{\top} \psi'(0) + \sum_{i=1}^{d} (p - x)_i \nabla^2 \psi_i(0) > 0.$$
 (35)

Then there exists some  $\overline{\sigma} > 0$  such that for all  $\sigma \in (0, \overline{\sigma})$ 

(i) it holds that

$$\|\mathbb{E}\nu_{x,\sigma} - p\| \le \frac{2(E_1(\sigma; f_0, h, r)\|p - x\| + \sqrt{d}E_1(\sigma; f_1, h, r) + (\|p - x\| + \|\mu\|_1)\Upsilon(\sigma; \Sigma))}{\mu(p)}$$

(ii) it holds for  $P_0 = \psi'(0)\Sigma^{-1}\psi'(0)^{\top}$  that

$$\begin{split} \left\| \frac{1}{\sigma^2} \operatorname{Cov}(\nu_{x,\sigma}) - P_0 \right\| \\ &\leq \frac{4d}{\mu(p)^2} \frac{E_2(\sigma; \overline{f}, \overline{h}, r)}{\sigma^2} \\ &\quad + \left( \frac{4(\|\mu\|_2 + \|\mu\|_1^2)}{\mu(p)^2} \frac{\Upsilon(\sigma; \Sigma)^2}{\sigma^2} + 12 \frac{E_1(\sigma; f_0, h, r) + \Upsilon(\sigma; \Sigma)}{\mu(p)} \right) \|P_0\| \end{split}$$

with h,  $f_0$  and  $f_1$  given in (36),  $\overline{f}_2$  and  $\overline{h}$  in (37),  $\|\mu\|_i = \mathbb{E}_{y \sim \mu} \|y - x\|^i$ ,  $\Upsilon(\sigma; \Sigma) = \frac{\sqrt{\det \Sigma}}{Z_{\sigma}} \exp(-\frac{1}{2\sigma^2}\delta)$ ,

$$E_1(\sigma; f_1, h, r) = \max_{l=1,\dots,d} E_1(\sigma; f_1^l, h, r),$$

$$E_2(\sigma; \overline{f}_2, \overline{h}, r) = \max_{l,j=1,\dots,d} E_2(\sigma; \overline{f}_2^{jl}, \overline{h}, r),$$

with  $E_1(\sigma; f_1^l, h, r)$  given in (31) and  $E_2(\sigma; \overline{f}_2^{jl}, \overline{h}, r)$  in (34), respectively.

**Remark 16.** This theorem only requires  $\mathcal{M}$  to be locally a manifold, namely at  $\mathcal{U}$ .

*Proof.* We denote by  $\lambda = \psi^{-1} \# \mu$  be corresponding pullback measure on  $\mathcal V$  and denote its positive density again by  $\lambda \in C^3(\mathcal V)$ . Note that  $\lambda(0) = \mu(p)$  by (20). Moreover, let v = x - p and  $P_0 = \psi'(0) \Sigma^{-1} \psi'(0)^{\top}$  for the rest of the proof. First we need to investigate the non-asymptotic convergence of following two integrals:

$$\begin{split} p_{\sigma}(x) &= \int_{\mathcal{M}} \varphi_{\sigma}(x-y) \, \mathrm{d}\mu(y) \\ &= \underbrace{\int_{\mathcal{V}} \lambda(z) \frac{1}{Z_{\sigma}} e^{-\frac{1}{2\sigma^2} \|x-\psi(z)\|^2} \, \mathrm{d}z}_{S_0} + \underbrace{\int_{\mathcal{M} \backslash \mathcal{U}} \varphi_{\sigma}(x-y) \, \mathrm{d}\mu(y)}_{R_2} \end{split}$$

and

$$p_{\sigma}(x)(\mathbb{E}\nu_{x,\sigma} - x) = \int_{\mathcal{M}} (y - x)\varphi_{\sigma}(x - y) \,\mathrm{d}\mu(y)$$

$$= \underbrace{\int_{\mathcal{V}} \lambda(z)(\psi(z) - x) \frac{1}{Z_{\sigma}} e^{-\frac{1}{2\sigma^{2}} \|x - \psi(z)\|^{2}} \,\mathrm{d}z}_{S_{1}} + \underbrace{\int_{\mathcal{M} \setminus \mathcal{U}} (y - x)\varphi_{\sigma}(x - y) \,\mathrm{d}\mu(y)}_{R_{2}}$$

We can write each  $S_i$  and  $R_i$  as

$$S_i = \frac{1}{Z_{\sigma}} \int_{\mathcal{V}} f_i(z) e^{-\frac{1}{\sigma^2} h(z)} \, \mathrm{d}z \,, \quad R_i = \int_{\mathcal{M} \setminus \mathcal{U}} g_i(y) \varphi_{\sigma}(x-y) \, \mathrm{d}\mu(y) \,,$$

where

$$h(z) = \frac{1}{2} ||x - \psi(z)||^2, \quad f_0(z) = \lambda(z), \quad f_1(z) = \lambda(z)(\psi(z) - x),$$
 (36)

as well as  $g_0(y) = 1$  and  $g_1(y) = y - x$ . Then  $f_0, f_1 \in C^3(\mathcal{V})$  and

$$h(0) = \frac{1}{2} ||v||^2, \quad \nabla h(0) = 0, \quad \nabla^2 h(0) = \Sigma,$$

where the second equality is due to  $p=\psi(0)$  being the closest point from  $\mathcal U$  to x. Also obviously  $f_0(0)=\lambda(0)$  and  $f_1(0)=\lambda(0)(p-x)$ . Applying Theorem 13 component-wise, we obtain for  $\sigma\in(0,\overline\sigma_i)$ , with  $\overline\sigma_i=\overline\sigma_i(\eta,C)$  given as in Theorem 13 with  $\eta,C>0$  dependent on  $f_i$  and h, the estimate

$$\underbrace{\left\| e^{\frac{1}{2\sigma^2} \|v\|^2} \sqrt{\det \Sigma} S_i - f_i(0) \right\|}_{=:\|F_i\|} \le D_i(\sigma),$$

with

$$D_0(\sigma) = E_1(\sigma; f_0, h, r), \quad D_1(\sigma) = \sqrt{d}E_1(\sigma; f_1, h, r) := \sqrt{d} \max_{l=1}^{max} E_1(\sigma; f_1^l, h, r),$$

and where  $E_1$  is given in Theorem 13 and  $f_1^l$  is the l-th component of  $f_1$ . Furthermore, we can estimate for i = 0, 1

$$\begin{aligned} \|R_i\| &= \left\| \int_{\mathcal{M} \setminus \mathcal{U}} g_i(y) \varphi_{\sigma}(x - y) \, \mathrm{d}\mu(y) \right\| \\ &= \frac{1}{Z_{\sigma}} \left\| \int_{\mathcal{M} \setminus \mathcal{U}} g_i(y) e^{\frac{1}{2\sigma^2} \|x - y\|^2} \, \mathrm{d}\mu(y) \right\| \\ &\leq e^{-\frac{1}{2\sigma^2} \|v\|^2} \frac{e^{-\frac{1}{2\sigma^2} \delta}}{Z_{\sigma}} \int_{\mathcal{M} \setminus \mathcal{U}} \|g_i(y)\| \, \mathrm{d}\mu(y) \,, \end{aligned}$$

where we have used that  $||x-y||^2 \ge \delta + ||v||^2$  for all  $y \in \mathcal{M} \setminus \mathcal{U}$ . Thus we get the estimate

$$\left\| e^{\frac{1}{2\sigma^2} \|v\|^2} R_i \right\| \le \frac{e^{-\frac{1}{2\sigma^2} \delta}}{Z_{\sigma}} \|\mu\|_i,$$

where the quantity  $\|\mu\|_i = \int_{\mathcal{M}} \|y - x\|^i \, \mathrm{d}\mu(y)$  denotes the *i*-th centered moment of  $\mu$ . Let us denote  $J = e^{\frac{1}{2\sigma^2}\|v\|^2} \sqrt{\det \Sigma}$  and express

$$S_i = f_i(0)/J + F_i/J$$

Now we can estimate the distance to p by

$$\mathbb{E}\nu_{x,\sigma} - p = \frac{S_1 + R_1}{S_0 + R_0} - (p - x) = \frac{f_1(0)/J + F_1/J + R_1}{f_0(0)/J + F_0/J + R_0} - (p - x)$$
$$= \frac{f_1(0) + F_1 + JR_1}{f_0(0) + F_0 + JR_0} - (p - x) = \frac{(p - x) + T_1}{1 + T_0} - (p - x) = \frac{T_1 - T_0(p - x)}{1 + T_0}.$$

where  $T_i = (F_i + JR_i)/\lambda(0)$ . Let us bound

$$||T_i|| \le \frac{1}{\lambda(0)} \left( D_i(\sigma) + \sqrt{\det \Sigma} \frac{e^{-\frac{1}{2\sigma^2}\delta}}{Z_{\sigma}} \right) ||\mu||_i \stackrel{\sigma \to 0}{\longrightarrow} 0.$$

Then  $|T_0| < 1/2$  if  $\sigma < \overline{\sigma}(f_0, h, r) := \min\{\overline{\sigma}_0, \widetilde{\sigma}\}$  with  $\widetilde{\sigma}$  depending on  $D_0$ ,  $\Sigma$  and  $\delta$ . For this  $\sigma \in (0, \overline{\sigma})$  we obtain

1353  $\|\mathbb{E}\nu_{x,\sigma} - p\|$ 

1354 
$$\leq 2|T_0|||p-x|| + 2||T_1||$$

$$\leq \frac{2}{\lambda(0)} \left( E_1(\sigma; f_0, h, r) \|p - x\| + \sqrt{d} E_1(\sigma; f_1, h, r) + (\|p - x\| + \|\mu\|_1) \sqrt{\det \Sigma} \frac{e^{-\frac{1}{2\sigma^2}\delta}}{Z_{\sigma}} \right).$$

This proves (i). The proof of (ii) is analogous. For  $(z,\widetilde{z})\in\mathcal{V}\times\mathcal{V}$  define

$$\overline{f}_{2}(z,\widetilde{z}) = \lambda(z)\lambda(\widetilde{z})((\psi(z) - x)(\psi(z) - x)^{\top} - (\psi(z) - x)(\psi(\widetilde{z}) - x)^{\top}),$$

$$\overline{h}(z,\widetilde{z}) = \frac{1}{2}\|x - \psi(z)\|^{2} + \frac{1}{2}\|x - \psi(\widetilde{z})\|^{2}.$$
(37)

Then  $f, h \in C^3(\mathcal{V} \times \mathcal{V}), \overline{f}_2(0,0) = 0$  and

$$\overline{h}(0,0) = ||v||^2$$
,  $\nabla \overline{h}(0,0) = 0$ ,  $\nabla^2 \overline{h}(0,0) = \overline{\Sigma} := \operatorname{diag}(\Sigma, \Sigma)$ .

Let  $\overline{f}_2^{jl}$  be the (j,l)-th entry of  $\overline{f}_2$ . Then it is an elementary, but very tedious exercise to show that  $\overline{f}$  and  $\overline{h}$  satisfy the conditions of Theorem 14 and that moreover the second order coefficient is given by

$$A_2(\overline{f}_2^{jl}, \overline{h}) = \lambda(0)^2 (\psi'(0) \Sigma^{-1} \psi'(0)^{\top})_{lj} = \lambda(0)^2 (P_0)_{lj},$$

i.e. 
$$\overline{A}_2 := (A_2(\overline{f}_2^{jl}, \overline{h}))_{l,j=1}^d = \lambda(0)^2 P_0$$
. Now split

$$\begin{aligned} p_{\sigma}(x)^2 \operatorname{Cov}(\nu_{x,\sigma}) &= p_{\sigma}(x)^2 (\operatorname{Cov}(\nu_{x,\sigma}) - \mathbb{E}\nu_{x,\sigma} (\mathbb{E}\nu_{x,\sigma})^{\top}) \\ &= \int_{\mathcal{M} \times \mathcal{M}} ((y-x)(y-x)^{\top} - (y-x)(\widetilde{y}-x)^{\top}) \varphi_{\sigma}(x-y) \varphi_{\sigma}(x-\widetilde{y}) \operatorname{d}(\mu \otimes \mu)(y,\widetilde{y}) \\ &= \underbrace{\int_{\mathcal{V} \times \mathcal{V}} \overline{f}_2(z,\widetilde{z}) \frac{1}{Z_{\sigma}^2} e^{-\frac{1}{2\sigma^2} (\|x-\psi(z)\|^2 + \|x-\psi(\widetilde{z})\|^2)} \operatorname{d}(z,\widetilde{z})}_{\overline{S}_2} \end{aligned}$$

$$+\underbrace{\int_{\mathcal{M}\times\mathcal{M}\setminus\mathcal{U}\times\mathcal{U}}((y-x)(y-x)^{\top}-(y-x)(\widetilde{y}-x)^{\top})\varphi_{\sigma}(x-y)\varphi_{\sigma}(x-\widetilde{y})\,\mathrm{d}(\mu\otimes\mu)(y,\widetilde{y})}_{\overline{R}_{2}}.$$

Applying Theorem 14 component-wise to  $\overline{f}_2$  we obtain for  $\sigma \in (0, \overline{\sigma}_2)$ , with  $\overline{\sigma}_2 = \overline{\sigma}_2(\eta, C)$  and  $\eta, C > 0$  dependent on  $\overline{f}_2$  and  $\overline{h}$ , that

$$\underbrace{\left\| (e^{\frac{1}{2\sigma^2} \|v\|^2} \sqrt{\det \Sigma})^2 \overline{S}_2 - \sigma^2 \overline{A}_2 \right\|}_{\|\overline{G}_2\|} \leq dE_2(\sigma; \overline{f}_2, \overline{h}, r) := d \max_{l, j = 1, \dots, d} E_2(\sigma; \overline{f}_2^{jl}, \overline{h}, r),$$

with  $E_2(\sigma; \overline{f}_2^{jl}, \overline{h}, r)$  given in Theorem 34. The term  $\overline{R}_2$  on the other hand can be estimated again by

$$\|\overline{R}_{2}\| \leq e^{-\frac{1}{\sigma^{2}}\|v\|^{2}} \frac{e^{-\frac{1}{\sigma^{2}}\delta}}{Z_{\sigma}^{2}} \int_{\mathcal{M}\times\mathcal{M}\setminus\mathcal{U}\times\mathcal{U}} \|(y-x)(y-x)^{\top} - (y-x)(\widetilde{y}-x)^{\top}\| d(\mu\otimes\mu)(y,\widetilde{y})$$

$$\leq e^{-\frac{1}{\sigma^{2}}\|v\|^{2}} \frac{e^{-\frac{1}{\sigma^{2}}\delta}}{Z_{\sigma}^{2}} (\|\mu\|_{2} + \|\mu\|_{1}^{2}).$$

For  $\overline{S}_2$  it holds

$$\overline{S}_2 = \sigma^2 \lambda(0)^2 P_0(x) / J^2 + \overline{G}_2 / J^2$$

and thus

$$\frac{1}{\sigma^2} \operatorname{Cov}(\nu_{x,\sigma}) - P_0 = \frac{\overline{S}_2/\sigma^2 + \overline{R}_2/\sigma^2}{(S_0 + R_0)^2} - P_0$$

$$= \frac{\lambda(0)^2 P_0/J^2 + \overline{G}_2/(\sigma^2 J^2) + \overline{R}_2/\sigma^2}{(\lambda(0)/J + F_0/J + R_0)^2} - P_0$$

$$= \frac{P_0(x) + \overline{T}_2}{(1 + T_0)^2} - P_0$$

$$= \frac{\overline{T}_2 - T_0(2 + T_0)P_0}{(1 + T_0)^2}$$

with  $\overline{T}_2 = (\overline{G}_2 + J^2 \overline{R}_2)/(\lambda(0)^2 \sigma^2)$ . Again we have the estimate

$$\|\overline{T}_2\| \le \frac{d}{\lambda(0)^2} \frac{E_2(\sigma)}{\sigma^2} + \frac{(\det \Sigma)(\|\mu\|_2 + \|\mu\|_1^2)}{\lambda(0)^2} \frac{e^{-\frac{1}{\sigma^2}\delta}}{Z_\sigma^2 \sigma^2}$$

and for  $\sigma \in (0, \overline{\sigma})$  (with  $\overline{\sigma}$  as before, guaranteeing  $|T_0| < 1/2$ )

$$\begin{split} & \left\| \frac{1}{\sigma^{2}} \operatorname{Cov}(\nu_{x,\sigma}) - P_{0} \right\| \\ & \leq 4 \|\overline{T}_{2}\| + 12 |T_{0}| \|P_{0}\| \\ & \leq \frac{4d}{\lambda(0)^{2}} \frac{E_{2}(\sigma)}{\sigma^{2}} + \left( \frac{4(\det \Sigma)(\|\mu\|_{2} + \|\mu\|_{1}^{2})}{\lambda(0)^{2}} \frac{e^{-\frac{1}{\sigma^{2}}\delta}}{Z_{\sigma}^{2}\sigma^{2}} + 12 \frac{D_{0}(\sigma) + \sqrt{\det \Sigma} \frac{e^{-\frac{1}{2\sigma^{2}}\delta}}{Z_{\sigma}}}{\lambda(0)} \right) \|P_{0}\| \end{split}$$

#### D.3 BOUNDS IN TERMS OF DISTRIBUTION AND MANIFOLD PARAMETERS

In this section we bound the constants appearing in Theorem 15 in terms of parameters of the distribution  $\mu$  and the manifold  $\mathcal{M} \in \mathcal{M}_k(\tau, M)$ , when the chart is given by the graph chart  $\psi = \psi_p : \mathcal{V} \to \mathcal{M}$  with  $\mathcal{V} := B_{\min\{\tau_{\mathcal{M}}, M\}/4}^{\mathrm{T}_p \mathcal{M}}(0), x \in \mathcal{T}(\tau), \tau \in (0, \tau_{\mathcal{M}})$  and  $p = \pi(x)$ . According to the proof of Theorem 15 need to consider the following maps

$$h: \mathcal{V} \to \mathbb{R}: z \mapsto \frac{1}{2} \|x - \psi_p(z)\|^2, \tag{38}$$

as well as

$$f_0: \mathcal{V} \to \mathbb{R}: z \mapsto \lambda(z)$$
, (39)

$$f_1: \mathcal{V} \to \mathbb{R}^d: z \mapsto \lambda(z)(\psi_n(z) - x)$$
, (40)

$$\overline{f}_2: \mathcal{V} \times \mathcal{V} \to \mathbb{R}^{d \times d}: (z, \widetilde{z}) \mapsto \lambda(z) \lambda(\widetilde{z}) ((\psi_p(z) - x) ((\psi_p(z) - x)^\top - (\psi_p(\widetilde{z}) - x)^\top)), \quad (41)$$

where  $\lambda = \psi_p^{-1} \# \mu = \operatorname{pr}_p \# \mu$ .

#### D.3.1 CONDITIONS (24) AND (26) FOR h

First we show how to express conditions (24) and (26) in terms of the manifold parameters for h given in (38). To see (24), note that for  $z \in \mathcal{V}$  we have by the chain and product rule as well as Lemma 7 (i) that

$$\begin{split} \|\nabla^{3}h(z)\| &\leq 3\|\psi_{p}''(z)\| \|\psi_{p}'(z)\| + \|x - \psi_{p}(z)\| \|\psi_{p}'''(z)\| \\ &\leq 3M^{2} + (\|x - p\| + \|p - \psi_{p}(z)\|)M \\ &\leq 3M^{2} + (\tau_{\mathcal{M}} + \frac{8}{7}\|z\|)M \\ &\leq (3M + \frac{9}{7}\tau_{\mathcal{M}})M \,. \end{split}$$

Thus (24) holds with, say  $C_{\tau_{\mathcal{M}},M}^{(1)} = \frac{1}{2}(M + \tau_{\mathcal{M}})M$ , and  $\eta = \min\{\tau_{\mathcal{M}},M\}/4$ . Next, for (26) and  $\eta = \min\{\tau_{\mathcal{M}},M\}/4$  we compute for  $z \in B_{\eta}^{\mathrm{T}_{p},\mathcal{M}}(0)$ 

$$h(z) - h(0) \ge \inf_{\widetilde{z} \in \mathcal{V} \setminus B_{\|z\|}^{\mathbf{T}_{p}, \mathcal{M}}(0)} h(\widetilde{z}) - h(0) = \inf_{q \in \psi_{p}(\mathcal{V} \setminus B_{\|z\|}^{\mathbf{T}_{p}, \mathcal{M}}(0))} \frac{1}{2} \|x - q\|^{2} - \frac{1}{2} \|x - p\|^{2}$$

$$\ge \inf_{q \in \mathcal{M} \setminus B_{\|z\|}(p)} \frac{1}{2} \|x - q\|^{2} - \frac{1}{2} \|x - p\|^{2}$$

$$\ge \frac{1}{2} \frac{\tau_{\mathcal{M}} - \tau}{\tau_{\mathcal{M}} + \tau} \|z\|^{2}.$$

In the first inequality we have used  $\psi_p(\mathcal{V}\setminus B_{\|z\|}^{\mathrm{T}_p\mathcal{M}}(0))\subseteq \mathcal{M}\setminus B_{\|z\|}(p)$ , which follows from Lemma 7 (i), whereas in the last inequality we have used Lemma 11 and  $x\in\mathcal{T}(\tau)$ . Thus (26) holds on  $\mathcal{V}$  with  $c=\frac{1}{2}\frac{\tau_{\mathcal{M}}-\tau}{\tau_{\mathcal{M}}+\tau}$  and implies (29).

# D.3.2 Lower bound on $\lambda_{\min}(\Sigma)$ and upper bound on $\sqrt{\det\Sigma}$

We first consider  $\lambda_{\min}(\Sigma)$  for  $\Sigma = \nabla^2 h(0)$  with h given in (38) and when  $x \in \mathcal{T}(\tau)$ . By Lemma 8 and (17) we directly obtain

$$\lambda_{\min}(\Sigma) = \|\Sigma^{-1}\|^{-1} \ge 1 - \|x - p\|\kappa_{\mathcal{M}} \ge 1 - \frac{\tau}{\tau_{\mathcal{M}}}.$$

Due to (17) and the definition of the shape operator  $S_p$ , an upper bound on  $\sqrt{\det \Sigma}$  is given by

$$\sqrt{\det \Sigma} = \sqrt{\det(I_{T_p \mathcal{M}} + S_p^{p-x})} \le \left(1 + \frac{\|p - x\|}{\rho(p, p - x)}\right)^{k/2} \le \left(1 + \frac{\tau}{\tau_{\mathcal{M}}}\right)^{k/2} \le 2^{k/2}.$$

# D.3.3 CONDITION (24) FOR $\overline{f}_2^{jl}$

An upper bound for  $\|\nabla^3 \overline{f}_2^{jl}(z,\widetilde{z})\|$  in terms of M and the third derivatives of  $\frac{d\mu}{dVol_M}$ , where

$$\overline{f}_2^{jl}(z,\widetilde{z}) = \lambda(z)\lambda(\widetilde{z})((\psi_p^j(z) - x)(\psi_p^l(z) - x) - (\psi_p^j(z) - x)(\psi_p^l(\widetilde{z}) - x)),$$

can be obtained by direct differentiation. We don't pursue the complete derivation here and instead say that (24) holds with, say,  $C_{\tau_M,M,\mu}^{(2)}$ .

D.3.4 Upper bound on 
$$||f_0||_{L^1(\mathcal{V})}$$
,  $||f_1^j||_{L^1(\mathcal{V})}$  and  $||\overline{f}_2^{jl}||_{L^1(\mathcal{V}\times\mathcal{V})}$  and  $||\nabla^2 \overline{f}_2^{jl}(0)||$ 

We clearly have  $||f_0||_{L^1(\mathcal{V})} = \lambda(\mathcal{V}) \le 1$ . Next, due to  $||y - x|| \le \operatorname{diam}(\mathcal{M}) + ||p - x|| \le \operatorname{diam}(\mathcal{M}) + \tau$  for  $y \in \mathcal{M}$ , we have

$$||f_1^j||_{L^1(\mathcal{V})} \le ||f_1||_{L^1(\mathcal{V};\mathbb{R}^d)} = \int_{\mathcal{U}} ||y - x|| \, \mathrm{d}\mu(y) = ||\mu||_1 \le \mathrm{diam}(\mathcal{M}) + \tau,$$

and

$$\begin{aligned} \|\overline{f}_{2}^{jl}\|_{L^{1}(\mathcal{V}\times\mathcal{V})} &\leq \|\overline{f}_{2}\|_{L^{1}(\mathcal{V}\times\mathcal{V};\mathbb{R}^{d\times d})} \\ &= \int_{\mathcal{U}\times\mathcal{U}} \|(y-x)(y-x)^{\top} - (y-x)(\widetilde{y}-x)^{\top}\| \,\mathrm{d}(\mu\otimes\mu)(y,\widetilde{y}) \\ &\leq \|\mu\|_{2} + \|\mu\|_{1}^{2} \\ &\leq 2(\mathrm{diam}(\mathcal{M}) + \tau)^{2} \,. \end{aligned}$$

Furthermore for  $v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} \in \mathbb{R}^k \times \mathbb{R}^k$  we have

$$\overline{f}_{2}''(0,0)[v,v] = \lambda(0)^{2} \cdot (2v_{1}v_{1}^{\top} - 2v_{1}v_{2}^{\top} + (p-x)(\mathbb{I}_{p}(v_{1},v_{2}) - \mathbb{I}_{p}(v_{2},v_{2}))^{\top}) + 2\lambda(0)\nabla\lambda(0)^{\top}(v_{1}+v_{2}) \cdot (p-x)(v_{1}-v_{2})^{\top}.$$

Hence, using (20) and the fact that  $||p-x|| ||\mathbb{I}_p(v_1,v_2)|| \leq \frac{\tau}{\tau_M} \leq 1$  for any unit vectors  $v_1,v_2 \in T_p \mathcal{M}$  we obtain

$$\|\nabla^2 \overline{f}_2^{jl}(0)\| \le \|\overline{f}_2''(0,0)\| \le 6\mu(p)^2 + 8\tau\mu(p)\|\operatorname{grad}_{\mathcal{M}}\mu(p)\|.$$

#### D.4 PROOF OF THEOREM 1

We apply Theorem 15 to every point  $x \in \mathcal{T}(\tau)$  and  $p = \pi(x)$  with the graph chart  $\psi = \psi_p$ . It remains to provide universal constants for the bounds in Theorem 15 (i) and (ii). First we bound  $E_1(\sigma; f_0, h, r, \theta)$  with  $E_1$  given in (31). By Section D.3 we can take  $r_{\tau_M, M} = \eta = \min\{\tau_M, M\}/4$ ,  $c_{\tau_M, M} = \frac{1}{2} \frac{\tau_M - \tau}{\tau_M + \tau}$  and  $C_{\tau_M, M, \mu} = \max\{C_{\tau_M, M}^{(1)}, C_{\tau_M, M, \mu}^{(2)}\}$  and hence, due to  $\|f_0\|_{L^1(\mathcal{V})} \leq 1$ ,

$$\begin{split} E_1(\sigma; f_0) &\leq \frac{2^{k/2}}{Z_{\sigma}} \exp\left(-c_{\tau_{\mathcal{M},M}} \log(\sigma)^2\right) \\ &+ \mu(p) \left(2C_{\tau_{\mathcal{M},M,\mu}} \sigma |\log(\sigma)|^3 + \mathrm{G}_0((1 - \frac{\tau}{\tau_{\mathcal{M}}}) \min\{r_{\tau_{\mathcal{M},M}} \sigma^{-1}, |\log(\sigma)|\})\right) \\ &+ C_{\tau_{\mathcal{M},M,\mu}} \sigma |\log(\sigma)| (1 + 2C_{\tau_{\mathcal{M},M,\mu}} \sigma |\log(\sigma)|^3) \,. \end{split}$$

Similarly, due to  $||f_1^{\jmath}||_{L^1(\mathcal{V})} \leq ||\mu||_1$ , we have the bound

$$E_{1}(\sigma; f_{1}) \leq \frac{2^{k/2}}{Z_{\sigma}} \exp\left(-c_{\tau_{\mathcal{M}}, M} \log(\sigma)^{2}\right) \left(\operatorname{diam}(\mathcal{M}) + \tau\right)$$

$$+ \mu(p) \left(2C_{\tau_{\mathcal{M}}, M, \mu} \sigma |\log(\sigma)|^{3} + G_{0}\left(\left(1 - \frac{\tau}{\tau_{\mathcal{M}}}\right) \min\{r_{\tau_{\mathcal{M}}, M} \sigma^{-1}, |\log(\sigma)|\}\right)\right)$$

$$+ C_{\tau_{\mathcal{M}}, M, \mu} \sigma |\log(\sigma)| \left(1 + 2C_{\tau_{\mathcal{M}}, M, \mu} \sigma |\log(\sigma)|^{3}\right).$$

Further, due to  $\|\overline{f}_2^{jl}\|_{L^1(\mathcal{V}\times\mathcal{V})} \leq \|\mu\|_2 + \|\mu\|_1^2$  and the bound on  $\|\nabla^2 \overline{f}_2^{jl}(0)\|$  we have

$$E_{2}(\sigma; \overline{f}_{2}) \leq 2 \frac{2^{k/2}}{Z_{\sigma}} \exp(-c_{\tau_{\mathcal{M}},M} \log(\sigma)^{2}) (\operatorname{diam}(\mathcal{M}) + \tau)^{2}$$

$$+ \mu(p) \frac{6\mu(p) + 8\tau \|\operatorname{grad}_{\mathcal{M}} \mu(p)\|}{1 - \tau/\tau_{\mathcal{M}}} \left( kC_{\tau_{\mathcal{M}},M,\mu} \sigma |\log(\sigma)|^{3} + G_{2} \left( (1 - \frac{\tau}{\tau_{\mathcal{M}}}) \min\{r_{\tau_{\mathcal{M}},M} \sigma^{-1}, |\log(\sigma)| \right) \right)$$

$$+ C_{\tau_{\mathcal{M}},M,\mu} \sigma |\log(\sigma)|^{3} (1 + 2C_{\tau_{\mathcal{M}},M,\mu} \sigma |\log(\sigma)|^{3})$$

Now by Lemma 11 we can pick  $\delta = \frac{\tau_{\mathcal{M}} - \tau}{\tau_{\mathcal{M}} + \tau} \eta^2 = 2c_{\tau_{\mathcal{M}},M} r_{\tau_{\mathcal{M}},M}^2$  and thus

$$\Upsilon(\sigma; \Sigma) \leq \frac{2^{k/2}}{Z_{\sigma}} \exp\left(-\frac{c_{\tau_{\mathcal{M}}, M} r_{\tau_{\mathcal{M}}, M}^2}{\sigma^2}\right).$$

Finally,  $||P_0|| \le (1 - \tau/\tau_{\mathcal{M}})^{-1}$ . Note that since  $\mu(\cdot) \in C^3(\mathcal{M})$  and  $\operatorname{supp} \mu = \mathcal{M}$ , there exist positive lower and upper bounds on  $\mu(\cdot)$  and its derivatives on  $\mathcal{M}$ . This shows that there K and  $\overline{\sigma}$  depending only on  $\tau$ , M and  $\operatorname{diam}(\mathcal{M})$  (as k) and  $\mu$  such that (5) hold and finishes the proof.

# D.5 A USEFUL COROLLARY TO THEOREM 1

The following elementary consequence of Theorem 1 will be useful in the proofs of the results from Section 5.

**Corollary 17.** Suppose that for some  $s: \mathbb{R}^d \to \mathbb{R}^d$  and  $\tau \in (0, \tau_M)$  and  $\epsilon \leq \tau$  we have

$$||s(x) - \pi(x)|| < \epsilon \text{ for all } x \in \mathcal{T}(\tau).$$
 (42)

1561 Then  $s(\mathcal{T}(\tau)) \subseteq \mathcal{T}(\epsilon)$ . Moreover, for  $x \in \mathcal{T}(\epsilon)$  we have  $||s(x) - x|| \le 2\epsilon$  and for  $x \in \mathcal{T}(\tau)$  we have  $||s(x) - x|| \le 2\tau$ .

*Proof.* From (42) and the fact that  $\pi(x) \in \mathcal{M}$  we clearly have  $s(x) \in \mathcal{T}(\epsilon)$  for any  $x \in \mathcal{T}(\tau)$ , i.e.  $s(\mathcal{T}(\tau)) \subseteq \mathcal{T}(\epsilon)$ . The last claims follows then from the triangle inequality  $||s(x) - x|| \le ||s(x) - \pi(x)|| + ||\pi(x) - x||$ .

# E APPROXIMATE RIEMANNIAN GRADIENT FLOW WITH LANDING

In this section  $C = \|\nabla f|_{\mathcal{T}(\tau_{\mathcal{M}})}\|_{\infty}$  and  $L = \operatorname{Lip}(\nabla f|_{\mathcal{T}(\tau_{\mathcal{M}})})$  denote the supremum of  $\nabla f$  and the Lipschitz constant of  $\nabla f$  on  $\mathcal{T}(\tau_{\mathcal{M}})$  for some manifold  $\mathcal{M}$ . We will often assume the following approximation condition on a function  $s \in C^1(\mathbb{R}^d; \mathbb{R}^d)$ :

$$||s(x) - \pi(x)|| < \epsilon, \quad ||s'(x) - P_0(x)|| < \epsilon \text{ for all } x \in \mathcal{T}(\tau).$$

We will also abbreviate the (negative) right hand side of the dynamics (8) by

$$G_s^{\eta}(x) = s'(x)\nabla f(x) + \eta(x - s(x)) \tag{44}$$

## E.1 STATIONARY POINTS OF (8) AND OPTIMALITY CRITERIA

We need to analyze the meaning of approximate stationary points of (8) for the optimization problem (1).

**Lemma 18.** Suppose that  $\tau \in (0, \tau_{\mathcal{M}})$  and  $\sigma > 0$  are such that (43) is satisfied. Moreover, suppose that  $s(\mathcal{T}(\tau)) \subseteq \mathcal{T}(\tau)$ . Suppose that  $\widetilde{\tau} \in (0, \tau]$ ,  $\delta > 0$  and that  $x_* \in \mathcal{T}(\widetilde{\tau})$  is a  $\delta$ -approximate stationary point of (8), i.e.

$$||G_s^{\eta}(x_*)|| \le \delta. \tag{45}$$

Then for  $p_* = \pi(x_*)$  it holds that

$$\|\operatorname{grad}_{\mathcal{M}} f(p_*)\| \le 2(L + C + \eta)\epsilon + 2\frac{\widetilde{\tau}/\tau_{\mathcal{M}}}{1 - \widetilde{\tau}/\tau_{\mathcal{M}}}C + \delta.$$
 (46)

*Proof.* We have the estimates

$$\|\eta(s(x_*) - x_*) - \eta(\pi(x_*) - x_*)\| \le \eta\epsilon$$

and

$$\begin{aligned} \|P_{0}(\pi(x_{*}))\nabla f(\pi(x_{*})) - s'(x_{*})\nabla f(s'(x_{*}))\| \\ &\leq \|P_{0}(\pi(x_{*}))\nabla f(\pi(x_{*})) - P_{0}(\pi(x_{*}))\nabla f(s(x_{*}))\| \\ &+ \|P_{0}(\pi(x_{*}))\nabla f(s(x_{*})) - P_{0}(x_{*})\nabla f(s(x_{*}))\| \\ &+ \|P_{0}(x_{*})\nabla f(s(x_{*})) - s'(x_{*})\nabla f(s(x_{*}))\| \\ &\leq (\operatorname{Lip}(\nabla f) + \|\nabla f|_{\mathcal{T}(\tau)}\|_{\infty})\epsilon + \|P_{0}(\pi(x_{*})) - P_{0}(x_{*})\|\|\nabla f|_{\mathcal{T}(\tau)}\|_{\infty}, \end{aligned}$$

where we have used that  $P_0(\pi(x_*))$  is an orthogonal projection and hence has unit norm. By Lemma 23 applied to  $z_1 = s'(x_*)\nabla f(s(x_*))$  and  $z_2 = \eta(s(x_*) - x_*)$ , Lemma 9 and the fact that  $P_0(\pi(x_*))\nabla f(\pi(x_*)) = \operatorname{grad}_{\mathcal{M}} f(\pi(x_*))$ , the inequality (46) follows.

## E.2 LIE DERIVATIVE OF MANIFOLD DISTANCE

**Lemma 19.** Suppose that  $\tau \in (0, \tau_M)$  and  $\sigma > 0$  is such that (43) is satisfied. Then it holds

$$-\langle \nabla d(x), G_s^{\eta}(x) \rangle \leq -2\eta d(x) + \epsilon (C+\eta) \sqrt{2 d(x)} \text{ for all } x \in \mathcal{T}(\tau).$$

*Proof.* We have for  $x \in \mathcal{T}(\tau)$ 

$$\langle x - \pi(x), -G_s^{\eta}(x) \rangle$$

$$= -2\eta \, d(x) + \langle \pi(x) - x, (s'(x) - P_0(x)) \nabla f(s(x)) \rangle$$

$$+ \eta \, \langle x - \pi(x), s(x) - \pi(x) \rangle$$

$$\leq -2\eta \, d(x) + \|\pi(x) - x\| \|s'(x) - P_0(x)\| \|\nabla f(s(x))\|$$

$$+ \eta \|\pi(x) - x\| \|s(x) - \pi(x)\|$$

$$\leq -2\eta \, d(x) + (\|\nabla f|_{\mathcal{T}(\mathcal{T})}\|_{\infty} \epsilon + \eta \epsilon) \sqrt{2 \, d(x)} .$$

## E.3 Exact landing with $\sigma = 0$

Here we establish the following noiseless version of Theorem 3.

**Theorem 20.** Consider the flow (8) for  $\sigma = 0$  and  $\eta \ge 0$  with  $x(0) \in \mathcal{T}(\tau)$  for some  $\tau \in (0, \tau_{\mathcal{M}})$ . Then the solution x(t) exists for all times  $t \ge 0$  and is contained in  $\mathcal{T}(\tau)$ . Moreover every accumulation point  $x_*$  of this flow satisfies  $\operatorname{grad}_{\mathcal{M}} f(x_*) = 0$  and is at most  $\tau$  away from the point  $p_* = \pi(x_*) \in \mathcal{M}$  with  $\|\operatorname{grad}_{\mathcal{M}} f(p_*)\| \le L\tau$ . Moreover, if  $\eta > 0$ , then  $x_* = p_* \in \mathcal{M}$  and  $\operatorname{grad}_{\mathcal{M}} f(p_*) = 0$ , i.e. every accumulation point is critical.

*Proof.* Note first that

$$\frac{\mathrm{d}}{\mathrm{d}t}\,\mathrm{d}(x) = \langle x - \pi(x), \dot{x} \rangle = -\eta \|x - \pi(x)\|^2 = -2\eta\,\mathrm{d}(x)\,,$$

i.e.  $d(x(t)) = e^{-2\eta t} d(x(0))$  and hence  $x(t) \in \mathcal{T}$  for all  $t \ge 0$ . Moreover, if  $\eta > 0$  then  $d(x(t)) \to 0$  for  $t \to \infty$ . Let us now consider  $f \circ \pi$  as a Lyapunov function for (9) with  $\sigma = 0$ . We observe that

$$\begin{split} \frac{\mathrm{d}}{\mathrm{d}\,t} f(\pi(x)) &= -\nabla f(\pi(x))^{\top} P_0(x)^{\top} P_0(x) \nabla f(\pi(x)) \\ &= -\left\langle \mathrm{P}_{\mathrm{T}_{\pi(x)}\,\mathcal{M}} \, \nabla f(x), H_x^{-2} \, \mathrm{P}_{\mathrm{T}_{\pi(x)}\,\mathcal{M}} \, \nabla f(x) \right\rangle \\ &= -\|\mathrm{grad}\, f(x)\|_{H_x}^2 \,, \end{split}$$

with  $\|v\|_{H_x}^2 = \|H_x^{-1}v\|_{\mathrm{T}_{\pi(x)}\mathcal{M}}^2$ . This shows that f(x(t)) is non-increasing for  $t \to \infty$ . Since  $\mathcal{T}$  is bounded,  $f_* = \lim_{t \to \infty} f(x(t))$  is finite. Moreover, we have

$$\int_0^\infty \|\operatorname{grad} f(x(t))\|_{H_{x(t)}}^2 dt = f(x(0)) - f_* < \infty.$$
(47)

Notice that  $\{x(t) \mid t \geq 0\} \subseteq \overline{\mathcal{T}}(\mathrm{dist}_{\mathcal{M}}(x(0)))$  being compact. Since  $\mathcal{T} \to \mathbb{R} : z \mapsto \|\mathrm{grad}\, f(z)\|_{H_z}^2$  is continuous, it is uniformly continuous on  $\overline{\mathcal{T}}(\mathrm{dist}_{\mathcal{M}}(x(0)))$ . This, together with (47) implies, by Barbalat's lemma Farkas & Wegner (2016), that  $\lim_{t \to \infty} \|\mathrm{grad}\, f(x(t))\|_{H_{x(t)}}^2 = 0$ . Clearly any accumulation point  $x_*$  of  $\{x(t) \mid t \geq 0\}$  satisfies  $H_{x_*}^{-1} \operatorname{grad} f(x_*) = 0$ , i.e.  $\operatorname{grad} f(x_*) = 0$ . Then  $\pi(x_*) \in \mathcal{M}$  is an accumulation point of  $\{\pi(x(t)) \mid t \geq 0\} \subseteq \mathcal{M}$  and

$$\|\operatorname{grad} f(\pi(x_*))\| = \|\operatorname{grad} f(\pi(x_*)) - \operatorname{grad} f(x_*)\| \le L\|\pi(x_*) - x_*\| = L\tau.$$

## E.4 PROOF OF THEOREM 3

Note first that, by Theorem 1, if (7) hold with  $\epsilon \to \epsilon'$ , then 43 are satisfied with  $\epsilon = \epsilon' + K\sigma |\log(\sigma)|^3$  for  $\sigma \in (0, \overline{\sigma}(\tau, \mathcal{M}, \mu))$ . In the following we will write  $\epsilon > 0$  for the latter quantity. First let us analyze the manifold distance of (8). We have by Lemma 19 whenever  $x(t) \in \mathcal{T}(\tau)$  that

$$\frac{\mathrm{d}}{\mathrm{d}t}\,\mathrm{d}(x) = -\left\langle \nabla\,\mathrm{d}(x), G_s^{\eta}(x) \right\rangle \le -2\eta\,\mathrm{d}(x) + \epsilon(C+\eta)\sqrt{2\,\mathrm{d}(x)}\,,$$

where the right hand side is non-positive iff

$$\operatorname{dist}_{\mathcal{M}}(x) \geq \frac{\epsilon}{2} \left( \frac{C}{\eta} + 1 \right) =: \tau_0.$$

Note that  $\tau_0 < \tau$  if  $\epsilon < 2\tau/(1+C/\eta)$ . In particular if  $\mathrm{dist}_{\mathcal{M}}(x(0)) \in [\tau_0,\tau)$ , then  $\overline{T}(\mathrm{dist}_{\mathcal{M}}(x(0)))$  is invariant w.r.t. the flow (8) and x(t) exists for all  $t \geq 0$ . Moreover, by a similar argument as in the standard proof of Lyapunov's direct method Khalil & Grizzle (2002), for each  $\delta > 0$  there exists some T>0 such that for all  $t \geq T$  it holds that  $x(t) \in \mathcal{T}(\tau_0+\delta)$ . If s is a gradient field, i.e.  $s=\nabla g$  for some function  $g\in C^1(\mathbb{R}^d)$ , then  $V(x):=f(s(x))+\eta(\|x\|^2/2-g(x))$  satisfies  $\nabla V(x)=G_s^\eta(x)$ . Similar as in the proof of Theorem 20 we can take  $G_s$  as a Lyapunov function to obtain

$$\frac{\mathrm{d}}{\mathrm{d}t}V(x) = -\|G_s^{\eta}(x)\|^2$$

along the dynamics (8). By a similar argument V(x) is non-increasing and every accumulation point  $x_*$  of  $\{x(t) \mid t \geq 0\}$  satisfies  $G^\eta_s(x)(x_*) = 0$  and belongs to  $\overline{\mathcal{T}}(\tau_0)$ . Note that the condition  $s(\mathcal{T}(\tau)) \subseteq \mathcal{T}(\tau)$  in Lemma 18 is satisfied when  $\epsilon \leq \tau$ . Applying Lemma 18 with  $\delta = 0$  and  $\widetilde{\tau} = \tau_0$  yields (11) and proves Theorem 3 for the case when s is a gradient field. Now consider the more general case when s is not necessarily a gradient field. In this case consider  $V(x) = f(s(x)) + \eta \operatorname{d}(x)$  to obtain

$$\frac{\mathrm{d}}{\mathrm{d}\,t}V(x) = -\langle s'(x)\nabla f(s(x)) + \eta(x - \pi(x)), G_s^{\eta}(x)\rangle$$

$$\leq -\|G_s^{\eta}(x)\|^2 + \eta\|s(x) - \pi(x)\|\|G_s^{\eta}(x)\|$$

$$< -\|G_s^{\eta}(x)\|^2 + \eta\epsilon\|G_s^{\eta}(x)\|.$$

Thus, if  $\|G^\eta_s(x)\| \geq \eta\epsilon$  we have  $\frac{\mathrm{d}}{\mathrm{d}\,t}V(x) \leq 0$ . By a barrier function argument Khalil & Grizzle (2002) this implies that every accumulation point  $x_*$  of  $\{x(t) \mid t \geq 0\}$  satisfies  $\|G^\eta_s(x_*)\| \leq \eta\epsilon$  and belongs to  $\overline{\mathcal{T}}(\tau_0)$ . Again applying Lemma 18 with  $\delta = \eta\epsilon$  and  $\widetilde{\tau} = \tau_0$  yields (11) and finishes the proof.

## F DISCRETIZED RIEMANNIAN GRADIENT FLOW AND DESCENT

In this section we provide proof of Theorem 5 as well as analysis of the approximate Riemannian gradient descent and the discretized landing flow. As before  $C = \|\nabla f|_{\mathcal{T}(\tau_{\mathcal{M}})}\|_{\infty}$  and  $L = \operatorname{Lip}(\nabla f|_{\mathcal{T}(\tau_{\mathcal{M}})})$  denote the supremum of  $\nabla f$  and the Lipschitz constant of  $\nabla f$  on  $\mathcal{T}(\tau_{\mathcal{M}})$  for some manifold  $\mathcal{M}$ . Moreover, the following bounds will be useful:

**Lemma 21.** If  $\tau \in (0, \tau_M)$  and (43) holds for some  $\epsilon > 0$ , then

$$\sup_{x \in \mathcal{T}(\tau)} \|s'(x)\|_{\infty} \le \epsilon + \frac{1}{1 - \tau/\tau_{\mathcal{M}}}.$$

Additionally, if  $\epsilon \in (0, \tau]$ , then

$$||G_s^{\eta}(x)|| \le C||s'(x)|| + \eta ||s(x) - x|| \le C(\epsilon + \frac{1}{1 - \tau/\tau_M}) + 2\eta \tau \text{ for } x \in \mathcal{T}(\tau)$$

*Proof.* Via triangle inequality, Lemma 8 and the definition of C.

## F.1 DISCRETIZED APPROXIMATE RIEMANNIAN GRADIENT FLOW

In this section we analyze the discretized version of (9), specifically the corresponding gradient descent

$$x_{k+1} = x_k - \gamma_k \nabla F_\sigma^{\eta}(x_k) \,, \tag{48}$$

for some sequence of step sizes  $\{\gamma_k\}_{k=1}^\infty$ . The selection of  $\gamma_k$  can be inferred from any standard analysis of gradient descent for  $F_\sigma^\eta$  to guarantee that all accumulation points of the resulting sequence  $\{x_k\}_{k=0}^\infty$  are stationary points of  $F_\sigma^\eta$ . Since we can only interpret stationary points of  $F_\sigma^\eta$  in terms of our original problem (1) when they are contained in a tubular neighborhood  $\mathcal{T}(\tau)$  for some  $\tau \in (0, \tau_\mathcal{M})$  (see Lemma 18), we drive conditions on the step-size to ensure  $\{x_k\}_{k=0}^\infty \subseteq \mathcal{T}(\tau)$ .

**Theorem 22.** Let  $\tau \in (0, \tau_{\mathcal{M}}/2)$  and  $\sigma > 0$  be such that (43) holds for some  $\epsilon \in (0, \frac{\eta \tau}{2(C+\eta)}]$ . Then if  $\gamma_k \in [0, \gamma_{\text{tubular}}]$  with

$$\gamma_{\rm tubular}(\epsilon,\tau,\eta) = \tau \cdot \min \left\{ \frac{1}{2(C(\epsilon+2)+2\eta\tau)}, \frac{\frac{1}{4}\eta\tau - \frac{1}{2}(C+\eta)\epsilon}{4(C(\epsilon+4)+3\eta\tau)^2} \right\}$$

and  $x_0 \in \mathcal{T}(\tau)$ , the iterates  $x_k$  of the discretized flow (48) belong to  $\mathcal{T}(\tau)$ .

*Proof.* First we show that  $x_k \in \mathcal{T}(\tau)$  implies  $x_{k+1} \in \mathcal{T}(\tau)$ . If  $x_k \in \mathcal{T}(\tau/3)$ , then, since  $\gamma_k \|\nabla F_\sigma^\eta(x_k)\| \leq \tau/2$  by Lemma 21, it follows  $x_{k+1} \in \mathcal{T}(\tau)$ . Hence assume  $x_k \in \mathcal{T}(\tau) \setminus \mathcal{T}(\tau/3)$ , i.e.  $\sqrt{2 \operatorname{d}(x_k)} \geq \tau/2$ . Let us write

$$d(x_{k+1}) = d(x_k) - \gamma_k \langle \nabla d(x_k), \nabla F_{\sigma}^{\eta}(x_k) \rangle + \frac{1}{2} \gamma_k^2 ||(I - P_0)||_{\mathcal{T}(3\tau/2)} ||_{\infty} ||\nabla F_{\sigma}^{\eta}(x_k)||^2,$$

where we have used  $\nabla^2 d = I - P_0$  and  $x_{k+1} \in \mathcal{T}(3\tau/2)$ , since again  $\gamma_k \|\nabla F_{\sigma}^{\eta}(x_k)\| \leq \tau/2$ . Then by Lemma 19 

 $-\langle \nabla d(x), \nabla F_{\sigma}^{\eta}(x) \rangle \leq -2\eta d(x) + \epsilon (C+\eta) \sqrt{2 d(x)}$ .

Now, Lemma 21 and the fact that  $(1 - 3\tau/(2\tau_M))^{-1} < 4$  imply 

$$\delta = \frac{1}{2} \| (I - P_0)|_{\mathcal{T}(3\tau/2)} \|_{\infty} \| \nabla F_{\sigma}^{\eta}|_{\mathcal{T}(\tau)} \|_{\infty}^2 \le 4(C(\epsilon + 4) + 3\eta\tau)^2$$

Therefore we have

$$d(x_{k+1}) - d(x_k) \le \gamma_k(-2\eta d(x_k) + \epsilon(C+\eta)\sqrt{2 d(x_k)} + \gamma_k \delta),$$

where right hand side is non-positive for all  $\sqrt{2 d(x_k)} \ge \tau/2$  if

$$-\frac{1}{4}\eta\tau^2 + \frac{1}{2}(C+\eta)\tau\epsilon + \gamma_k\delta \le 0.$$

or, equivalently,

$$\gamma_k \le \frac{1}{\delta} \left( \frac{1}{4} \eta \tau^2 - \frac{1}{2} (C + \eta) \tau \epsilon \right).$$

If these are satisfied, then  $d(x_{k+1}) \leq d(x_k) \leq \tau^2/8$  and therefore  $x_{k+1} \in \mathcal{T}(\tau)$ . In all,  $\{x_k\}_{k=0}^{\infty} \subseteq \mathcal{T}(\tau)$  $\mathcal{T}(\tau)$  provided that  $x_0 \in \mathcal{T}(\tau)$ .

#### Proof of Theorem 5

Let us write  $\epsilon' = \epsilon + K(\tau, \mathcal{M}, \mu) \sigma |\log(\sigma)|^3 \le \tau/2$ . First we note that  $\{x_k\}_{k=0}^{\infty} \subseteq \mathcal{T}(\epsilon')$  as soon as  $x_0 \in \mathcal{T}(\tau/2)$ , because  $s(\mathcal{T}(\tau)) \subseteq \mathcal{T}(\epsilon')$  and  $x_k - \gamma_k s'(x_k) \nabla f(x_k) \in \mathcal{T}(\tau)$ , since  $\gamma_k \le \frac{\tau}{2C(\epsilon' + (1 - \tau/\tau_{\mathcal{M}})^{-1})} \le \frac{\tau}{2\|s'(x_k) \nabla f(x_k)\|},$ 

$$\gamma_k \le \frac{\tau}{2C(\epsilon' + (1 - \tau/\tau_{\mathcal{M}})^{-1})} \le \frac{\tau}{2\|s'(x_k)\nabla f(x_k)\|}$$

where first inequality is due to  $(1 - \tau/\tau_M)^{-1} \le 2$  and the second inequality due to Lemma 21. For brevity let us denote  $y_k = x_k - \gamma_k s'(x_k) \nabla f(x_k)$  and  $z_k = x_k - \gamma_k P_0(x_k) \nabla f(\pi(x_k))$ . Also let  $L_0 = \operatorname{Lip}(\nabla (f \circ \pi)|_{\mathcal{T}(\tau)})$ . By Lemma 10 and  $(1 - \tau/\tau_{\mathcal{M}})^{-1} \leq 2$  it follows that

$$L_0 \le C \|P_0'|_{\mathcal{T}}(\tau)\| + (1 - \tau/\tau_{\mathcal{M}})^{-1}L \le 8C \left(2(\frac{3}{\tau_{\mathcal{M}}} + \tau M) + \frac{1}{\tau_{\mathcal{M}}}\right) + 2L.$$

Then, since  $||x_k - \pi(x_k)|| \le \epsilon'$  for  $k \ge 1$ , we have

$$f(x_{k+1}) - f(x_k) = f(s(y_k)) - f(x_k)$$

$$= f(\pi(z_k)) - f(\pi(x_k)) + f(\pi(x_k)) - f(x_k)$$

$$+ f(s(y_k)) - f(\pi(y_k)) + f(\pi(y_k)) - f(\pi(z_k))$$

$$\leq -\gamma_k \left\langle P_0(x_k) \nabla f(\pi(x_k)), P_0(x_k) \nabla f(\pi(x_k)) \right\rangle + \gamma_k^2 \frac{L_0}{2} \|P_0(x_k) \nabla f(\pi(x_k))\|^2$$

$$+ C \|x_k - \pi(x_k)\| + C\epsilon + L_0 \gamma_k (C\epsilon + L \|x_k - \pi(x_k)\|)$$

$$\leq -\gamma_k (1 - \frac{L_0}{2} \gamma_k) \|P_0(x_k) \nabla f(\pi(x_k))\|^2 + (2C + L_0 \gamma_k (C + L))\epsilon'.$$

Now let  $\gamma_k \in [\gamma_{\min}, \gamma_{\max}] \subseteq (0, \frac{2}{I_{i0}})$ . Then summing over  $k = 1, \dots, N$  yields

$$\gamma_{\min}(1 - \frac{1}{2}\gamma_{\max}L_0)\frac{1}{N}\sum_{k=1}^{N} \|P_0(x_k)\nabla f(\pi(x_k))\|^2 \le \frac{f(x_1) - f(x_{N+1})}{N} + (2C + L_0\gamma_k(C+L))\epsilon'.$$

Now note that by Lemma 9 and the fact that  $(1 - \epsilon'/\tau_M)^{-1} \le 2$  it holds

$$\|\operatorname{grad}_{\mathcal{M}} f(\pi(x_k))\|^2 \le 2\|(P_0(\pi(x_k)) - P_0(x_k))\nabla f(\pi(x_k))\|^2 + 2\|P_0(x_k)\nabla f(\pi(x_k))\|^2$$

$$< 8C^2(\epsilon'/\tau_{\mathcal{M}})^2 + 2\|P_0(x_k)\nabla f(\pi(x_k))\|^2,$$

which implies

$$\gamma_{\min}(1 - \frac{1}{2}\gamma_{\max}L_0)\frac{1}{N}\sum_{k=0}^{N}\|\operatorname{grad}_{\mathcal{M}}f(\pi(x_k))\|^2 \le \frac{4D}{N} + 8C^2(\epsilon'/\tau_{\mathcal{M}})^2 + 2(2C + L_0\gamma_k(C+L))\epsilon'.$$

This finishes the proof.

F.3 AUXILIARY RESULTS

#### F.3.1 A LEMMA ON NORMS OF ORTHOGONAL VECTORS

**Lemma 23.** Let  $x, y \in \mathbb{R}^n$  be orthogonal and  $\epsilon > 0$ . If there exists some  $z \in \mathbb{R}^n$  with  $||x - z|| \le \epsilon$  and  $||y - z|| \le \epsilon$ , then  $||x|| \le 2\epsilon$ ,  $||y|| \le 2\epsilon$  and  $||z|| \le \sqrt{2}\epsilon$ . If on the other hand for some  $\delta > 0$  there exist some  $z_1, z_2 \in \mathbb{R}^n$  with  $||x - z_1|| \le \epsilon$ ,  $||y - z_2|| \le \epsilon$  and  $||z_1 - z_2|| \le \delta$ , then  $||x|| \le 2\epsilon + \delta$ ,  $||y|| \le 2\epsilon + \delta$  and  $||z|| \le \sqrt{2}\epsilon + \delta/\sqrt{2}$ .

*Proof.* The first claim follows by geometric considerations. The second follows from the first by noting that the midpoint  $z=(z_1+z_2)/2$  satisfies  $\|z-z_1\| \le \delta/2$  and  $\|z-z_2\| \le \delta/2$  and hence  $\|x-z\| \le \epsilon + \delta/2$  and  $\|y-z\| \le \epsilon + \delta/2$ , i.e.  $\|x\| \le 2\epsilon + \delta$ ,  $\|y\| \le 2\epsilon + \delta$  and  $\|z\| \le \sqrt{2}\epsilon + \delta/\sqrt{2}$ .  $\square$ 

#### F.3.2 GAUSSIAN TAIL BOUNDS

The following elementary lemmas give simplifications for some of the constants appearing in the proof of Theorem 13 and Theorem 14.

**Lemma 24.** It holds for R > 0 that

$$G_0(R) := \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n \setminus B_R(0)} e^{-\frac{1}{2}||u||^2} du \le 2e^{-\frac{1}{2n}R^2}$$

and for  $R \geq 2n$  that

$$\mathrm{G}_2(R) := \frac{1}{(2\pi)^{n/2}} \int_{\mathbb{R}^n \backslash B_R(0)} \|u\|^2 e^{-\frac{1}{2}\|u\|^2} \, du \leq \frac{2^{2-n/2}}{\Gamma(n/2)} e^{-\frac{1}{4}R^2} \, .$$

In particular  $G_0(|\log(\sigma)|) = O(\sigma^l)$  and  $G_2(|\log(\sigma)|) = O(\sigma^l)$  for any  $l \ge 1$  as  $\sigma \to 0$ .

*Proof.* See Majerski (2015).

## G FURTHER NUMERICAL EXPERIMENTS AND IMPLEMENTATION DETAILS

#### G.1 OPTIMIZATION OVER O(n)

## G.1.1 IMPLEMENTATION DETAILS

In all of our experiments, we discretize the flow equation 8 using the Euler scheme with a step size of  $t_{\rm step} = 1 \cdot 10^{-4}$  and set the landing gain  $\eta = 3 \cdot 10^3$ .

**Data generation:** In our experiments, we take  $Q = \operatorname{diag}(1, \dots, n)$  and a randomly sampled symmetric  $A \in \mathbb{S}^{n \times n}$  with  $\mathcal{N}(0, 1)$ -entries.

**Score architecture:** We use the following score architecture:

$$s_{\sigma}(X) = \frac{1}{\sigma} \widetilde{s}_{\sigma}(X) \,,$$

with  $X=(x_1 \quad \cdots \quad x_n) \in \mathbb{R}^{n \times n}$  and  $Y=(y_1 \quad \cdots \quad y_n)=\widetilde{s}_{\sigma}(X)$ , where

$$y_i = \mathrm{MLP}_{l,w}([r; x_i; \sigma]) \text{ for } i = 1, \dots, n,$$

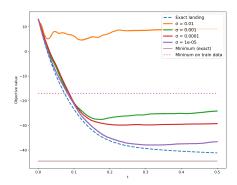
and MLP a fully connected multi-layer perceptron with ReLU activation function, l layers of width w. The features  $r = r(X) \in \mathbb{R}^m$  are

$$r_j(X) = \operatorname{tr}(Q_j X K_j X^\top), \quad j = 1, \dots, m$$

where  $Q_j, K_j \in \mathbb{R}^{n \times n}$  are learnable weight matrices (shared for all i = 1, ..., n). For n = 10 we take l = 4, w = 512, m = 128 and for n = 20 we take l = 4, w = 2048, m = 512.

**Diffusion and training parameters:** We train minimizing equation 23 with the Adam optimizer, early stopping (i.e.  $t \sim \mathrm{Unif}[\epsilon, T]$ ) with T=3 and  $\epsilon=10^{-4}$  and a cosine learning rate scheduling from  $1 = 10^{-3}$  to  $1 = 5 \cdot 10^{-5}$ . We use  $N_{\mathrm{epochs}} = 10000$  and  $N_{\mathrm{epochs}} = 50000$  epochs for n=10 and n=20, respectively.

#### G.1.2 EXPERIMENTS



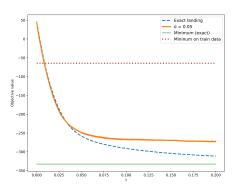


Figure 3: Objective value vs. flow time t for the orthogonal manifold for n=10 (left, different  $\sigma > 0$ ) and for n=20 (right,  $\sigma = 0.05$ )

## G.2 DATA-DRIVEN REFERENCE TRACKING

#### G.2.1 BENCHMARK SYSTEMS AND TRACKING GOALS

**Benchmark systems:** We consider two classical benchmark systems: The double pendulum and the unicycle car model LaValle (2006). For the double pendulum the state is  $x=(\theta_1,\omega_1,\theta_2,\omega_2)$  with  $\omega_i=\dot{\theta}_i$ , gravity g, masses  $m_1,m_2$ , lengths  $l_1,l_2$ , dampings  $d_1,d_2$ , control torque u applied at joint 1 (first pendulum), and  $\Delta\theta:=\theta_2-\theta_1$ .

$$\dot{ heta}_1 = \omega_1,$$
  $\dot{ heta}_2 = \omega_2,$   $\mathbf{M}( heta) \, \ddot{ heta} + \mathbf{C}( heta, \dot{ heta}) + \mathbf{G}( heta) + \mathbf{D} \, \dot{ heta} = oldsymbol{ au}, \quad oldsymbol{ heta} = \begin{pmatrix} heta_1 \\ heta_2 \end{pmatrix}, \quad oldsymbol{ au} = \begin{pmatrix} u \\ 0 \end{pmatrix},$ 

with

$$\mathbf{M}(\theta) = \begin{pmatrix} (m_1 + m_2)l_1^2 & m_2l_1l_2\cos\Delta \\ m_2l_1l_2\cos\Delta\theta & m_2l_2^2 \end{pmatrix}, \qquad \mathbf{D} = \begin{pmatrix} d_1 & 0 \\ 0 & d_2 \end{pmatrix},$$

$$\mathbf{C}(\theta, \dot{\theta}) = \begin{pmatrix} -m_2l_1l_2\sin\Delta\theta & \omega_2^2 \\ m_2l_1l_2\sin\Delta\theta & \omega_1^2 \end{pmatrix}, \qquad \mathbf{G}(\theta) = \begin{pmatrix} (m_1 + m_2)gl_1\sin\theta_1 \\ m_2gl_2\sin\theta_2 \end{pmatrix}.$$

We pick the output  $y=(\theta_1,\theta_2)$  and set  $m_1=l_1=g=1, m_2, l_2=0.5$  and  $d_1,d_2=0.1$ . For the unicycle car model the dynamics is given by  $x=(x,y,\theta)$  with

$$\dot{\mathbf{x}} = v \cos(\theta), \quad \dot{\mathbf{y}} = v \sin(\theta), \quad \dot{\theta} = \omega$$

and the input  $u=(v,\omega)$  and output  $y=x=(x,y,\theta)$ . Here  $(x,y), v, \theta, \omega$  is the car's position, velocity, angle and angular velocity, respectively.

**Tracking goals:** For the double pendulum system the goal is to track a reference trajectory r via the first joint angle  $\theta_1$  and we pick the optimal control objective f to be 15 with

$$Q = \begin{pmatrix} 10 & 0 \\ 0 & 10 \end{pmatrix}, \quad R = 0.01,$$

while for the unicycle car model the goal is to track a positional reference  $r = (r_x, r_y)$ , i.e. we pick

$$Q = \begin{pmatrix} 10 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 0 \end{pmatrix} \,, \quad R = \begin{pmatrix} 0.01 & 0 \\ 0 & 0.01 \end{pmatrix} \,.$$

Here R is some small penalty on the input u to keep it bounded during the optimization.

#### G.2.2 IMPLEMENTATION DETAILS

 **Discretization:** We discretize both continuous-time dynamics  $\dot{x} = f_{\text{cont}}(x, u)$  via the RK4-method and discretization step  $\Delta t$  to obtain the discrete-time dynamics (13) as

$$f(x, u) = x + \frac{\Delta t}{6} (k_1 + 2k_2 + 2k_3 + k_4) \text{ where } \begin{cases} k_1 &= f_{\text{cont}}(x, u) \\ k_2 &= f_{\text{cont}}(x + \frac{1}{2}\Delta t k_1, u) \\ k_3 &= f_{\text{cont}}(x + \frac{1}{2}\Delta t k_2, u) \\ k_4 &= f_{\text{cont}}(x + \Delta t k_3, u) \end{cases}$$

We use  $\Delta t = 0.1$  for the double pendulum and  $\Delta t = 0.05$  for the unicycle model.

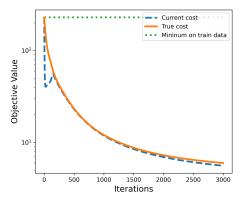
**Data generation:** To generate trajectories, we use i.i.d. random inputs  $u_k \sim \mathrm{Unif}[-5,5]$  for the double pendulum and  $u_k = (v_k, \omega_k) \sim \mathrm{Unif}[0,1] \otimes \mathcal{N}(0,25)$  and a horizon of  $N_h = 100$  for both system. We use  $N_{\mathrm{data}} = 50000$  trajectories for the double pendulum and  $N_{\mathrm{data}} = 20000$  trajectories for the unicycle model.

**Score architecture:** The score architecture is a 1-dimensional version of the standard UNet architecture Ronneberger et al. (2015) with a sin-cos time-embedding Song et al. (2020) and residual connections, where the different input-, state- and output dimensions are concatenated and treated as additional channels. The down- and upsampling convolutions are done w.r.t. the temporal dimension and channels.

**Diffusion and training parameters:** Same as in Appendix G.1.1 with this time  $N_{\rm epochs} = 50000$  training epochs. For the DRGD step-size we pick a fixed step-size of  $\gamma = 0.001$ .

#### G.2.3 EXPERIMENTS

In Figure 4 we present the objective value evolution for the experiment from Section 6.2.



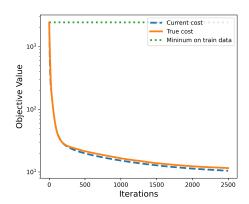
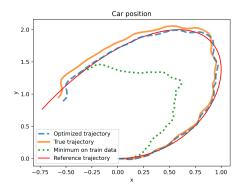


Figure 4: Denoising Riemannian gradient descent: Objective value f vs the iteration count j for double pendulum (left) and unicycle car model (right). Note the logarithmic scale on the y-axis. The current cost (blue, dashed) is the objective value  $f(\boldsymbol{u}_j, \boldsymbol{y}_j)$  at the current (in general infeasible  $(\boldsymbol{u}_j, \boldsymbol{y}_j) \notin \mathcal{M}_{\text{IO}}$ ) iterate, while the true cost (orange) is the value  $f(\boldsymbol{u}_j, \boldsymbol{y}_j^{\text{true}})$ , with  $\boldsymbol{y}_j^{\text{true}}$  obtained by simulating (13) with input  $\boldsymbol{u}_j$ .

In Figure 5 we show optimized trajectories for two other reference trajectories. Note that we have set our iteration budget at N=4000, while the objective is still decreasing. How to accelerate the denoising Riemannian gradient descent without losing feasibility is a core question for future work.



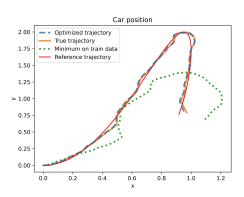


Figure 5: Denoising Riemannian gradient descent: Unicycle car position (right) with the optimized output trajectory  $y^*$  (blue, dashed), the true system trajectory  $y^{\text{true}}$  (orange), the initial trajectory  $y_0$  (green, dotted) and the reference trajectory  $y_0$  (red)