



CENTRE FOR **STOCHASTIC GEOMETRY**  
AND ADVANCED **BIOIMAGING**



Anton Mallasto and Aasa Feragen

## **Wrapped Gaussian Process Regression on Riemannian Manifolds**

No. 03, February 2018

# Wrapped Gaussian Process Regression on Riemannian Manifolds

Anton Mallasto and Aasa Feragen

Department of Computer Science, University of Copenhagen, {mallasto,aasa}@di.ku.dk

## Abstract

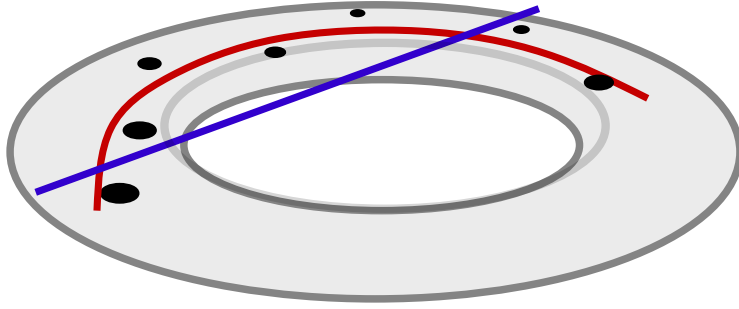
Gaussian process (GP) regression is a powerful tool in non-parametric regression providing uncertainty estimates. However, it is limited to data in vector spaces. In fields such as shape analysis and diffusion tensor imaging, the data often lies on a manifold, making GP regression non-viable, as the resulting predictive distribution does not live in the correct geometric space. We tackle the problem by defining wrapped Gaussian processes (WGP) on Riemannian manifolds, using the probabilistic setting to generalize GP regression to the context of manifold-valued targets. The method is validated empirically on diffusion tensor imaging (DTI) data and in the Kendall shape space, endorsing WGP regression as an efficient and flexible tool for manifold-valued regression.

## 1 Introduction

Regressing curves when the training data  $\{(x_i, y_i)\}_{i=1}^N$  is Euclidean is well studied. When  $y_i$  are manifold-valued, on the other hand, poses difficulties, due to the lack of the vector space structure. Applying Euclidean statistical methods to manifold-valued data does not respect the intrinsic structure, hence the product of inference might not represent the same object category anymore. For example, see Fig. 1, where not all the points on the linearly regressed curve live on the disc anymore.

Sometimes the data observed is inherently uncertain. In this case, it is favorable to estimate a distribution over possible regressed functions, which is what Gaussian process (GP) regression achieves, yielding uncertainty estimates of the resulting inference in a tractable manner. Furthermore, GP regression is an example of Bayesian inference, where it is possible to incorporate prior knowledge to aid the inference. These qualitative properties are the main motivation for us to consider a generalization of GP regression to Riemannian manifolds.

**Related work:** Fletcher [5] generalized linear regression to handle manifold-valued data with one dimensional Euclidean covariates by *geodesic regression*. Geodesic regression was then extended to include multiple dimensional covariates [11, 18]. Furthermore, [18] considered a kriging (GP regression) method, that takes advantage of the multivariate geodesic regression model to form a reference coordinate



**Figure 1: Why geometrically intrinsic regression is important.** Imagine your data points (black dots) lie on the grey disc. The blue line is the result of carrying out linear regression viewing the data points as elements in a Euclidean space. This clearly escapes the natural geometric space the data objects lie in. Therefore intrinsic geodesic regression (red curve) should be considered.

system, which is used to compute residuals of the manifold-valued data points. Regular GP regression is then applied on the residuals and the result is mapped back onto the manifold. The procedure, however, depends heavily on the localization of the problem to a single tangent space, and does not offer an intrinsic probabilistic interpretation.

This paper extends the aforementioned work by allowing kriging that is not localized to a single tangent space, also providing a probabilistic framework offering interpretability. Furthermore, the kriging method in [18] took advantage of the geodesic submanifold regression to initialize a reference coordinate system. Our method, instead, enables the user to take advantage of more general priors, including the use of geodesic submanifold regression.

Other examples of work on regression on manifolds with Euclidean independent variables include kernel-based approaches [2, 3], generalized polynomial regression [8], and regression model that is stochastically developed onto the manifold [12]. Steinke and Hein [23] consider the problem of approximating a function between manifolds via minimizing regularized empirical risk. In this setting, also the independent variables are manifold-valued. The WGP regression proposed in this paper can be extended to this setting, as long as a kernel can be defined on the domain.

**The contribution can be summarized as follows:** We generalize GPs to Riemannian manifolds as wrapped Gaussian processes (WGPs), and provide a novel framework for non-parametric regression with uncertainty estimates using WGP regression. We demonstrate the method in Section 5 on a toy-example on a 2-sphere, in the context of diffusion tensor imaging (DTI), and on a data set of Corpus Callosum shapes. The method is analytically tractable for manifolds with infinite injectivity radius, such as manifolds with non-positive curvature. Otherwise, we suggest the approximation in Remark 2. Computationally, the method is relatively cheap, as the only addition compared to GP regression is the computation of exponential and logarithmic maps.

## 2 Preliminaries

We briefly summarize the mathematical prerequisites needed. First, we recall how GPs are used in non-parametric regression in the Euclidean case, after which we turn to basic concepts in Riemannian geometry and briefly discuss geodesic submanifold regression.

### 2.1 Gaussian process regression

Denote by  $\mathcal{N}(\mu, \Sigma)$  the multivariate Gaussian distribution with mean vector  $\mu \in \mathbb{R}^n$  and covariance matrix  $\Sigma \in \mathbb{R}^{n \times n}$ , and write the probability density function  $p$  as  $p(v) = \mathcal{N}(v \mid \mu, \Sigma)$  for  $v \in \mathbb{R}^n$ .

A *Gaussian process* (GP) [19] is a collection  $f$  of random variables, such that any finite subcollection  $(f(\omega_i))_{i=1}^N$  has a joint Gaussian distribution, where  $\omega_i \in \Omega \subset \mathbb{R}^l$ , and  $\Omega$  is the *index set*. A GP is entirely characterized by the pair

$$m(\omega) = \mathbb{E}[f(\omega)], \quad (2.1)$$

$$k(\omega, \omega') = \mathbb{E}[(f(\omega) - m(\omega))(f(\omega') - m(\omega'))^T], \quad (2.2)$$

where  $m$  and  $k$  are called the *mean function* and *covariance function*, respectively. We denote such a GP by  $f \sim \mathcal{GP}(m, k)$ . It follows from the definition that the covariance function (*kernel*)  $k$  is symmetric and positive semidefinite.

Let  $\mathbf{D} = \{(x_i, y_i) \mid x_i \in \mathbf{x} \subset \mathbb{R}^l, y_i \in \mathbf{y} \subset \mathbb{R}^n\}$  be the training data. The GP predictive distribution for outputs  $\mathbf{y}_*$  at the test inputs  $\mathbf{x}_*$ , given in vector form, is

$$p(\mathbf{y}_* \mid \mathbf{D}, \mathbf{x}_*) = \mathcal{N}(\boldsymbol{\mu}_*, \Sigma_*), \quad (2.3)$$

$$\boldsymbol{\mu}_* = \mathbf{k}_*^T (\mathbf{k} + K_{\text{err}})^{-1} \mathbf{y}, \quad (2.4)$$

$$\Sigma_* = \mathbf{k}_{**} - \mathbf{k}_*^T (\mathbf{k} + K_{\text{err}})^{-1} \mathbf{k}_*, \quad (2.5)$$

where, given a kernel  $k : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  we use the notation  $\mathbf{k} = k(\mathbf{x}, \mathbf{x})$ ,  $\mathbf{k}_* = k(\mathbf{x}, \mathbf{x}_*)$ ,  $\mathbf{k}_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$  and  $K_{\text{err}}$  is the measurement error variance. In the notation above, the function  $k$  is applied elementwise on the vectors  $\mathbf{x}, \mathbf{x}_*$ .

Typically in model selection, the kernel  $k$  is picked from a parametric family  $\{k_\theta \mid \theta \in \Theta\}$  of covariance functions, such as the *radial basis function* (RBF) kernels

$$k_{\sigma^2, \lambda}(x, y) = \sigma^2 \exp\left(-\frac{\|x - y\|^2}{2\lambda}\right), \quad \sigma^2, \lambda > 0, \quad (2.6)$$

choosing the parameters  $(\sigma^2, \lambda)$  so that the *marginal likelihood*  $\mathbb{P}\{\mathbf{y} \mid (\sigma^2, \lambda)\}$  is maximized.

### 2.2 Riemannian geometry

To fix notation, we briefly present the essentials of Riemannian geometry. For a thorough presentation, see [4]. A *Riemannian manifold* is a smooth manifold  $M$  with a smoothly varying inner product  $g_p(\cdot, \cdot)$  (we will often use the notation  $\langle \cdot, \cdot \rangle_p$ ) on the tangent space  $T_p M$  at each  $p \in M$ , called a *Riemannian metric*, inducing

the distance function  $d$  between points on the  $M$ . Each element  $(p, v)$  in the tangent bundle  $TM = \bigcup_{p \in M} (p \times T_p M)$  defines a geodesic  $\gamma$  (a curve locally minimizing distance between two points) on  $M$ , so that  $\gamma(0) = p$  and  $\frac{d}{dt}\gamma(t)|_{t=0} = v$ . The *exponential map*  $\text{Exp} : TM \rightarrow M$  given by  $(p, v) \mapsto \text{Exp}_p(v) = \gamma(1)$ , where  $\gamma$  is the geodesic corresponding to  $(p, v)$ . The exponential map  $\text{Exp}_p$  at  $p$  is a diffeomorphism between a neighborhood  $0 \in U \subset T_p M$  and neighbourhood  $p \in V \subset M$ , which is chosen in a maximal way, so if  $V \subsetneq V'$ , then a diffeomorphism between  $V'$  and a neighborhood in the tangent space cannot be defined anymore. We also call  $V$  the *area of injectivity*.

We can define the inverse map  $\text{Log}_p : V \rightarrow T_p M$ , characterized by  $\text{Exp}_p(\text{Log}_p(p')) = p'$ . Outside of  $V$ , we use  $\text{Log}_p(p')$  to denote a smallest  $v \in T_p M$  chosen in a *measurable*, consistent way. We call the the minimum distance from  $p$  to the boundary of a maximal  $V$  the *injectivity radius* of  $\text{Exp}_p$  and the complement of  $V$  in  $M$  the *cut-locus* at  $p$  denoted by  $\mathcal{C}_p$ . The manifolds with non-positive curvature form an important class of manifolds with infinite injectivity radius, that is, they have an empty cut-locus  $\mathcal{C}_p$  for every  $p \in M$ .

Let  $M_i$  be Riemannian manifolds with metrics  $g_i$ , exponential maps  $\text{Exp}^i$  and logarithmic maps  $\text{Log}^i$  for  $i = 1, 2$ . Then  $M = M_1 \times M_2$  turns into a Riemannian manifold when endowed with the metric  $g = g_1 + g_2$ , which has the component-wise computed exponential map  $\text{Exp}_{(p_1, p_2)}((v_1, v_2)) = (\text{Exp}_{p_1}^1(v_1), \text{Exp}_{p_2}^2(v_2))$ , akin to the logarithmic map  $\text{Log}$  on the product manifold.

### 2.2.1 Probabilistic notions

Let  $X$  be a random point on a Riemannian manifold  $M$ , the set

$$\mathbb{E}[X] := \{p \mid p \in \arg \min_{q \in M} (\mathbb{E}[d(q, X)^2])\}. \quad (2.7)$$

is called the *Fréchet means* of  $X$ . If there is a unique mean  $\bar{p}$ , then by abuse of notation we write  $\mathbb{E}[X] = \bar{p}$ . Given a data set  $\mathbf{p} = \{p_i \in M\}_{i=1}^N$ , an *empirical Fréchet mean* is a minimizer of the quantity

$$\min_{q \in M} \sum_{i=1}^N d(q, p_i)^2. \quad (2.8)$$

The set of empirical Fréchet means is denoted by  $\mathbb{E}[\mathbf{p}]$ .

Given two probability spaces  $(\mathcal{X}_i, \mathcal{S}_i, \nu_i)$  for  $i = 1, 2$  and a measurable map  $F : \mathcal{X}_1 \rightarrow \mathcal{X}_2$ , we say that the measure  $\nu_2$  is the push-forward of the measure  $\nu_1$  with respect to  $F$ , if  $\nu_2(A) = \nu_1(F^{-1}(A))$  for every  $A$  in the sigma-algebra  $\mathcal{S}_2$ . We denote this by  $\nu_2 = F_{\#}\nu_1$ .

For more about intrinsic statistics on manifolds, see [16].

### 2.2.2 Geodesic submanifold regression

Geodesic regression on a Riemannian manifold  $M$  was introduced by Fletcher [5]. It is a generalization of linear regression, that seeks the geodesic parametrized by

$(p, v) \in TM$  that minimizes the quantity

$$E(p, v) = \frac{1}{2} \sum_{i=1}^N d(\text{Exp}_p(t_i v), p_i)^2, \quad (2.9)$$

given the training data  $(t_i, p_i) \in \mathbb{R} \times M$  for  $i = 1, \dots, N$ .

This framework has been generalized to deal with more covariates [11]; assume we are given data  $(x_i, p_i) \in \mathbb{R}^l \times M$  for  $i = 1, \dots, N$ . Then, we want to solve for the submanifold  $\gamma$  parametrized by  $(p, v_1, \dots, v_l)$  that minimizes

$$E(p, v_1, \dots, v_l) = \frac{1}{2} \sum_{i=1}^N d\left(\text{Exp}_p\left(\sum_{j=1}^l x_i(j) v_j\right), p_i\right)^2. \quad (2.10)$$

This is analogous to fitting a hyperplane in the Euclidean case. Another generalization for multiple independent variables was carried out in [18]. Later on in this work, we propose a way to construct priors for the GP regression on manifolds by regressing a geodesic model.

*Tangent space geodesic regression* is a Naïve generalization of linear regression, achieved by linearizing the space by picking  $p \in M$ , transforming the data set  $(x_i, p_i) \in \mathbb{R}^l \times M$  for  $i = 1, \dots, N$  into images of the Riemannian logarithmic map at  $p$ . Then, one can carry out linear regression in the tangent space and map the result onto the manifold using the exponential map, yielding a quick approximation of geodesic submanifold regression.

### 3 Wrapped Gaussian processes

We are now ready to introduce *wrapped Gaussian distributions* (WGDs), computing the conditional distribution of two jointly WGD random points on the manifold. This is an essential part of wrapped Gaussian process (WGP) regression on manifolds introduced in the next chapter, alike in the Euclidean case. In this chapter we also introduce WGDs in a formal way, without studying their properties further.

#### 3.1 Wrapped Gaussian distributions

Wrapped Gaussian distributions (WGDs) originated in directional statistics [13]. There exist multiple different ways of generalizing Gaussian distributions to manifolds. For example, Sommer [21] uses an intrinsic, anisotropic diffusion process for the generalization. Pennec [15], on the other hand, generalizes the Gaussian as the distribution maximizing entropy with a fixed mean and covariance. WGDs rely on linearizing the manifold through a wrapping function, in our case the Riemannian exponential map.

Let  $(M, d)$  be an  $n$ -dimensional Riemannian manifold. We say that a random point  $X$  on  $M$  follows a *wrapped Gaussian distribution* (WGD), if for some  $\mu \in M$  and symmetric positive definite matrix  $K \in \mathbb{R}^{n \times n}$

$$X \sim (\text{Exp}_\mu)_\# (\mathcal{N}(0, K)), \quad (3.1)$$

denoted by  $X \sim \mathcal{N}_M(\mu, K)$ . To sample from this distribution, draw  $v$  from  $\mathcal{N}(0, K)$  and map the sample to the manifold by  $\text{Exp}_\mu(v)$ . Now, define *the basepoint* and *tangent space covariance* of  $X$  as

$$\mu_{\mathcal{N}_M}(X) := \mu, \quad \text{Cov}_{\mathcal{N}_M}(X) := K. \quad (3.2)$$

In the case of infinite injectivity radius  $\mu_{\mathcal{N}_M}(X) \in \mathbb{E}[X]$ , but not in general [14, Prop. 2.11]. The random points  $X_i \sim \mathcal{N}_{M_i}(\mu_i, K_i)$ ,  $i = 1, 2$ , are *jointly WGD*, if the random point  $(X_1, X_2)$  on  $M_1 \times M_2$  is WGD, that is,

$$(X_1, X_2) \sim \mathcal{N}_{M_1 \times M_2} \left( \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} K_1 & K_{12} \\ K_{21} & K_2 \end{pmatrix} \right), \quad (3.3)$$

for some matrix  $K_{12} = K_{21}^T$ .

We now compute the conditional distribution of two jointly WGD random points, which is the core of WGP regression in Section 4.

**Theorem 1.** *Assume  $X_1, X_2$  are jointly WGD as in (3.3), then we have the conditional distribution*

$$X_1 \mid (X_2 = p_2) \sim (\text{Exp}_{\mu_1})_\# \left( \sum_{v \in A} \lambda_v \mathcal{N}(\mu_v, K_v) \right), \quad (3.4)$$

where

$$\begin{aligned} \mu_v &= K_{12} K_2^{-1} v, \\ K_v &= K_1 - K_{12} K_2^{-1} K_{12}^T, \\ \lambda_v &= \frac{\mathcal{N}(v \mid \mathbf{0}, K_2)}{\mathbb{P}\{A\}}, \\ A &= \{v \in T_{\mu_2} M \mid \text{Exp}_{\mu_2}(v) = p_2\}, \\ \mathbb{P}\{A\} &= \sum_{v \in A} \mathcal{N}(v \mid \mathbf{0}, K_2). \end{aligned} \quad (3.5)$$

*Proof.* Pick  $p_1 \in M$ . Let  $B = \text{Exp}_{\mu_1}^{-1}(p_1)$  be the preimage of  $p_1$  in  $T_{\mu_1} M$ , similarly  $A = \text{Exp}_{\mu_2}^{-1}(p_2)$  as above for  $p_2$ , and furthermore  $K$  be the tangent space covariance of  $(X_1, X_2)$  given in (3.3), then

$$\begin{aligned} &\mathbb{P}\{X_1 = p_1 \mid (X_2 = p_2)\} \\ &= \frac{\mathbb{P}\{u \in B, v \in A\}}{\mathbb{P}\{v \in A\}} \\ &= \sum_{v \in A, u \in B} \frac{\mathcal{N}(v \mid \mathbf{0}, K_2)}{\mathbb{P}\{A\}} \frac{\mathcal{N}((u, v) \mid \mathbf{0}, K)}{\mathcal{N}(v \mid \mathbf{0}, K_2)} \\ &= \sum_{v \in A, u \in B} \lambda_v \mathcal{N}(u \mid \mu_v, K_v) \\ &= \mathbb{P}\{Z = p_1\}, \end{aligned} \quad (3.6)$$

where  $Z \sim (\text{Exp}_{\mu_1})_\# \left( \sum_{v \in A} \lambda_v \mathcal{N}(\mu_v, K_v) \right)$ , and  $\mathcal{N}(u \mid \mu_v, K_v)$  is the predictive distribution calculated as in the Euclidean case in (2.3).  $\square$

**Remark 2.** *If the injectivity radius of the exponential map is infinite, then*

$$\begin{aligned} X_1 \mid (X_2 = p_2) \\ \sim (\text{Exp}_{\mu_1})_{\#}(\mathcal{N}(\mu_{\text{Log}_{\mu_2}(p_2)}, K_{\text{Log}_{\mu_2}(p_2)})), \end{aligned} \quad (3.7)$$

following the notation in (3.5). Furthermore, if the probability mass on the area of injectivity of the exponential map is large enough, we can use this expression as a reasonable approximation for the predictive distribution, as the Gaussian mixture distribution in the tangent space can be well approximated by a single Gaussian.

### 3.2 Wrapped Gaussian processes

A collection  $f$  of random points on a manifold  $M$  indexed over a set  $\Omega$  is a *wrapped Gaussian process* (WGP), if every finite subcollection  $(f(\omega_i))_{i=1}^N$  is jointly WGD on  $M^N$ . We define

$$m(\omega) := \mu_{\mathcal{N}_M}(f(\omega)) \quad (3.8)$$

$$k(\omega, \omega') := \text{Cov}_{\mathcal{N}_M}(f(\omega), f(\omega')), \quad (3.9)$$

called the *basepoint function* (BPF) and *tangent space covariance function* (TSCF) of  $f$ , respectively. The restriction we have on  $\Omega$ , is being able to define a kernel on it.

A WGP  $f$  can be viewed as a WGD on the possibly infinite-dimensional product manifold  $M^{|\Omega|}$ . To elaborate, formally one can state

$$f \sim (\text{Exp}_m)_{\#}(\mathcal{GP}(0, k)). \quad (3.10)$$

The difference is, that the tangent space distribution is a GP instead of a GD. The WGP is entirely characterized by the pair  $(m, k)$ , similar to the Euclidean case. Therefore, we introduce the notation  $f \sim \mathcal{GP}_M(m, k)$ .

## 4 Gaussian process inference on manifolds

In the following, we discuss two different methods of GP regression on a Riemannian manifold  $M$  with infinite injectivity radius (or using the approximation in Remark 2), given the noise-free training data

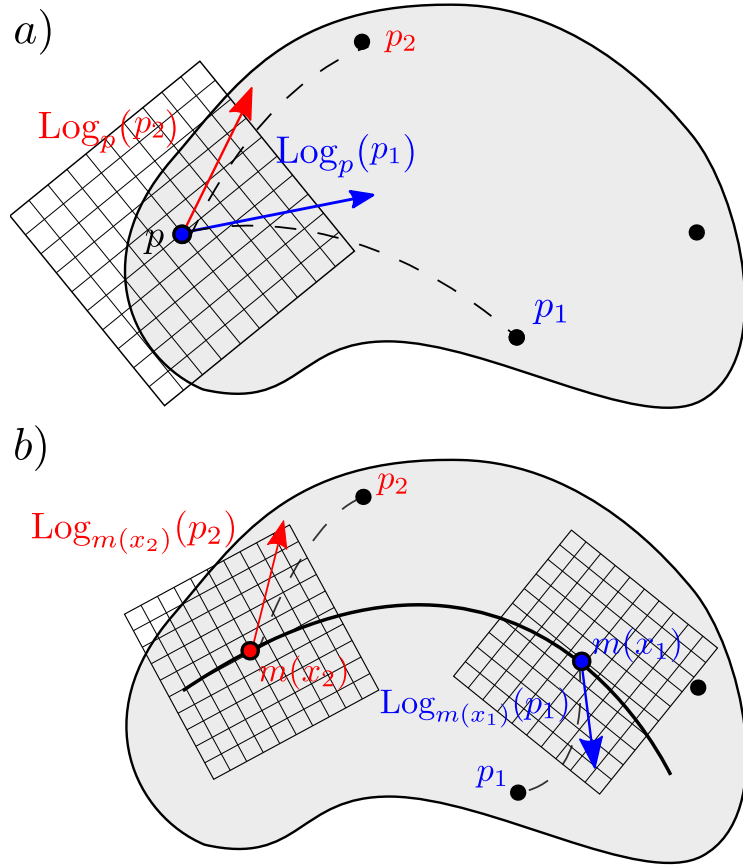
$$\mathbf{D}_M = \{(x_i, p_i) \mid x_i \in \mathbb{R}^l, p_i \in M, i = 1, \dots, N\}. \quad (4.1)$$

For shorthand notation, we denote  $\mathbf{x} = (x_i)_{i=1}^N$  and  $\mathbf{p} = (p_i)_{i=1}^N$ . Additionally,  $\mathbf{x}_*$  is used for the test inputs, and  $\mathbf{p}_*$  for the test outputs. Later, we remark that the first approach is actually a special case of the latter one, see Fig. 2.

### 4.1 Naïve tangent space approach

Choose  $p \in M$  (typically  $p \in \mathbb{E}[\mathbf{p}]$ ), and transform the training data  $\mathbf{D}_M$  into  $\mathbf{D}_{T_p M}$  by

$$\mathbf{D}_{T_p M} = (\mathbf{x}, \mathbf{y}) := \{(x_i, y_i) \mid y_i = \text{Log}_p(p_i)\}, \quad (4.2)$$



**Figure 2:** a) Tangent space GP data transformation. Data point  $p_i$  (in black) is transformed into  $\text{Log}_p(p_i) \in T_p M$ . This can be seen as a special case of WGP regression, with a fixed prior BPF  $m(x) = p$ . In b), the data transformation is visualized with a more general prior BPF  $m$  (black curve).

see Fig. 2 a). As  $\mathbf{D}_{T_p M} \subset \mathbb{R}^l \times T_p M$  now lives in a Euclidean space, fit a GP  $f_{\text{euc}} \sim \mathcal{GP}(m_{\text{euc}}, k_{\text{euc}})$  to the data using GP regression, resulting in the predictive distribution  $\mathbf{y}_* | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_*, \Sigma_*)$ . Then, reversing the previous data transformation, we can map the random vector to a random point  $\mathbf{p}_* | \mathbf{p}$  on the manifold  $M$ , resulting in

$$\mathbf{p}_* | \mathbf{p} = \text{Exp}_p(\mathbf{y}_*) \sim (\text{Exp}_p)_\#(\mathcal{N}(\boldsymbol{\mu}_*, \Sigma_*)). \quad (4.3)$$

## 4.2 Wrapped Gaussian process regression

Now we generalize GP regression inside a probabilistic framework, relying on the results presented in Section 3, by assuming a WGP prior  $f_{\text{prior}} \sim \mathcal{GP}_M(m, k)$ . According to the prior, the joint distribution between the training outputs  $\mathbf{p}$  and test outputs  $\mathbf{p}_*$  at  $\mathbf{x}_*$  is given by

$$\begin{pmatrix} \mathbf{p}_* \\ \mathbf{p} \end{pmatrix} \sim \mathcal{N}_{M_1 \times M_2} \left( \begin{pmatrix} \mathbf{m}_* \\ \mathbf{m} \end{pmatrix}, \begin{pmatrix} \mathbf{k}_{**} & \mathbf{k}_* \\ \mathbf{k}_*^T & \mathbf{k} \end{pmatrix} \right), \quad (4.4)$$

where  $\mathbf{m} = m(\mathbf{x})$ ,  $\mathbf{m}_* = m(\mathbf{x}_*)$ ,  $\mathbf{k} = k(\mathbf{x}, \mathbf{x})$ ,  $\mathbf{k}_* = k(\mathbf{x}_*, \mathbf{x})$ , and  $\mathbf{k}_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$ . Therefore, by Theorem 1 and using the approximation in Remark 2 (if necessary)

$$\begin{aligned} \mathbf{p}_* | \mathbf{p} &\sim (\text{Exp}_{\mathbf{m}_*})_\#(\mathcal{N}(\boldsymbol{\mu}_*, \Sigma_*)), \\ \boldsymbol{\mu}_* &= \mathbf{k}_* \mathbf{k}^{-1} \text{Log}_{\mathbf{m}} \mathbf{p}, \\ \Sigma_* &= \mathbf{k}_{**} - \mathbf{k}_* \mathbf{k}^{-1} \mathbf{k}_*^T. \end{aligned} \quad (4.5)$$

The predictive distribution  $\mathbf{p}_* | \mathbf{p}$  is not necessarily WGD, as  $\boldsymbol{\mu}_*$  might be non-zero. The distribution can be sampled from, but computing exactly quantities such as  $\mathbb{E}[\mathbf{p}_* | \mathbf{p}]$  is not trivial. As in [7, Sect. 3.1.1], the distribution can be approximated via Riemannian unscented transform or by using a WGD with the basepoint at  $\text{Exp}_{\mathbf{m}_*}(\boldsymbol{\mu}_*)$  and parallel transporting the tangent space covariance to this point along the geodesic  $\gamma(t) = \text{Exp}_{\mathbf{m}_*}(t\boldsymbol{\mu}_*)$ .

**Remark 3.**  $\text{Exp}_{\mathbf{m}_*}(\boldsymbol{\mu}_*)$  is not necessarily a Fréchet mean of  $\mathbf{p}_* | \mathbf{p}$ . However, it is the maximum a posteriori (MAP) estimate. For this reason, we will use  $\text{Exp}_{\mathbf{m}_*}(\boldsymbol{\mu}_*)$  as a point prediction in Section 5.

### 4.2.1 Choosing a prior

The prior WGP  $f_{\text{prior}} \sim \mathcal{GP}_M(m, k)$  indexed over  $\Omega$  is chosen by picking a kernel  $k$  on  $\Omega$  to be the TSCF, and picking a BPF  $m$  so that  $p$  and  $m(x_i)$  live in the same connected component of  $M$  for every data-point  $(x_i, p_i)$ .

In Section 5, two kinds of prior BPFs are used. The first BPF  $m_1$  is a generalization of a centered GP, given by  $m_1(\omega) = \bar{p}$ , for all  $x \in \Omega$  and a  $\bar{p} \in \mathbb{E}[\mathbf{p}]$ . The second kind  $m_2$ , uses a previous regression (such as geodesic submanifold regression)  $\gamma$  on the dataset  $\mathbf{D}_M$ . That is,  $m_2(\omega) = \gamma(\omega)$  for all  $\omega \in \Omega$ . For computational reasons, we only consider TSCFs that assume each tangent space coordinate independent, resulting in the *diagonal RBF* kernel

$$k(\mathbf{x}, \mathbf{x}') = \text{diag}(k_1(\mathbf{x}, \mathbf{x}'), k_2(\mathbf{x}, \mathbf{x}'), \dots, k_n(\mathbf{x}, \mathbf{x}')), \quad (4.6)$$

where each  $k_i$  are chosen to be RBF kernels,  $\text{diag}(A, B)$  is a block-diagonal matrix with blocks  $A$  and  $B$ ,  $\mathbf{x}, \mathbf{x}' \in \Omega$ , and  $n$  is the dimension of  $M$ . The diagonal RBF yields uncertainty estimates, but not a generative model, as this would need covariance between coordinates.

**Optimizing hyperparameters.** We choose the TSCF from a parametric family of kernels  $\{k_\theta\}_{\theta \in \Theta}$  maximizing the *marginal likelihood*, as in the Euclidean case. In the setting of WGP, the marginal likelihood becomes

$$\mathbb{P}\{\mathbf{p} \mid \theta\} = \sum_{v \in \text{Exp}_m^{-1}(\mathbf{p})} \mathcal{N}(v \mid \mathbf{0}, K_\theta), \quad (4.7)$$

where  $K_\theta = k_\theta(\mathbf{x}, \mathbf{x})$ . To improve the approximation discussed in Remark 2, we propose to maximize the quantity

$$\mathbb{P}\{\mathbf{p} \mid \theta\} \approx \mathcal{N}(\text{Log}_m(\mathbf{y}) \mid \mathbf{0}, K_\theta), \quad (4.8)$$

as maximizing this quantity increases the probability mass given by the prior distribution to the area of injectivity. The diagonal RBF kernel (Eq. (4.6)) can be optimized by choosing each  $k_i$  to maximize the marginal likelihood of the respective tangent space coordinate independently. That is,  $k_i$  is chosen to maximize the marginal likelihood of the data set  $\{(x_j, \pi_i(\text{Log}_{m(x_j)}(p_j)))\}_{j=1}^N$ , where  $\pi_i$  is the projection onto the  $i$ th component.

A part of engineering the kernel is to pick a frame for the manifold. A frame is a smooth map  $\rho : M \rightarrow \mathbb{R}^{n \times n}$ , so that the columns of  $\rho(p)$  form an orthonormal basis for  $T_p M$ . This way, there is a relation between tangent vectors in different tangent spaces, and so the covariance becomes meaningful.

The WGP regression process is summarized in Alg. 4.

**Algorithm 4** (WGP regression.). *The following describes step-by-step how to carry out WGP regression.*

**Input** Manifold-valued training data  $\mathbf{D}_M = \{(x_i, p_i)\}_{i=1}^n$ .

**Output** Predictive distribution for  $\mathbf{p}_* \mid \mathbf{p}$  at  $\mathbf{x}_*$ .

- i. Choose a prior BPF  $m$ .
- ii. Transform  $\mathbf{D}_{T_m M} \leftarrow \{(x_i, \text{Log}_{m(x_i)}(p_i))\}_{i=1}^N$ .
- iii. Choose a prior TSCF  $k$  from a parametric family by optimizing the hyperparameters.
- iv. Using GP prior  $\mathcal{GP}(0, k)$ , carry out Euclidean GP regression for the transformed data  $\mathbf{D}_{T_m M}$ , yielding the mean and covariance  $(\boldsymbol{\mu}_*, \Sigma_*)$ .
- vi. End with the predictive distribution  $\mathbf{p}_* \mid \mathbf{p} \sim (\text{Exp}_{\mathbf{m}_*})_\#(\mathcal{N}(\boldsymbol{\mu}_*, \Sigma_*))$

### 4.2.2 Observations with noise

A difficulty arises, when introducing a noise model on our observations. In the Euclidean case, a popular noise model on the observations  $(x_i, p_i)$  is given by  $p_i = f(x_i) + \epsilon$ , where  $f$  is the function we approximate and  $\epsilon \sim \mathcal{N}(0, K_{\text{err}})$  is the noise term. In [5], this model is generalized to the manifold setting implicitly as

$$p_i = \text{Exp}_{f(x_i)}(\epsilon), \quad (4.9)$$

which is also supported by the central limit theorem provided in [10]. However, this makes the WGP analytically intractable. To allow computations, we propose the error model  $\text{Log}_{m(x_i)}(p_i) = \text{Log}_{m(x_i)}(f(x_i)) + \epsilon$ , that is, the error lives in the tangent space of the prior mean at  $x_i$ . This can be viewed as a first order approximation of (4.9) around  $m(x_i)$ . Introduction of this error changes the regression procedure only slightly; the joint distribution of  $\mathbf{p}$  and  $\mathbf{p}_*$  changes into

$$\begin{pmatrix} \mathbf{p}_* \\ \mathbf{p} \end{pmatrix} \sim \mathcal{N}_{M_1 \times M_2} \left( \begin{pmatrix} \mathbf{m}_* \\ \mathbf{m} \end{pmatrix}, \begin{pmatrix} \mathbf{k}_{**} & \mathbf{k}_* \\ \mathbf{k}_*^T & \mathbf{k} + K_{\text{err}} \end{pmatrix} \right). \quad (4.10)$$

Rest of the computations are then carried out similarly, with the replacement of  $\mathbf{k}$  with  $\mathbf{k} + K_{\text{err}}$  everywhere.

## 5 Experiments

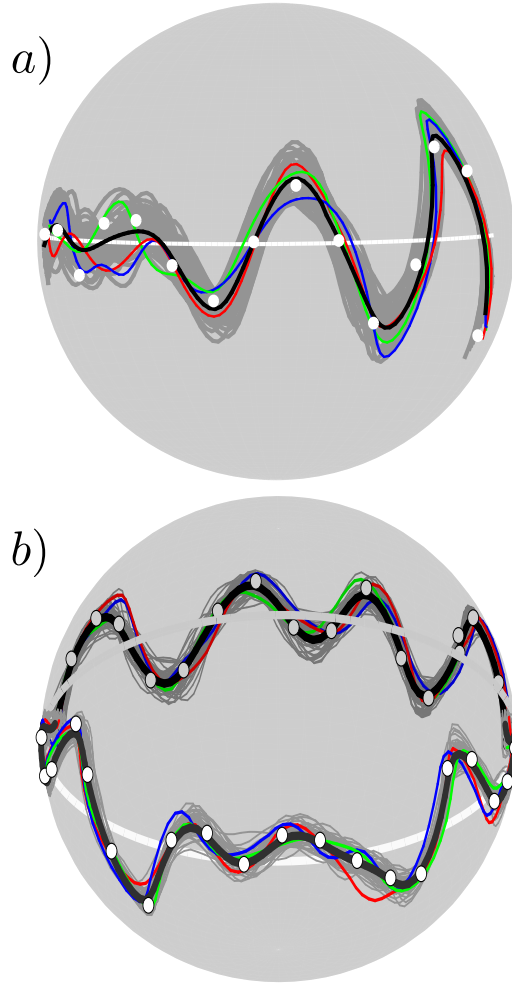
We demonstrate WGP regression in three different contexts; on a manufactured data set on the 2-sphere, on data from a DTI slice, and finally on Corpus Callosum shapes.

### 5.1 Data on 2-sphere

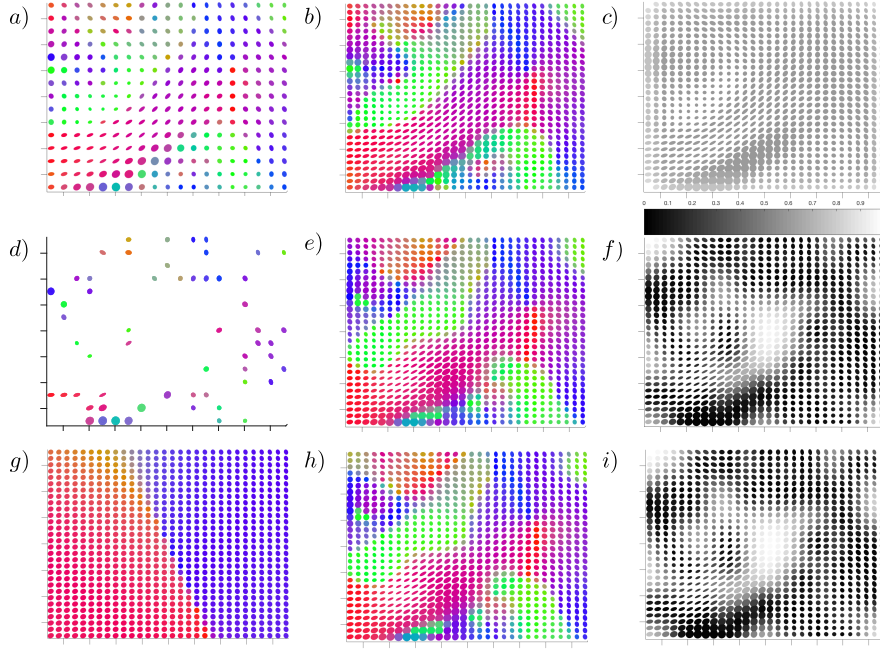
For an easy visualization, we show results of WGP regression on two different manufactured datasets on a 2-sphere labeled by univariate independent variables, see Fig. 3. Note, how the curves sampled from the predictive distribution given by the WGP regression lie on the manifold, respecting the intrinsic geometry. The metric used for the 2-sphere is the one with Euclidean inner product as the Riemannian metric at each tangent space. As this Riemannian manifold does not have infinite injectivity radius, we use the approximation explained in Remark 2.

In Fig. 3 a) geodesic regression  $\gamma$  is computed to be used as the prior BPF, and a diagonal RBF kernel (described in Eq. (4.6)) with optimized hyperparameters is chosen as the prior TSCF.

In Fig. 3 b), we view the independent variables to live on the 1-sphere  $S^1$ , so that samples from the predictive distribution are loops. To compute the geodesic regression, we treat the independent variables as elements of  $\mathbb{R}$ , not  $S^1$ . The prior TSCF is chosen to be the diagonal RBF on  $S^1$ , that is, each component is of the form  $k(t, t') = \sigma^2 \exp(-\frac{d(t, t')^2}{2\lambda})$ , where  $t, t' \in S^1$  and  $d(t, t') = \min(\|t - t'\|, 2\pi - \|t - t'\|)$ . This kernel is not positive-definite for all  $\sigma^2, \lambda$ , as is implied in [7]. However, for the data set in Fig. 3, the pairs  $(\sigma^2, \lambda)$  maximizing the marginal likelihood were feasible pairs for each coordinate. For more on periodic kernels, see [20].



**Figure 3:** Depicted in a) is WGP regression using a prior BPF given by geodesic regression (white curve) on a manufactured data set (white dots) on a 2-sphere  $S^2$ . The predictive distribution is visualized using the MAP estimate (thick black line, see Remark 3) and multiple samples from the distribution (in gray) with three samples emphasized (in red, green and blue). In b), a data set going all around the sphere was used, and the index set of the prior WGP considered as  $S^1$ , so that the sampled paths are closed loops.



**Figure 4:** WGP regression on DTI data. Colors depict the direction of the largest eigenvector of the respective tensor. **First row:** a) The original DTI data set, b) MAP estimate of the predictive distribution of WGP regression on the original data set, c) relative error in b) (white = maximum error). **Second row:** d) 20% of data points in a) randomly subsampled to be the training set, e) MAP estimate of the predictive distribution, f) relative error in e). **Third row:** g) Geodesic submanifold regression on d), h) MAP estimate of the predictive distribution using g) as prior BPF, i) relative error in h).

## 5.2 Diffusion tensor imaging data

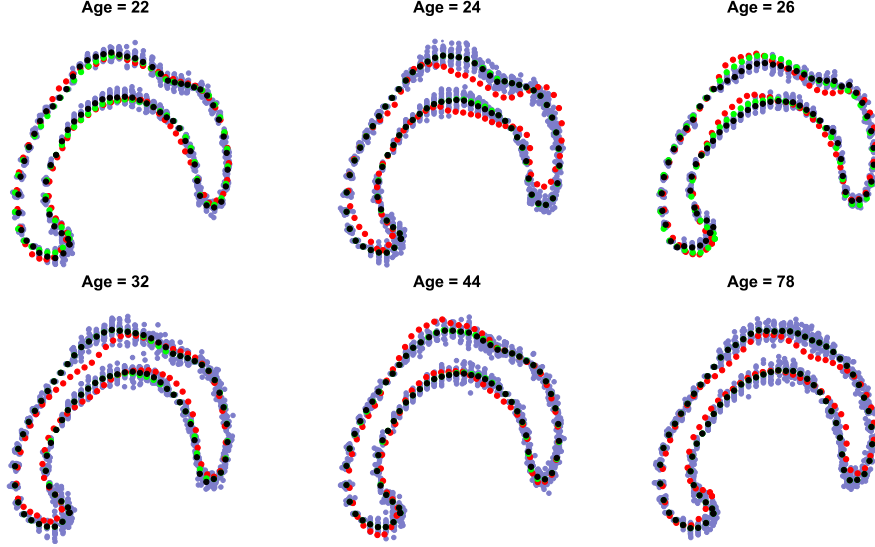
Now, we consider a DTI coronal slice of a HCP data set [6, 22, 24], which lives on the manifold  $\mathbb{R}^2 \times \text{PD}(3)$ , where  $\text{PD}(n)$  is the set of positive definite matrices of dimension  $n$ .  $\text{PD}(n)$  forms a Riemannian manifold when endowed with the affine-invariant metric [17]. The manifold has non-positive curvature, thus the WGP regression is exact. The original data set consists of  $15 \times 19$  tensors (elements of  $\text{PD}(3)$ ) sampled isotropically, see Fig. 4 a). For another experiment, we subsample the original data set picking only a fifth of the original data points randomly, see Fig. 4 d). Then, for both data sets the interpolation is carried out on a  $30 \times 30$  grid. As a measure of uncertainty of the result, we calculate the sum of variances of each tangent space coordinate at the interpolated points. In Figs. 4 b), e) are shown the corresponding MAP estimates of the predictive distribution (see Remark 3). Empirical Fréchet means were used as the prior BPFs and diagonal RBFs with optimized hyperparameters as the prior TSCFs.

In the next experiment, the prior BPF is chosen to be the geodesic submanifold regression on the subsampled dataset of the DTI coronal slice, see Fig. 4 g), h). The MAP estimates in Figs. 4 e) and h) do not differ vastly, although different BPFs were used. They yield a different result in the upper-left corner area, where the subsampled dataset is not dense, hence the regressed result approaches the prior BPF. In the middle, where we also lack information, the resulting tensor fields look similar. The error structures are very similar, seen in Figs. 4 f), i). Both MAP estimates grasp the general trend of the original tensor field, as can be seen by comparing Figs. 4 a), b), e) and h).

## 5.3 Corpus Callosum data

Next, we turn to a dataset of landmark representations of Corpus Callosum (CC) shapes [5]. A landmark representation is a set of  $k$  points in  $\mathbb{R}^2$ , so that length, translation and rotation factors have been quotiented out, resulting in a point in the *Kendall's shape space* [9], that is equivalent to the complex projective plane  $\mathbb{C}P^{k-1}$ . The dataset consists of 65 shapes, of which we pick randomly 6 to be the test set, the rest are used for training.

Results are presented in Fig. 5. A tangent space geodesic regression is used as the prior BPF, and a diagonal RBF kernel with optimized hyperparameters is used as the prior TSCF. As the CC shapes vary considerably even in the same age group, the WGP predictive mean does not yield notable gains on the tangent space geodesic regression used as prior BPF. However, it provides uncertainty estimates of the shape. Notably, the results imply that aging brings about wider variation in the upper-right part of the CC.



**Figure 5:** WGP regression applied to a population of Corpus Callosum shapes labeled by age. Red shapes are data points from the test set, not used for training. In black, the MAP estimates of the predictive distributions, in green values of the prior BPF at corresponding ages. Drawn in blue are 20 samples from the predictive distribution.

## 6 Conclusion and discussion

This paper introduced WGP regression on Riemannian manifolds in a novel Bayesian inference framework. The method relies on WGP, defined using WGDs. The conditional distribution of two jointly WGD random points was computed to be used in WGP regression. The method was demonstrated in three cases; on toy data lying on 2-sphere, tensor data originating from DTI and on a set of Corpus Callosum shapes. The results of the experiments imply that WGP regression can be used effectively on Riemannian manifolds, providing meaningful uncertainty estimates.

This being the first step, there are still open questions; how to engineer prior distributions effectively, and how to treat the predictive distribution. The predictive distribution admits an explicit expression, but the prediction is not a WGP anymore. Therefore, we do not have same closure properties of the family of distributions as in the Euclidean case. This leaves open the question, whether one should consider other generalizations of GDs than the wrapped one when carrying out GP regression on manifolds.

We suggested an approximation in Remark 2, not quantifying how reliable it is in the case of non-infinite injectivity radius. In practice the approximation seems plausible (see Fig. 3), but should be studied in more detail. Furthermore, it is of interest, in which cases the computations can be carried out analytically, when the injectivity radius is non-infinite.

The central limit theorem presented in [10] suggests to use WGD distributed error terms, but this poses the difficulty of incorporating the noise term into the prior, when the noise term might live in a different tangent space. The workaround used in this paper was to approximate this error term linearly in the tangent space of the prior BPF, however, other models should also be considered.

Finally, GP regression could be generalized to a broader family of spaces than Riemannian manifolds. In WGP regression, the key is having a wrapping function from a model vector space onto the manifold. For example, another context where such structure appears, is the weak Riemannian structure of the space of probability measures under the Wasserstein metric [1].

## 7 Acknowledgements

This research was supported by Centre for Stochastic Geometry and Advanced Bioimaging, funded by a grant from the Villum Foundation. Data were provided [in part] by the Human Connectome Project, WU-Minn Consortium (Principal Investigators: David Van Essen and Kamil Ugurbil; 1U54MH091657) funded by the 16 NIH Institutes and Centers that support the NIH Blueprint for Neuroscience Research; and by the McDonnell Center for Systems Neuroscience at Washington University.

## References

- [1] L. Ambrosio and N. Gigli. A user’s guide to optimal transport. In *Modelling and optimisation of flows on networks*, pages 1–155. Springer, 2013.
- [2] M. Banerjee, R. Chakraborty, E. Ofori, D. Vaillancourt, and B. C. Vemuri. Nonlinear regression on Riemannian manifolds and its applications to Neuro-image analysis. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 719–727. Springer, 2015.
- [3] B. C. Davis, P. T. Fletcher, E. Bullitt, and S. Joshi. Population shape regression from random design data. *International journal of computer vision*, 90(2):255–266, 2010.
- [4] M. P. Do Carmo and J. Flaherty Francis. *Riemannian geometry*, volume 115. Birkhäuser Boston, 1992.
- [5] P. T. Fletcher. Geodesic regression and the theory of least squares on Riemannian manifolds. *International journal of computer vision*, 105(2):171–185, 2013.
- [6] M. F. Glasser, S. N. Sotiropoulos, J. A. Wilson, T. S. Coalson, B. Fischl, J. L. Andersson, J. Xu, S. Jbabdi, M. Webster, J. R. Polimeni, et al. The minimal preprocessing pipelines for the Human Connectome Project. *Neuroimage*, 80:105–124, 2013.
- [7] S. Hauberg, F. Lauze, and K. S. Pedersen. Unscented Kalman filtering on Riemannian manifolds. *Journal of mathematical imaging and vision*, 46(1):103–120, 2013.
- [8] J. Hinkle, P. Muralidharan, P. T. Fletcher, and S. Joshi. Polynomial regression on Riemannian manifolds. In *European Conference on Computer Vision*, pages 1–14. Springer, 2012.
- [9] D. G. Kendall. Shape manifolds, procrustean metrics, and complex projective spaces. *Bulletin of the London Mathematical Society*, 16(2):81–121, 1984.
- [10] W. S. Kendall, H. Le, et al. Limit theorems for empirical Fréchet means of independent and non-identically distributed manifold-valued random variables. *Brazilian Journal of Probability and Statistics*, 25(3):323–352, 2011.

- [11] H. J. Kim, N. Adluru, M. D. Collins, M. K. Chung, B. B. Bendlin, S. C. Johnson, R. J. Davidson, and V. Singh. Multivariate general linear models (MGLM) on Riemannian manifolds with applications to statistical analysis of diffusion weighted images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2705–2712, 2014.
- [12] L. Kühnel and S. Sommer. Stochastic development regression on non-linear manifolds. In *International Conference on Information Processing in Medical Imaging*, pages 53–64. Springer, 2017.
- [13] K. V. Mardia and P. E. Jupp. *Directional statistics*, volume 494. John Wiley & Sons, 2009.
- [14] J. Oller, J. M. Corcuera, et al. Intrinsic analysis of statistical estimation. *The Annals of Statistics*, 23(5):1562–1581, 1995.
- [15] X. Pennec. Probabilities and statistics on Riemannian manifolds: Basic tools for geometric measurements. In *NSIP*, pages 194–198, 1999.
- [16] X. Pennec. Intrinsic statistics on Riemannian manifolds: Basic tools for geometric measurements. *Journal of Mathematical Imaging and Vision*, 25(1):127–154, 2006.
- [17] X. Pennec, P. Fillard, and N. Ayache. A Riemannian framework for tensor computing. *International Journal of Computer Vision*, 66(1):41–66, 2006.
- [18] D. Pigoli, A. Menafoglio, and P. Secchi. Kriging prediction for manifold-valued random fields. *Journal of Multivariate Analysis*, 145:117–131, 2016.
- [19] C. E. Rasmussen and C. K. Williams. *Gaussian processes for machine learning*, volume 1. MIT press Cambridge, 2006.
- [20] A. Solin and S. Särkkä. Explicit link between periodic covariance functions and state space models. In *Artificial Intelligence and Statistics*, pages 904–912, 2014.
- [21] S. Sommer. Anisotropic distributions on manifolds: template estimation and most probable paths. In *International Conference on Information Processing in Medical Imaging*, pages 193–204. Springer, 2015.
- [22] S. Sotiropoulos, S. Moeller, S. Jbabdi, J. Xu, J. Andersson, E. Auerbach, E. Yacoub, D. Feinberg, K. Setsompop, L. Wald, et al. Effects of image reconstruction on fiber orientation mapping from multichannel diffusion MRI: reducing the noise floor using SENSE. *Magnetic resonance in medicine*, 70(6):1682–1689, 2013.
- [23] F. Steinke and M. Hein. Non-parametric regression between manifolds. In *Advances in Neural Information Processing Systems*, pages 1561–1568, 2009.
- [24] D. C. Van Essen, S. M. Smith, D. M. Barch, T. E. Behrens, E. Yacoub, K. Ugurbil, W.-M. H. Consortium, et al. The WU-Minn human connectome project: an overview. *Neuroimage*, 80:62–79, 2013.