

SLEIPNIR: Deterministic and Provably Accurate Feature Expansion for Gaussian Process Regression with Derivatives

Emmanouil Angelis*

ANGELISE@ETHZ.CH

*Learning and Adaptive Systems Group
ETH Zürich*

Philippe Wenk*

WENKPH@ETHZ.CH

*Learning and Adaptive Systems Group
ETH Zürich and Max Planck ETH Center for Learning Systems*

Bernhard Schölkopf

BERNHARD.SCHOELKOPF@TUEBINGEN.MPG.DE

*Empirical Inference Group
Max Planck Institute for Intelligent Systems, Tübingen*

Stefan Bauer

STEFAN.BAUER@TUEBINGEN.MPG.DE

*Empirical Inference Group
Max Planck Institute for Intelligent Systems, Tübingen*

Andreas Krause

KRAUSEA@ETHZ.CH

*Learning and Adaptive Systems Group
ETH Zürich*

Abstract

Gaussian processes are an important regression tool with excellent analytic properties which allow for direct integration of derivative observations. However, vanilla GP methods scale cubically in the amount of observations. In this work, we propose a novel approach for scaling GP regression with derivatives based on quadrature Fourier features. We then prove deterministic, non-asymptotic and exponentially fast decaying error bounds which apply for both the approximated kernel as well as the approximated posterior. To furthermore illustrate the practical applicability of our method, we then apply it to ODIN, a recently developed algorithm for ODE parameter inference. In an extensive experiments section, all results are empirically validated, demonstrating the speed, accuracy, and practical applicability of this approach.

1. Introduction

Gaussian process (GP) regression (Rasmussen and Williams, 2006) is an important machine learning model with many desirable properties. Due to their inherently Bayesian nature, GPs naturally provide uncertainty estimates which are of crucial importance in Bayesian optimization (Mockus, 2012) or active learning (Seo et al., 2000). Furthermore, since a derivative of a GP is again a GP, there is a natural extension to derivative observations (Solak et al., 2003), which can be leveraged in the context of probabilistic numerics (Hennig et al., 2015) or Gaussian-process-based gradient matching (Calderhead et al., 2009; Wenk et al., 2020). Despite these appealing properties, a clear drawback of classic GP regression is the fact that they scale cubic in the amount of observation points.

*. The first two authors contributed equally.

Scaling standard GP regression For standard, unconstrained GP regression, there exist several ideas on how to tackle this problem. One family of approaches focuses on summarizing the data set with a fixed amount of pseudo-observations, the so-called inducing points (Quiñero-Candela and Rasmussen, 2005; Snelson and Ghahramani, 2006; Titsias, 2009). Complementary to these methods, special structure in either the kernel function (Wilson et al., 2014) or the input data (Cunningham et al., 2008) can be exploited to speed up the necessary matrix-vector multiplications needed to perform GP regression. Wilson and Nickisch (2015) combine these two ideas, creating an efficient approximation scheme linear in the amount of observations. Independently, Sarkka et al. (2013) propose a SDE-based reformulation and connect GP regression to Kalman filtering. Finally, there is a family of approaches approximating the kernel function via a finite dimensional scalar product of feature vectors. These features have been obtained using the Nyström method (Williams and Seeger, 2001), MC samples (Rahimi and Recht, 2008), sparse spectrum approximations (Lazaro-Gredilla et al., 2010), variational optimization (Hensman et al., 2017) or a quadrature scheme (Mutny and Krause, 2018). Furthermore, Solin and Särkkä develop a deterministic feature expansion based on the Fourier transform of the Laplace operator. While they are able to provide deterministic, non-asymptotic error bounds, their error decays linearly with the size of the domain of the operator expansion L , which cannot grow faster than the amount of features. This essentially means that their approximation decays at best linearly.

Scaling GP regression with derivatives In the context of Gaussian process regression with derivatives though, scalable methods seem to have received little attention (Eriksson et al., 2018), despite their practical relevance. While some approaches like inducing points could be applied by just using the same inducing points, it is not obvious how such changes would affect the quality of the resulting estimates. Similarly, an empirical extension of the approach of Solin and Särkkä to the derivative case is presented by Solin et al. (2018), without quantifying the errors of the approximation scheme. For random Fourier features (RFF), Szabó and Sriperumbudur (2018) present a rigorous theoretical analysis. However, due to the probabilistic nature of RFF, it is not possible to provide deterministic, non-asymptotic guarantees for any given data set of fixed size.

To obtain deterministic (i.e. hold with probability 1) and non-asymptotic (i.e. for data sets with a fixed amount of observations) error bounds, we thus turn to quadratic Fourier features. Mutny and Krause (2018) recently derived exponentially fast decaying bounds for standard GP regression. Building on their work, we derive approximations and bounds for kernel derivatives. As we will demonstrate, these bounds can be used directly to quantify the absolute error of the predictive posterior mean and covariance of a GP with derivative information, leading to deterministic, non-asymptotic error bounds that can guide the choice of the amount of features needed to obtain a desired accuracy.

However, as we shall demonstrate, our bounds are so powerful that they can easily be extended to more applications. We will demonstrate this using ODE-informed regression (ODIN) as an example. In ODE-Informed regression, Gaussian processes are used to infer the parameters of a system of ODEs governing the dynamics of a time-continuous system that is observed at discrete time points under additive Gaussian noise. The idea to use GPs for this task goes back to the pioneering work of Calderhead et al. (2009), whose theoretical models were later refined by Dondelinger et al. (2013) and Wenk et al. (2018). As shown e.g. by Gorbach et al. (2017) and Abbati et al. (2019), such GP based inference schemes scale almost linearly in the dimension of the ODE. Since applications of these algorithms involve deducing novel scientific knowledge (Dony et al., 2019; Macdonald and

Husmeier, 2015), any approximation scheme would require a rigorous error analysis. Fortunately, our algorithm leads to deterministic, finite sample error bounds for the risk it optimizes, again guiding the choice of the complexity of the approximation for a desired accuracy level.

Contributions: In summary, we

- extend the QFF framework to Gaussian process regression with derivative information, naming the extension SLEIPNIR,
- derive, prove and empirically validate deterministic, finite-sample guarantees for the accuracy of said approximation,
- demonstrate how these theoretical insights can be used to control the approximation error in GPR with derivatives,
- expand these insights to consistently reduce the cubic run time of ODE informed regression (Wenk et al., 2020),
- verify the theoretical results empirically on four different systems with both locally linear and more involved, nonlinear dynamics, demonstrating a significant reduction in computational complexity without noticeable loss of accuracy.

All code needed to recreate our results can be found at <https://github.com/sdi1100041/SLEIPNIR> to facilitate future research and reproducibility.

2. Background

In this section, we provide a high-level overview of the background for our work. For an in-depth introduction, see Rasmussen and Williams (2006) (GPs), Rahimi and Recht (2008) (RFF) and Mutny and Krause (2018) (QFF).

2.1 Feature Expansions for normal GPR

Unfortunately, any reasoning in such a model requires the calculation and inversion of the covariance matrices C_ϕ and A . Without any tricks or approximations, this will scale as $\mathcal{O}(N^3)$. Random Fourier features (RFF) introduced by Rahimi and Recht (2008) and quadrature Fourier features (QFF) introduced by Mutny and Krause (2018) are two approximations that can be used to reduce this complexity for standard GP regression. Both approaches are based on the following observation: Any kernel represents a scalar product in a potentially infinite dimensional feature space, but if we allow for a small error, this could be approximated by a finite-dimensional feature vector.

For readability, we will introduce all concepts using a kernel with scalar inputs $k(t_i, t_j)$. However, it should be noted that this is by no means a necessary condition, and the concept generalizes nicely to higher dimensional inputs.

For a stationary, scalar kernel, Bochner’s theorem (see e.g. Rudin, 1976) guarantees the existence of a density $p(\omega)$ such that we can write

$$k(|t_i - t_j|) = \int_{-\infty}^{\infty} p(\omega) \begin{pmatrix} \cos(\omega t_i) \\ \sin(\omega t_i) \end{pmatrix}^T \begin{pmatrix} \cos(\omega t_j) \\ \sin(\omega t_j) \end{pmatrix} d\omega. \quad (1)$$

This equation can be interpreted as a scalar product of infinitely many features given by $\cos(\omega t_i)$ and $\sin(\omega t_i)$. But how do we find the most important features such that the kernel k can be reasonably well approximated by a finite feature vector

$$k(|t_i - t_j|) \approx \phi(t_i)^T \phi(t_j)? \quad (2)$$

To obtain random Fourier features, Rahimi and Recht (2008) propose two different Monte Carlo sampling schemes. The first approximation just obtains MC samples from $p(\omega)$ and for each sample adds $[\sin(\omega t_i) \cos(\omega t_i)]$ to the feature vector $\phi(t_i)$. In the second scheme, they observe that

$$\begin{pmatrix} \cos(\omega t_i) \\ \sin(\omega t_i) \end{pmatrix}^T \begin{pmatrix} \cos(\omega t_j) \\ \sin(\omega t_j) \end{pmatrix} = \mathbb{E}_b \{ \sqrt{2} \cos(\omega t_i + b) \}, \quad (3)$$

where $b \sim \text{Unif}([0, 2\pi])$. Thus, one can obtain a second approximation by sampling concurrently b and ω and for each sample adding $\cos(\omega t_i + b)$ to the feature vector $\phi(t_i)$. Accentuating the presence of the bias term b , we will refer to this approximation in the following as RFF-B, while we refer to the first one as RFF.

One main drawback of RFFs however is the fact that the intermediate sampling step makes it impossible to obtain any deterministic bounds on their approximation quality. Mutny and Krause (2018) thus propose to approximate the integral of Equation (1) using Hermitian quadrature. While RFF and RFF-B use sampling to determine the locations of the ω_i , in the context of Hermitian quadrature, these locations are fully determined by the functional form of the kernel. Thus, quadrature Fourier features (QFF) are deterministic in nature, leading to a deterministic, theoretical analysis. Furthermore, as a numerical integration scheme, they are especially efficient in low dimensions. As we shall see both in our theoretical and empirical analysis, similar properties can be achieved if we want to include derivative observations.

2.2 Gaussian Process Regression with Derivatives

Assume there exists a scalar-valued function $x(t)$. At N distinct time points $\mathbf{t} = [t_0, \dots, t_{N-1}]$, we obtain noisy observations of the function itself and its derivatives, assuming Gaussian noise with standard deviations σ and $\sqrt{\gamma}$. Following Wenk et al. (2020), these observations are represented as the vectors \mathbf{y} and \mathbf{F} . Using GP regression, we aim to find estimates for $\mathbf{x} = [x(t_0), \dots, x(t_{N-1})]$ given \mathbf{y} and \mathbf{F} . In GP regression, we assume that $x(t)$ is drawn from a Gaussian process with a kernel $k(t_i, t_j)$ that is parameterized by hyper-parameters ϕ . For fixed ϕ , this leads to tractable Gaussian priors over the function \mathbf{x} and its derivatives $\dot{\mathbf{x}}$. These priors can be combined with the observation models to obtain the generative model summarized in Figure 1.



Figure 1: Generative model for Gaussian process regression with derivative observations.

Here, the matrices \mathbf{C}_ϕ , \mathbf{D} and \mathbf{A} are fully determined by the kernel k and the choice of \mathbf{t} . Since all relevant probability densities are Gaussian, the posteriors can be calculated analytically, as they

will be Gaussian as well. In the rest of this paper, we will focus on the predictive posterior mean and variance for both state and variance, i.e.

$$p(x(\tau)|\mathbf{y}, \mathbf{F}, \phi, \gamma) = \mathcal{N}(x(\tau)|\mu(\tau), \Sigma(\tau)) \quad (4)$$

and

$$p(x'(\tau)|\mathbf{y}, \mathbf{F}, \phi, \gamma) = \mathcal{N}(x'(\tau)|\mu'(\tau), \Sigma'(\tau)). \quad (5)$$

$p(\mathbf{x}(\tau)|\mathbf{y}, \mathbf{F}, \phi, \gamma)$ These quantities are of particular importance for any task involving smoothing or interpolation. As we shall demonstrate, our scheme allows for exponentially fast decaying error bounds for these quantities as well. For more details regarding notation, please refer to Section C.1.

2.3 Parameter Inference for ODEs with GPs

One important application of this regression technique is parameter inference of differential equations. To simplify notation, we briefly recap the theory for one dimensional systems of ODEs, noting that there is no significant difference to the multidimensional case.

Assume that we are given the noisy observations \mathbf{y} of a dynamical system with known parametric form $\dot{x} = f(x, \theta)$. The goal is to infer the unknown parameters θ . In ODIN, this equation is now used as a constraint for the probabilistic model in Figure 1. Parameter and states are then found by minimizing the following objective, where \mathbf{F} has already been substituted by the constraints:

$$\mathcal{R}(\mathbf{x}, \mathbf{F}, \mathbf{y}) = \mathbf{x}^T \mathbf{C}_\phi^{-1} \mathbf{x} \quad (6)$$

$$+ (\mathbf{x} - \mathbf{y})^T \sigma^{-2} (\mathbf{x} - \mathbf{y}) \quad (7)$$

$$+ (\mathbf{F} - \mathbf{D}\mathbf{x})^T (\mathbf{A} + \gamma \mathbf{I})^{-1} (\mathbf{F} - \mathbf{D}\mathbf{x}). \quad (8)$$

To infer the hyper-parameters from data, Wenk et al. (2020) propose a preprocessing step, in which for each dimension, ODIN maximizes the marginal likelihood by solving

$$\max \quad p(\mathbf{y}|\phi, \sigma) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{C}_\phi + \sigma^2 \mathbf{I}) \quad (9)$$

$$\text{w.r.t.} \quad \phi, \sigma \quad (10)$$

The model mismatch parameter γ can either be hand-tuned (as done e.g. by Wenk et al. (2018) or Gorbach et al. (2017)) or inferred at run time using ODIN as well. In this case, the risk term of Equations (6) - (8) is extended by the additional summand $\log \det(\mathbf{A} + \gamma \mathbf{I})$, representing the contribution of the determinant term of the conditional $p(\mathbf{F}|\mathbf{x}, \mathbf{y}, \phi, \gamma)$. In this work, we will refer to this case whenever we speak about learning γ . Clearly, ODIN scales similarly to standard GP regression with derivatives cubically in the amount of observations due to the inversions of \mathbf{C}_ϕ and \mathbf{A} . Again, all quantities that are not defined in detail can be looked up in Section C.1.

3. QFF for Derivative Information

For standard GP regression, the main computational challenge lies in calculating any terms involving the covariance matrix \mathbf{C}_ϕ . In the context of Hermitian quadrature, the work of Mutny and Krause (2018) effectively reduces this complexity from $\mathcal{O}(N^3)$ to $\mathcal{O}(NM^2 + M^3)$, where $M \ll N$ is a fixed constant representing the amount of features used in the feature expansion.

However, as soon as we introduce derivative information, we obtain an additional term depending on the matrices \mathbf{D} and \mathbf{A} . Similar to \mathbf{C}_ϕ , the size of these matrices grows linearly with the amount of observations N , again leading to a computational complexity of $\mathcal{O}(N^3)$. Since the calculation of these matrices involves the first two derivatives of the kernel, we need to extend the standard QFF framework to include such derivatives.

Since the functional form of the kernel dictates the QFF features, any analysis would need to be repeated for every kernel. We will thus restrict ourselves to the RBF kernel, even though it could in principle be extended to any stationary kernel. The RBF kernel is defined as $k(t_i, t_j) = k(t_i - t_j) = k(r_{ij}) = \rho \exp(-\frac{r_{ij}^2}{2l^2})$, where $r_{ij} := t_i - t_j$. Here, the variance ρ and lengthscale l are the hyper-parameters of the kernel, which need to be learned in a preprocessing step. The RBF kernel is known for its excellent smoothing properties. Since scaling up the amount of observations is mainly relevant for densely observed trajectories, this makes it an ideal choice for the experiments we will show in Section 5.

3.1 Approximate Derivative of the Kernel

Since both the integral and the derivative operator are linear, we can use their commutativity to obtain new quadrature schemes for the kernel derivatives. For ease of readability, let us define $I(f(\omega)) := \int_{-\infty}^{+\infty} e^{-\omega^2} f(\omega) d\omega$ and fix $\rho = \sqrt{\pi}$. This assumption is simply to avoid cluttered notation and will could be lifted by introducing a constant factor instead. For ease of notation, let $r := t_i - t_j$.

Applying Bochner's theorem (1) to the RBF kernel with $\rho = \sqrt{\pi}$, we get, similarly to Mutny and Krause (2018),

$$k(t_i, t_j) = I(\cos(\omega r \frac{\sqrt{2}}{l})). \quad (11)$$

Using the linearity of the integral operator, we can differentiate the above equality to obtain

$$\frac{\partial}{\partial t_i} k(t_i, t_j) = \frac{d}{dr} I(\cos(\omega r \frac{\sqrt{2}}{l})) = I(-\frac{\sqrt{2}}{l} \omega \sin(\omega r \frac{\sqrt{2}}{l})). \quad (12)$$

Similarly,

$$\frac{\partial}{\partial t_j} k(t_i, t_j) = I(\omega \frac{\sqrt{2}}{l} \sin(\omega r \frac{\sqrt{2}}{l})) \quad (13)$$

and

$$\frac{\partial^2}{\partial t_i \partial t_j} k(t_i, t_j) = I(\omega^2 \frac{2}{l^2} \cos(\omega r \frac{\sqrt{2}}{l})). \quad (14)$$

These calculations reduce the problem of approximating the kernel derivatives to approximating integrals. Thus, similar to the derivative free case, we can now leverage the powerful framework of Gauss-Hermite quadrature (Hildebrand, 1987).

Let

$$Q_m(f(\omega)) = \sum_{i=1}^m W_i^m f(\omega_i^m) \quad (15)$$

denote the Gauss-Hermite quadrature scheme of order m for the function f , where W_i^m are its weights and ω_i^m its abscissas. Construct the $2m$ dimensional feature vector $\phi(x)$ by adding for each i the two components

$$\begin{bmatrix} \sqrt{W_i^m} \cos(\omega_i^m \frac{\sqrt{2}}{l} t_i) & \sqrt{W_i^m} \sin(\omega_i^m \frac{\sqrt{2}}{l} t_i) \end{bmatrix}. \quad (16)$$

Given this feature vector, we first observe that

$$\phi(t_i)^T \phi(t_j) = \sum_{i=1}^m W_i^m \cos(\omega_i^m \frac{\sqrt{2}}{l} r) = Q_m(\cos(\omega \frac{\sqrt{2}}{l} r)), \quad (17)$$

as desired when recovering the QFF approximation for $k(t_i, t_j)$ for calculating C_ϕ .

However, this feature expansion can also be differentiated w.r.t. t_i to obtain

$$\phi(t_i)' = \begin{bmatrix} -\sqrt{W_i^m} \frac{\sqrt{2}}{l} \omega_i^m \sin(\omega_i^m \frac{\sqrt{2}}{l} t_i) \\ \sqrt{W_i^m} \frac{\sqrt{2}}{l} \omega_i^m \cos(\omega_i^m \frac{\sqrt{2}}{l} t_i) \end{bmatrix}_{i=1 \dots m}. \quad (18)$$

Using trigonometric identities, it can be shown that this feature expansion indeed yields

$$\phi(t_i)^T \phi(t_j) = Q_m(-\frac{\sqrt{2}}{l} \omega \sin(\omega r \frac{\sqrt{2}}{l})), \quad (19)$$

$$\phi(t_i)^T \phi(t_j)' = Q_m(\frac{\sqrt{2}}{l} \omega \sin(\omega r \frac{\sqrt{2}}{l})), \quad (20)$$

$$\phi(t_i)'^T \phi(t_j)' = Q_m(\omega^2 \frac{2}{l^2} \cos(\omega r \frac{\sqrt{2}}{l})). \quad (21)$$

Thus, the features given by (18) represent the quadrature schemes for the integrals shown in Equations (12) and (14). Thus, the kernel derivatives involved in calculating \mathbf{A} and \mathbf{D} can be approximated by

$$\frac{\partial}{\partial x} k(x, y) \approx \phi(x)'^T \phi(y) \quad (22)$$

and

$$\frac{\partial^2}{\partial x \partial y} k(x, y) \approx \phi(x)^T \phi(y)'. \quad (23)$$

As we shall prove and empirically validate in the rest of this paper, this feature expansion is efficient in the sense that the approximation error decays *exponentially* in the amount of features. This allows us to obtain accurate approximations of the kernel derivatives even for a small amount of *deterministically* chosen features.

3.2 Application to GPR with Derivatives

To apply this approximation to GPR with derivatives, we make use of the matrix inversion lemma.

Lemma 1 (Matrix Inversion Lemma) *Let A and C be invertible matrices. Then for any matrices U , V of appropriate dimensions, it holds that $(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}$.*

In the context of standard GP regression (e.g. Rahimi and Recht, 2008; Mutny and Krause, 2018), this can be used as follows. Let $\Phi \in \mathbb{R}^{2m \times N}$ be such that its columns are given by the feature vectors $\phi(t_i)$ as defined in (16) for a quadrature scheme of order m . After adding a small jitter λ on the diagonal of C_ϕ , we can use Lemma 1 to obtain

$$(C_\phi + \lambda I)^{-1} \approx \frac{1}{\lambda} (I - \Phi^T (\Phi \Phi^T + \lambda \mathbb{I})^{-1} \Phi) \quad (24)$$

Here, \approx should be read as an approximation up to an exponentially decaying error in M , which we will quantify in Section 4.

Since all matrices required to calculate the quantities in Equations (4) and (5) allow for a direct feature expansion, the matrix inversion lemma can be directly applied to approximate these quantities. Since we work directly with the posterior mean and variance, we do not even need a jitter term due to the presence of the noise terms.

3.3 SLEIPNIR

However, extending this concept to Equation (8) is slightly more tricky, as neither D nor A allow for a direct feature expansion. Nevertheless, it is still possible to derive a scalable approximation using the previously introduced feature expansions:

First, start by collecting the appropriate feature vectors from Equations (16) and (18) to write $C'_\phi \approx \Phi'^T \Phi$, $C_\phi \approx \Phi^T \Phi'$ and $C''_\phi \approx \Phi'^T \Phi'$. These approximations can then be inserted into the term given by Equation (8). By applying Theorem 1 multiple times, it is possible to finally obtain

$$\begin{aligned} z^T (A + \gamma I)^{-1} z &\approx \\ &\frac{1}{\gamma} z^T (\mathbb{I} - \Phi'^T (\Phi' \Phi'^T + \frac{\gamma}{\lambda} \Phi \Phi^T + \gamma \mathbb{I})^{-1} \Phi') z, \end{aligned} \quad (25)$$

where

$$z := f(\theta, x) - Dx \approx f(\theta, x) - \Phi'^T (\Phi \Phi^T + \lambda \mathbb{I})^{-1} \Phi x.$$

Combining these approximations with the one given by Equation (24), we obtain a computationally efficient scheme for calculating all objective functions in ODIN. Furthermore, both the additional term $\log \det(A + \gamma I)$ obtained when learning γ and learning the hyperparameters via Equation (9) can be scaled similarly using these approximations and directly applying the matrix inversion lemma.

Overall, this means that the original complexity of $\mathcal{O}(N^3)$ has been successfully reduced to $\mathcal{O}(NM^2 + M^3)$. For $M < N$, this is significantly accelerating ODIN, which is why we name our approximation scheme SLEIPNIR. The resulting algorithm will be referred to as ODIN-S.

4. Theoretical Results

It is clear that this bound is especially appealing, since the error will exhibit exponential decay for $m > l^2$. In this section, we state that this behavior can be generalized to our approximations of the kernel derivatives. As we have demonstrated in the previous section, our feature approximation can be used to efficiently reduce the computational complexity of GPRD and ODIN. In this section, we will provide an in-depth error analysis of this approximation, demonstrating the key theoretical result of this paper: The favorable properties of standard QFF carry over to the derivative case and can be efficiently deployed in practically relevant algorithms.

4.1 Exponentially Decaying Kernel Approximation Errors

Let us define

$$E_m := \sqrt{\pi} \frac{1}{m^m} \left(\frac{e}{4l^2} \right)^m \quad (26)$$

In the context of standard GP regression, Mutny and Krause (2018) have shown that

$$|k(t_i, t_j) - \phi(t_i)^T \phi(t_j)| \leq E_m \quad (27)$$

The results are summarized in Theorem 2, with a detailed proof given in the appendix in Section A.

Theorem 2 *Let $k(t_i, t_j)$ be defined as in Equation (11) and consider the Gauss-Hermite quadrature scheme of order m as described by Equation (15), defining $\phi(t_i)$ and $\phi(t_i)'$ as in Equations (16) and (18). Then, for $|r| = |t_i - t_j| \leq 1$, it holds that*

$$(i) \quad \left| \frac{\partial}{\partial t_i} k(t_i, t_j) - \phi(t_i)^T \phi(t_j) \right| \leq \frac{2e}{l^2} E_{m-2}$$

$$(ii) \quad \left| \frac{\partial^2}{\partial t_i \partial t_j} k(t_i, t_j) - \phi(t_i)^T \phi(t_j)' \right| \leq \frac{2e}{l^4} E_{m-3}$$

where E_m is defined as in Equation (26)

Not that in the above Theorem, $|r| < 1$ can always be achieved by rescaling the data and adapting the lengthscale l accordingly.

While the error is slightly larger for the derivative approximations than for the kernel itself, it is important to note that we are still getting the same exponential decaying behavior.

4.2 GPR with Derivatives Bounds

As we shall see in this section, this exponential decay carries over nicely to the case of GPR with derivatives. Define $e_{\tilde{\mu}}$, $e_{\tilde{\Sigma}}$, $e_{\tilde{\mu}'}$ and $e_{\tilde{\Sigma}'}$ as the absolute error between the feature approximations and the corresponding accurate quantities of the means and covariances of Equations (4) and (5). For each $\tau \in \mathbb{R}$, define e_{tot} as the maximum of these four errors. Using these definitions, we can show the exponential relation between feature approximation order m and the corresponding approximation error, as summarized in Theorem 3, with proof in the appendix.

Theorem 3 *Let us consider an RBF kernel with hyperparameters (ρ, l) and domain $[0, 1]$. Define $c := \min(\gamma, \sigma^2)$ and $R := \max(\|\mathbf{y}\|_\infty, \|\mathbf{F}\|_\infty)$. Let $C > 0$. Let us consider a QFF approximation scheme of order $m \geq 3 + \max\left(\frac{e}{2l^2}, \log\left(\frac{270n^2\rho^3 R}{l^8 c^2 C}\right)\right)$. Then, it holds for all $\tau \in [0, 1]$ that $e_{\text{tot}} \leq C$.*

From this theorem, we can also observe the following fact: If we decrease the acceptable worst case performance C , we clearly need more features. However, due to the logarithm in the theorem, an exponential decay in C only leads to linear growth in m , all other things being equal.

4.3 SLEIPNIR Bounds

Similarly, we can observe that when applying SLEIPNIR to ODIN for a fixed feature vector length M , we will always create a small error in the objective function. However, the deterministic nature of the QFF approximation still allows us to choose the amount of features in a way such that the approximation error will always be guaranteed to be smaller than a pre-chosen threshold. This result is summarized in Theorem 4 and proven in the appendix, Section E.

Theorem 4 *Let \mathcal{R} be the ODIN-objective as defined in Equations (6)-(8) and let $\tilde{\mathcal{R}}$ denote its counterpart obtained by approximating the matrices \mathbf{C}_ϕ , \mathbf{A} and \mathbf{D} as described in Section 3.3. Assume the parameters λ , γ , $\phi = (\rho, l)$ and N to be fixed. Suppose $N \geq 60$ and let $1 > \epsilon > 0$. Let m denote the order of the quadrature scheme. Then*

$$m \geq 10 + \max\left\{\frac{e}{2l^2}, \log_2\left(\frac{\rho^2 n^3}{\lambda^2 \gamma l^4 \epsilon}\right)\right\} \quad (28)$$

implies

$$\frac{|R_{\lambda\gamma\phi}(x, \theta) - \tilde{R}_{\lambda\gamma\phi}(x, \theta)|}{R_{\lambda\gamma\phi}(x, \theta)} \leq \epsilon \quad (29)$$

for any configuration of the variables \mathbf{x} and $\boldsymbol{\theta}$.

Again, it should be noted that the threshold ϵ appears inside the logarithm. Thus, an *exponential* decrease of the allowed error only requires a *linear* increase in the amount of features.

In this theorem, the bound on N is purely aesthetically motivated, to reduce the amount of terms. From a practical perspective however, one would probably want to use ODIN without any approximations anyways if only < 60 observations are available. Furthermore, the bound on ϵ just requires that investigating an approximation scheme with more than 100% relative error is of little interest.

5. Empirical Validation

In this section, we will provide empirical validation for all theoretical claims made in this paper. We will start by showing the exponentially decaying behavior of the kernel approximation error in Section 5.1. In Sections 5.2 and 5.3, we show that this exponentially decaying behavior carries over to both GPR with derivatives and ODIN. Ultimately, we will conclude by showing how ODIN with SLEIPNIR can be used on a realistic dataset, demonstrating its practical applicability in Section 5.4.

Wherever applicable, we will use three standard benchmark systems from the GP-based gradient matching literature, namely the Lotka Volterra (LV) system (Lotka, 1932), the Protein Transduction (PT) system (Vyshemirsky and Girolami, 2007) and the Lorenz system (Lorenz, 1963). While the LV system is quite easy to fit, both PT and Lorenz offer interesting challenges. The non-stationary dynamics of PT make it a formidable challenge for collocation methods. Classic literature (e.g. Dondelinger et al., 2013; Gorbach et al., 2017; Wenk et al., 2020) sidestep this issue by using a non-stationary sigmoid kernel. However, as we shall demonstrate in this section, this is not a problem in the case of densely observed trajectories, since we will be using the RBF kernel in all experiments. Finally, the Lorenz system is interesting due to its chaotic behavior. Chaotic systems are an interesting challenge for many parameter inference schemes due to the potentially high sensitivity to parameter changes and the presence of many local optima. For a full description of the experimental setup

and all metrics used please refer to Section F. For all experiments, we created our data-set using numerical integration of the ground truth parameters. We then added 25 different noise realizations to obtain 25 different data sets. This allows us to quantify robustness w.r.t. noise by showing median as well as 20% and 80% quantiles over these noise realizations for each experiment. In all experiments, we trained γ and learned the kernel hyperparameters from the data using the scalable approximations described in the previous section. To compare with ODIN-S, we chose to combine ODIN with the RFF and RFF-B feature expansions detailed in Section 2.

5.1 Exponentially Decaying Kernel Approximation Errors

In Figure 2, we evaluate the maximum error of the feature approximations over an interval $r \in [0, 1]$. From Theorem 2, we would expect an exponential decay in the approximation error for both the kernel and the first and second order derivatives. This exponential decay is clearly visible and continues until we hit numerical boundaries. In Figure 2, the l was chosen to be 0.1. However, as can be seen on the additional plots in the appendix (Sec. B), this behavior is robust across different lengthscales. Figure 2 allows for the interesting observation that next to the exponential decay of QFF, the error of the RFF almost looks constant, even though it decays linearly.

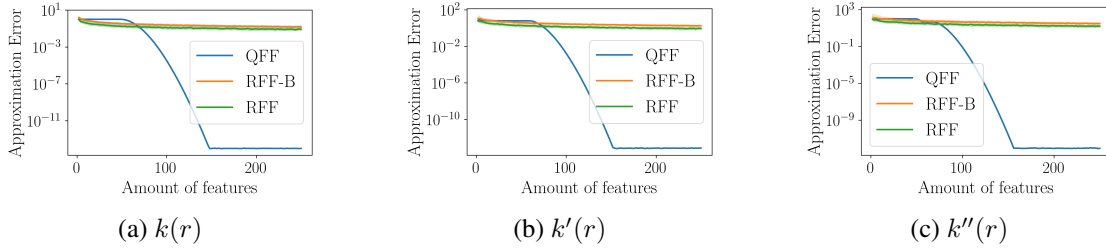


Figure 2: Comparing the maximum error of different feature expansions over $r \in [0, 1]$. For the random feature expansions, we show median as well as 12.5% and 87.5% quantiles over 100 random samples. Due to the exponential decay of the error of the QFF approximation, this stochasticity is barely visible. As given by the theoretical analysis, the error is a bit higher for the derivatives, but still decaying exponentially.

5.2 GPR with Derivatives

A similar behavior can be observed for the absolute error of the GP posterior we introduced in Section 2.2. As can be seen in Figure 3, the exponential decay of the error predicted in Section 4.2 also appears in practice. While we only show one state of one experiment, this behavior is consistent across experiments and noise settings, as can be seen by looking at the additional plots presented in the appendix, Section D.

5.3 ODIN-S on Standard Benchmark Systems

In Figure 4, we compare the performance of ODIN-S against accurate ODIN as well as ODIN augmented with RFF and RFF-B on three standard benchmark systems. In the top row, we keep the total amount of observations fixed to 1000 for LV and 2000 for PT and Lorenz, while varying the length of the feature vector. In the bottom row, we keep the amount of features fixed to 40 for LV, 300 for PT and 150 for Lorenz. All the data was created using observation noise with $\sigma^2 = 0.1$ for

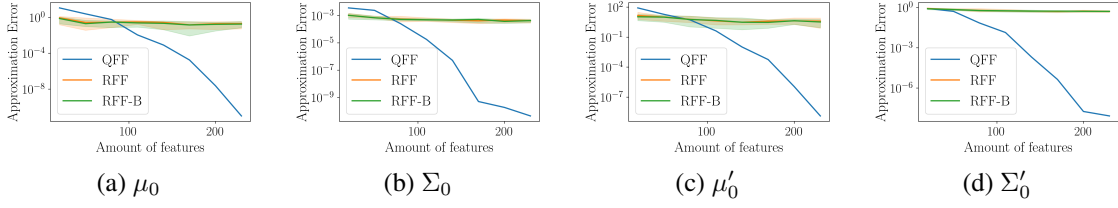


Figure 3: Approximation error of the different feature approximations compared to the accurate GP, evaluated at $t = 0.8$ for the Lorenz system with 1000 observations and an SNR of 100. For each feature, we show the median as well as the 12.5% and 87.5% quantiles over 10 independent noise realizations for the first state dimension.

LV, $\sigma^2 = 0.01$ for PT and a signal-to-noise ratio of 5 for Lorenz. Due to computational restrictions, it was not possible to evaluate accurate ODIN beyond what is shown in the plots. To provide an idea of the robustness of the evaluation, different noise, feature and observation settings are investigated in the appendix, Section G.

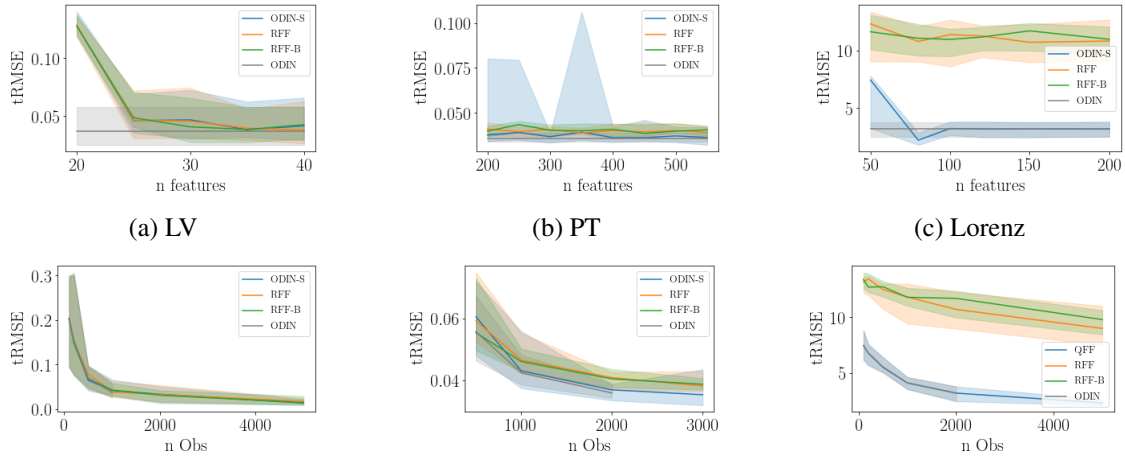


Figure 4: Comparing the tRMSE of the ODE parameters obtained via the three different approximation schemes. In the top row, we show how the results obtained by ODIN-S converge to the accurate results when increasing the amount of features. In the second row, we show the learning curves for a fixed feature vector length, demonstrating that increasing the amount of observations actually improves the parameter estimates and decreases the tRMSE.

As expected by our theoretical analysis, the trajectory RMSE of the SLEIPNIR-based ODIN-S eventually converges to the approximation-free ODIN in both median and quantiles if we increase the number of Fourier features. This clearly validates the main result of Theorem 4, expecting an exponentially small error. Also, it is interesting to observe that the MCMC-based RFF and RFF-B seem to struggle to keep up with the approximation quality of SLEIPNIR. While there is not much difference visible for the smooth and simple Lotka Volterra system, this is clearly visible for PT and especially striking for the chaotic Lorenz system. Complementing the results of Figure 2 and Figure 3, this further illustrates the power of the SLEIPNIR approximation, especially in this setting, where the kernel inputs are one-dimensional.

Table 1: Comparing the run time of accurate ODIN against the run time of the feature approximations per iteration in milliseconds. For each setting, we show the median \pm on standard deviation over 15 different iterations. As expected, the feature approximations lead to order of magnitude reductions in run time.

	accurate	feature
LV	973 ± 42.5	1.24 ± 0.243
PT	17900 ± 668	133 ± 12.8
Lorenz	10700 ± 466	13.2 ± 0.566

In Table 1, we finally compare the run time per iteration of ODIN with ODIN-S. Since the run time only depends on the amount of features and not on how the features were obtained, we omitted RFF and RFF-B. The run time was evaluated using the same amount of observation and amount of feature combinations as in Figure 4, namely (2000, 150) for LV with $\sigma^2 = 0.1$, (2000, 300) for PT with $\sigma^2 = 0.01$ and (1000, 40) for Lorenz with an SNR of 5. While it should already be clear from theoretical analysis that ODIN-S scales linearly in the amount of observations N and cubic in the length of the feature vector M , this was confirmed empirically as well. However, due to space restrictions, the plots have been moved to the appendix and are shown in Figure 29 and Figure 30.

5.4 Practical Applicability

In a final experiment, we show that ODIN-S is able to scale to realistic data sets. For this, we introduce a 12-dimensional ODE system representing the dynamics of a 6DOF quadcopter. We observe the system under Gaussian noise with SNR=10 over the time interval $t = [0, 15]$. We assume a sampling frequency of 1kHz, leading to 15'000 observations. We then run ODIN-S on a standard laptop (Lenovo Carbon X1) and obtain results in roughly 80min. Up to our knowledge, this is the first time that a system of such dimensions has been solved with a Gaussian process based parameter inference scheme, clearly demonstrating the power of our framework. The resulting trajectories are shown in Figure 5, including example observation points to visualize the noise level. The estimates of ODIN-S are so good that the ground truth is barely visible.

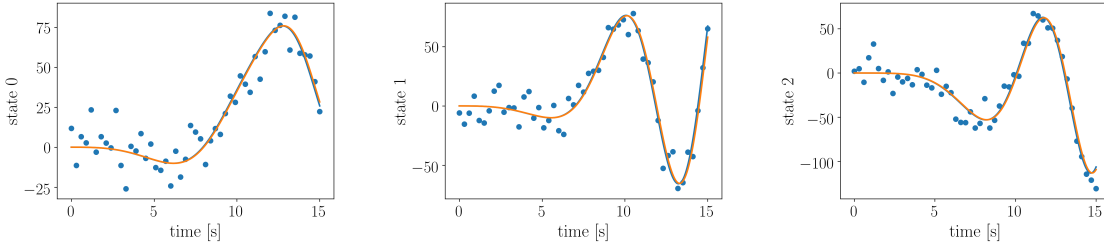


Figure 5: State trajectories of the first three states obtained by integrating the parameters inferred by ODIN-S (orange). The blue line represents the ground truth, while the blue dots show every 300-th observation for a signal-to-noise ratio of 10. States 5-12 have been moved to the appendix, Section G.3.

6. Conclusion

In this work, we introduced a new theoretical framework to scale up Gaussian process regression with derivative information, extending existing work based on quadrature Fourier features. We derived and proved deterministic, exponentially fast decaying error bounds for approximating the derivatives of an RBF kernel. We then combined these insights to create a computationally efficient approximation scheme for both standard GP regression with derivatives as well as the parameter inference scheme ODIN. The theoretical analysis of this approximation yielded deterministic, non-asymptotic error bounds. In an extensive empirical evaluation, we then showed orders of magnitude improvements on the run time without sacrificing accuracy. In future work, we are excited to see how SLEIPNIR could be deployed in other areas such as scaling up Bayesian optimization with derivatives (e.g. Wu et al., 2017) or probabilistic numerics (Hennig et al., 2015).

Acknowledgments

This research was supported by the Max Planck ETH Center for Learning Systems. This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme grant agreement No 815943.

References

- Gabriele Abbati, Philippe Wenk, Michael A. Osborne, Andreas Krause, Bernhard Schölkopf, and Stefan Bauer. AReS and MaRS adversarial and MMD-minimizing regression for SDEs. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1–10, Long Beach, California, USA, 09–15 Jun 2019. PMLR.
- Ben Calderhead, Mark Girolami, and Neil D Lawrence. Accelerating bayesian inference over nonlinear differential equations with gaussian processes. In *Advances in neural information processing systems*, pages 217–224, 2009.
- John P Cunningham, Krishna V Shenoy, and Maneesh Sahani. Fast gaussian process methods for point process intensity estimation. In *Proceedings of the 25th international conference on Machine learning*, pages 192–199. ACM, 2008.
- Frank Dondelinger, Dirk Husmeier, Simon Rogers, and Maurizio Filippone. Ode parameter inference using adaptive gradient matching with gaussian processes. In *Artificial intelligence and statistics*, pages 216–228, 2013.
- Leander Dony, Fei He, and Michael PH Stumpf. Parametric and non-parametric gradient matching for network inference: a comparison. *BMC bioinformatics*, 20(1):52, 2019.
- David Eriksson, Kun Dong, Eric Lee, David Bindel, and Andrew G Wilson. Scaling gaussian process regression with derivatives. In *Advances in Neural Information Processing Systems*, pages 6867–6877, 2018.
- Nico S Gorbach, Stefan Bauer, and Joachim M Buhmann. Scalable variational inference for dynamical systems. In *Advances in Neural Information Processing Systems*, pages 4806–4815, 2017.
- Philipp Hennig, Michael A Osborne, and Mark Girolami. Probabilistic numerics and uncertainty in computations. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 471(2179):20150142, 2015.
- James Hensman, Nicolas Durrande, Arno Solin, et al. Variational fourier features for gaussian processes. *Journal of Machine Learning Research*, 18(151):1–151, 2017.
- Francis Begnaud Hildebrand. *Introduction to numerical analysis*. Courier Corporation, 1987.
- Miguel Lazaro-Gredilla, Joaquin Quiñonero-Candela, Carl Edward Rasmussen, and Anibal R Figueiras-Vidal. Sparse spectrum gaussian process regression. *Journal of Machine Learning Research*, 11(Jun):1865–1881, 2010.
- Edward N Lorenz. Deterministic nonperiodic flow. *Journal of the atmospheric sciences*, 20(2): 130–141, 1963.
- Alfred J Lotka. The growth of mixed populations: two species competing for a common food supply. *Journal of the Washington Academy of Sciences*, 22(16/17):461–469, 1932.

- Benn Macdonald and Dirk Husmeier. Gradient matching methods for computational inference in mechanistic models for systems biology: a review and comparative analysis. *Frontiers in bioengineering and biotechnology*, 3:180, 2015.
- Jonas Mockus. *Bayesian approach to global optimization: theory and applications*, volume 37. Springer Science & Business Media, 2012.
- Mojmir Mutny and Andreas Krause. Efficient high dimensional bayesian optimization with additivity and quadrature fourier features. In *Advances in Neural Information Processing Systems*, pages 9005–9016, 2018.
- Joaquin Quiñonero-Candela and Carl Edward Rasmussen. A unifying view of sparse approximate gaussian process regression. *Journal of Machine Learning Research*, 6(Dec):1939–1959, 2005.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2008.
- CE. Rasmussen and CKI. Williams. *Gaussian Processes for Machine Learning*. Adaptive Computation and Machine Learning. MIT Press, Cambridge, MA, USA, January 2006.
- W. Rudin. *Principles of Mathematical Analysis*. International series in pure and applied mathematics. McGraw-Hill, 1976. ISBN 9780070856134.
- Simo Sarkka, Arno Solin, and Jouni Hartikainen. Spatiotemporal learning via infinite-dimensional bayesian filtering and smoothing: A look at gaussian process regression through kalman filtering. *IEEE Signal Processing Magazine*, 30(4):51–61, 2013.
- Sambu Seo, Marko Wallat, Thore Graepel, and Klaus Obermayer. Gaussian process regression: Active data selection and test point rejection. In *Mustererkennung 2000*, pages 27–34. Springer, 2000.
- Edward Snelson and Zoubin Ghahramani. Sparse gaussian processes using pseudo-inputs. In *Advances in neural information processing systems*, pages 1257–1264, 2006.
- Ercan Solak, Roderick Murray-Smith, William E Leithead, Douglas J Leith, and Carl E Rasmussen. Derivative observations in gaussian process models of dynamic systems. In *Advances in neural information processing systems*, pages 1057–1064, 2003.
- Arno Solin and Simo Särkkä. Hilbert space methods for reduced-rank gaussian process regression. *Statistics and Computing*, pages 1–28.
- Arno Solin, Manon Kok, Niklas Wahlström, Thomas B Schön, and Simo Särkkä. Modeling and interpolation of the ambient magnetic field by gaussian processes. *IEEE Transactions on robotics*, 34(4):1112–1127, 2018.
- Zoltán Szabó and Bharath K Sriperumbudur. On kernel derivative approximation with random fourier features. *arXiv preprint arXiv:1810.05207*, 2018.
- Michalis Titsias. Variational learning of inducing variables in sparse gaussian processes. In *Artificial Intelligence and Statistics*, pages 567–574, 2009.

- Vladislav Vyshemirsky and Mark A Girolami. Bayesian ranking of biochemical system models. *Bioinformatics*, 24(6):833–839, 2007.
- Philippe Wenk, Alkis Gotovos, Stefan Bauer, Nico Gorbach, Andreas Krause, and Joachim M Buhmann. Fast gaussian process based gradient matching for parameter identification in systems of nonlinear odes. *arXiv preprint arXiv:1804.04378*, 2018.
- Philippe Wenk, Gabriele Abbati, Stefan Bauer, Michael A Osborne, Andreas Krause, and Bernhard Schölkopf. Odin: Ode-informed regression for parameter and state inference in time-continuous dynamical systems. 2020.
- Christopher KI Williams and Matthias Seeger. Using the nyström method to speed up kernel machines. In *Advances in neural information processing systems*, pages 682–688, 2001.
- Andrew Wilson and Hannes Nickisch. Kernel interpolation for scalable structured gaussian processes (kiss-gp). In *International Conference on Machine Learning*, pages 1775–1784, 2015.
- Andrew G Wilson, Elad Gilboa, Arye Nehorai, and John P Cunningham. Fast kernel learning for multidimensional pattern extrapolation. In *Advances in Neural Information Processing Systems*, pages 3626–3634, 2014.
- Jian Wu, Matthias Poloczek, Andrew G Wilson, and Peter Frazier. Bayesian optimization with gradients. In *Advances in Neural Information Processing Systems*, pages 5267–5278, 2017.

Appendix A. Kernel Approximation Error Bounds

For any function f , let $I(f(\omega)) := \int_{-\infty}^{+\infty} e^{-\omega^2} f(\omega) d\omega$.

For any $r \in [0, 1]$, let $k(r) := \sqrt{\pi} e^{-\frac{r^2}{2l^2}} = \int_{-\infty}^{+\infty} e^{-\omega^2} \cos(\omega r \frac{\sqrt{2}}{l}) d\omega = I(\cos(\omega r \frac{\sqrt{2}}{l}))$.

Let $Q_m(f) = \sum_{i=1}^m W_i^m f(x_i^m)$ denote the Gauss-Hermite quadrature scheme of order m and function f with weights $W_i^m \geq 0$ and abscissas x_i^m and let $S^m := \{x_1^m, x_2^m, \dots, x_m^m\}$ denote the set of these abscissas.

From Gauss-Hermite quadrature, we know: If f is a polynomial of degree $2m - 1$ at most, then $Q_m(f(\omega)) = I(f(\omega))$, i.e the quadrature scheme exactly computes the integral.

Let $H_m(x)$ be the Hermite polynomial of order m and $h_m(x) := \frac{H_m(x)}{2^m}$ be its normalized version.

Since $I(H_i(\omega)H_j(\omega)) = r_{ij}2^j j! \sqrt{\pi}$, we know that $I(h_m(\omega)^2) = \sqrt{\pi} \frac{m!}{2^m}$.

Let $E_m := \sqrt{\pi} \frac{1}{m!} (\frac{e}{4l^2})^m$.

Using this definition, we can restate the error bounds derived by Mutny and Krause (2018) as

$$|I(\cos(\omega r \frac{\sqrt{2}}{l})) - Q_m(\cos(\omega r \frac{\sqrt{2}}{l}))| \leq E_m \quad (30)$$

and our bounds from Theorem 2 as

$$\frac{\sqrt{2}}{l} |I(\omega \sin(\omega r \frac{\sqrt{2}}{l})) - Q_m(\omega \sin(\omega r \frac{\sqrt{2}}{l}))| \leq 8(m-1)E_{m-1} \leq \frac{2e}{l^2} E_{m-2}, \quad (31)$$

$$\frac{2}{l^2} |I(\omega^2 \cos(\omega r \frac{\sqrt{2}}{l})) - Q_m(\omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))| \leq \frac{4}{l^2} (m-1)E_{m-2} \leq \frac{2e}{l^4} E_{m-3}. \quad (32)$$

Similar to normal quadrature, our proof technique is based on choosing polynomial q of degree $2m - 1$ whose values and derivatives agree with the original function f and its derivatives at a specific set of points. We can then define the remainder $s := f - q$. If we now approximate the integral of the function f by the integral of the polynomial q , bounding the approximation error is equivalent to bounding the integral of the remainder r .

We thus need two important components. After showing the existence of polynomials of a required degree that agree with f at specific points, we need some results regarding the residuals of polynomial approximations to functions. These components are demonstrated in the next two sections, which will then be combined in the final proof in Section A.3.

A.1 Polynomial approximation residuals

First, let us restate the following two well known Lemmata (see e.g. Hildebrand, 1987):

Lemma 5 *Let f be a real function n times continuously differentiable and $q_n(x)$ a polynomial of degree $n - 1$ that agrees with f at the set of distinct points $S = \{x_1, x_2, \dots, x_n\}$. Let $\pi(x) = \prod_{i=1}^n (x - x_i)$. Then $\forall x \in \mathbb{R}$ we have*

$$f(x) - q_n(x) = \frac{f^{(n)}(\xi)}{n!} \pi(x) \quad (33)$$

for some $\xi = \xi(x) \in I$, where I is the interval limited by the smallest and largest of the numbers x_1, x_2, \dots, x_n and x .

Lemma 6 *Let f be a real function $2n$ times continuously differentiable and $q_{2n}(x)$ a polynomial of degree $2n - 1$ that agrees with f at the set of distinct points $S = \{x_1, x_2, \dots, x_n\}$ and its derivative also agrees with f' at S . Let $\pi(x) = \prod_{i=1}^n (x - x_i)^2$. Then $\forall x \in \mathbb{R}$ we have*

$$f(x) - q_{2n}(x) = \frac{f^{(2n)}(\xi)}{(2n)!} \pi(x) \quad (34)$$

for some $\xi = \xi(x) \in I$, where I is the interval limited by the smallest and largest of the numbers x_1, x_2, \dots, x_n and x .

For our proof, we will need the following two extensions:

Lemma 7 *Let f be a real function $2n + 1$ times continuously differentiable and $q_{2n+1}(x)$ a polynomial of degree $2n$ that agrees with f at the set of distinct points $S = \{x_1, x_2, \dots, x_n, x_*\}$ and its derivative also agrees with f' at $S \setminus \{x_*\}$. Let $\pi(x) = (x - x_*) \prod_{i=1}^n (x - x_i)^2$. Then $\forall x \in \mathbb{R}$ we have*

$$f(x) - q_{2n+1}(x) = \frac{f^{(2n+1)}(\xi)}{(2n+1)!} \pi(x) \quad (35)$$

for some $\xi = \xi(x) \in I$, where I is the interval limited by the smallest and largest of the numbers $x_1, x_2, \dots, x_n, x_*$ and x .

Proof 1 *If $x \in S$ the property holds. Otherwise, fix \bar{x} . We define $F(x) = f(x) - q_{2n+1}(x) - K\pi(x)$, where K is chosen such that $F(\bar{x}) = 0$. We see that F has $n + 2$ distinct roots at I , so by Rolle's theorem F' has $n + 1$ distinct roots in I which are different from the points of S (because F' vanishes at intermediate points). Moreover, F' vanishes at all the points in $S \setminus \{x_*\}$ so in total F' has at least $2n + 1$ distinct roots in I . Consequently, by Rolle's theorem again, F'' has $2n$ distinct roots in I , ..., $F^{(2n+1)}$ has one root $\xi \in I$. At this point we have*

$$0 = F^{(2n+1)}(\xi) = f^{(2n+1)}(\xi) - (2n+1)!K \Leftrightarrow K = \frac{f^{(2n+1)}(\xi)}{(2n+1)!}. \quad (36)$$

So $F(\bar{x}) = 0 = f(\bar{x}) - q_{2n+1}(\bar{x}) - \frac{f^{(2n+1)}(\xi)}{(2n+1)!} \pi(\bar{x})^2 \Leftrightarrow f(\bar{x}) - q_{2n+1}(\bar{x}) = \frac{f^{(2n+1)}(\xi)}{(2n+1)!} \pi(\bar{x})^2$. And since \bar{x} was chosen arbitrarily the result holds for all x .

Lemma 8 *Let f be a real function $2n + 1$ times continuously differentiable and $q_{2n+1}(x)$ a polynomial of degree $2n$ that agrees with f at the set of distinct points $S = \{x_1, x_2, \dots, x_n\}$, its derivative also agrees with f' at S and its second derivative agrees with the second derivative of f at x_1 . Let $\pi(x) = (x - x_1) \prod_{i=1}^n (x - x_i)^2$. Then $\forall x \in \mathbb{R}$ we have*

$$f(x) - q_{2n+1}(x) = \frac{f^{(2n+1)}(\xi)}{(2n+1)!} \pi(x) \quad (37)$$

for some $\xi = \xi(x) \in I$, where I is the interval limited by the smallest and largest of the numbers x_1, x_2, \dots, x_n and x .

Proof 2 *If $x \in S$ the property holds. Otherwise, fix \bar{x} . We define $F(x) = f(x) - q_{2n+1}(x) - K\pi(x)$, where K is chosen such that $F(\bar{x}) = 0$. We see that F has $n + 1$ distinct roots at I , so by Rolle's*

theorem F' has n distinct roots in I which are different from the points of S (because F' vanishes at intermediate points). Moreover, F' vanishes at all the points in S so in total F' has at least $2n$ distinct roots in I . Consequently, by Rolle's theorem again, F'' has $2n - 1$ distinct roots in I which are different from the points of S (because F'' vanishes at intermediate points) and also F'' vanishes at x_1 , so in total F'' vanishes at $2n$ points (at least), ..., $F^{(2n+1)}$ has one root $\xi \in I$. At this point we have

$$0 = F^{(2n+1)}(\xi) = f^{(2n+1)}(\xi) - (2n+1)!K \Leftrightarrow K = \frac{f^{(2n+1)}(\xi)}{(2n+1)!}. \quad (38)$$

So $F(\bar{x}) = 0 = f(\bar{x}) - q_{2n+1}(\bar{x}) - \frac{f^{(2n+1)}(\xi)}{(2n+1)!}\pi(\bar{x})^2 \Leftrightarrow f(\bar{x}) - q_{2n+1}(\bar{x}) = \frac{f^{(2n+1)}(\xi)}{(2n+1)!}\pi(\bar{x})^2$. And since \bar{x} was chosen arbitrarily the result holds for all x .

A.2 Existence of polynomials that agree with function at certain points

Lemma 5 - 8 assume the existence of polynomials with certain properties. In the following, we will show that such polynomials actually exist and how they can be constructed.

Let $S = \{x_1, x_2, \dots, x_n\}$ be a set of n distinct points in \mathbb{R} and define $\pi(x) = \prod_{i=1}^n (x - x_i)$ and for $i = 1, \dots, n$, $l_i(x) = \frac{\pi(x)}{(x - x_i)\pi'(x_i)}$. It holds that $l_i(x_j) = r_{ij}$ and that $\deg(l_i) = n - 1$. Moreover, we define for $i = 1, \dots, n$, $h_i(x) = l_i^2(x)(-2l_i'(x_i)(x - x_i) + 1)$ and $\bar{h}_i(x) = l_i^2(x)(x - x_i)$. It holds that $h_i(x_j) = r_{ij}$, $h_i'(x_j) = 0$, $\bar{h}_i(x_j) = 0$ and $\bar{h}_i'(x_j) = r_{ij}$ and $\deg(h_i) = \deg(\bar{h}_i) = 2n - 1$ for $i = 1, \dots, n$. Thus, we can directly state the following four Lemmata

Lemma 9 Let f be a real function and define $y_1(x) = \sum_{i=1}^n f(x_i)l_i(x)$. Then y_1 is a polynomial with degree $n - 1$ that agrees with f at S .

Lemma 10 Let f be a real differential function and define $y_2(x) = \sum_{i=1}^n f(x_i)h_i(x) + \sum_{i=1}^n f'(x_i)\bar{h}_i(x)$. Then y_2 is a polynomial with degree $2n - 1$ that agrees with f at S and its derivative also agrees with the derivative in our proofs of f at S .

Lemma 11 Consider a point $x_0 \notin S$. Let f be a real differential function and define $y_3(x) = y_2(x) + \frac{f(x_0) - y_2(x_0)}{\pi^2(x_0)}\pi^2(x)$. Then y_3 is a polynomial with degree $2n$ that agrees with f at $S \cup \{x_0\}$ and its derivative also agrees with the derivative of f at S .

Lemma 12 Consider the point $x_1 \in S$. Let f be a real two times differential function and define $y_4(x) = y_2(x) + \frac{f''(x_1) - y_2''(x_1)}{2\pi'^2(x_1)}\pi^2(x)$. Then y_4 is a polynomial with degree $2n$ that agrees with f at S , its derivative also agrees with the derivative of f at S and its second derivative agrees with the second derivative of f at x_1 .

A.3 Derivation of bounds

In this section, we prove one of the two main theoretical results of our paper, namely the bounds on the error of our feature approximations (Theorem 2). We will start with Theorem 13, restating a well known result from Gauss-Hermite quadrature (Hildebrand, 1987). This theorem will inspire the proofs for Theorem 14 and 15, where we prove similar bounds for the more complex cases of first and second order derivatives.

Theorem 13 Consider the function $\cos(\frac{\sqrt{2}}{l}r\omega)$. If we approximate the integral $I(\cos(\frac{\sqrt{2}}{l}r\omega)) = \int_{-\infty}^{+\infty} e^{-\omega^2} \cos(\omega r \frac{\sqrt{2}}{l}) d\omega$ by $Q_m(\cos(\frac{\sqrt{2}}{l}r\omega))$, the Gauss-Hermite quadrature scheme of order m , we have

$$|I(\cos(\frac{\sqrt{2}}{l}r\omega)) - Q_m(\cos(\frac{\sqrt{2}}{l}r\omega))| \leq E_m. \quad (39)$$

Proof 3 Let $y_{2m}(\omega)$ be a polynomial of degree $2m - 1$ that agrees with $\cos(\frac{\sqrt{2}}{l}r\omega)$ at S^m and its derivate also agrees with the derivative of $\cos(\frac{\sqrt{2}}{l}r\omega)$ at S^m (we know that such a polynomial exists by Lemma 10). By Lemma 6 we have for every ω that

$$\cos(\frac{\sqrt{2}}{l}r\omega) - y_{2m}(\omega) = (\frac{\sqrt{2}r}{l})^{2m} \frac{\cos(\xi)}{(2m)!} h_m^2(\omega) \quad (40)$$

and since $r \leq 1$ and $|\cos(\xi)| \leq 1$ we get

$$|\cos(\frac{\sqrt{2}}{l}r\omega) - y_{2m}(\omega)| \leq (\frac{\sqrt{2}}{l})^{2m} \frac{h_m^2(\omega)}{(2m)!}. \quad (41)$$

Moreover, $Q_m(\cos(\frac{\sqrt{2}}{l}r\omega)) = Q_m(y_{2m}(\omega))$ (because $\cos(\frac{\sqrt{2}}{l}r\omega)$ and $y_{2m}(\omega)$ agree at S^m) and $Q_m(y_{2m}(\omega)) = I(y_{2m}(\omega))$ (because y_{2m} has degree less than $2m$) so $Q_m(\cos(\frac{\sqrt{2}}{l}r\omega)) = I(y_{2m}(\omega))$ and

$$|I(\cos(\frac{\sqrt{2}}{l}r\omega)) - Q_m(\cos(\frac{\sqrt{2}}{l}r\omega))| = |I(\cos(\frac{\sqrt{2}}{l}r\omega)) - I(y_{2m}(\omega))| \leq \quad (42)$$

$$I(|\cos(\omega r \frac{\sqrt{2}}{l}) - y_{2m}(\omega)|) \leq I((\frac{\sqrt{2}}{l})^{2m} \frac{h_m^2(\omega)}{(2m)!}) = \quad (43)$$

$$(\frac{\sqrt{2}}{l})^{2m} \frac{m! \sqrt{\pi}}{2^m (2m)!} \leq E_m. \quad (44)$$

In the preceding proof, the basic idea was to approximate the function with a polynomial (of degree less than $2m$) that agrees with the function at a specific set of points. This gives us a remainder that can be relatively efficiently bounded, leading to a tight error bound. In the following proofs, the main challenge lies in finding the right approximating polynomial (of degree less than $2m$) that yields an easy to handle and efficiently bounded remainder. Once such a polynomial is constructed, we can use the following idea: Let $f(\omega)$ be the function to be approximated by the approximating polynomial $p(\omega)$ and assume that the remainder can be absolutely bounded by the polynomial $s(\omega)$. Then

$$|I(f(\omega)) - Q_m(f(\omega))| = |I(f(\omega)) - I(p(\omega)) + Q_m(p(\omega)) - Q_m(f(\omega))| \leq \quad (45)$$

$$|I(f(\omega)) - I(p(\omega))| + |Q_m(p(\omega)) - Q_m(f(\omega))| \leq \quad (46)$$

$$I(|f(\omega) - p(\omega)|) + Q_m(|f(\omega) - p(\omega)|) \leq I(s(\omega)) + Q_m(s(\omega)). \quad (47)$$

If $s(\omega)$ has degree less than $2m$ so that $I(s(\omega)) = Q_m(s(\omega))$, then the final bound is $2I(s(\omega))$. Following standard practice, $s(\omega)$ will be chosen as the square of a polynomial. This is motivated by the following observation:

Consider $s(\omega)$ to be the square of a polynomial of degree n . W.l.o.g assume $s(\omega)$ to be monic (has leading coefficient 1). Then $s(\omega) = (h_n(\omega) + q(\omega))^2$ with $\deg(q) < n$. Consequently, $I(s(\omega)) = I(h_n^2(\omega)) + I(q^2(\omega)) + 2I(h_n(\omega)q(\omega)) = I(h_n^2(\omega)) + I(q^2(\omega)) \geq I(h_n^2(\omega))$. This means that $s(\omega) = h_n^2(\omega)$ minimizes $I(s(\omega))$ and suggests that in our proofs, the approximating polynomial should agree with the function at the set S^n , so that the remainder $s(\omega)$ is of the form $h_n^2(\omega)$ and gives us good values for $I(s(\omega))$.

We now have all the necessary tools to state and prove Theorem 2. We split the theorem into two parts: Theorem 14 restates the claim of Theorem 2 for the first order derivative, while Theorem 15 restates the claim for the second order derivative.

Theorem 14 *Consider the function $\frac{\sqrt{2}}{l}\omega \sin(\omega r \frac{\sqrt{2}}{l})$. If we approximate the integral $I(\frac{\sqrt{2}}{l}\omega \sin(\omega r \frac{\sqrt{2}}{l})) = \int_{-\infty}^{+\infty} e^{-\omega^2} \frac{\sqrt{2}}{l}\omega \sin(\omega r \frac{\sqrt{2}}{l}) d\omega$ by $Q_m(\frac{\sqrt{2}}{l}\omega \sin(\omega r \frac{\sqrt{2}}{l}))$, the Gauss-Hermite quadrature scheme of order m , we have*

$$\frac{\sqrt{2}}{l} |I(\omega \sin(\omega r \frac{\sqrt{2}}{l})) - Q_m(\omega \sin(\omega r \frac{\sqrt{2}}{l}))| \leq 8(m-1)E_{m-1}. \quad (48)$$

Proof 4 *Depending on whether m is odd or even we have:*

Suppose m is even. Then $0 \in S^{m-1}$. Let $y_{2m-3}(\omega)$ be a polynomial of degree $2m-4$ that agrees with $\sin(\frac{\sqrt{2}}{l}r\omega)$ at S^{m-1} and its derivative also agrees with the derivative of $\sin(\frac{\sqrt{2}}{l}r\omega)$ at $S^{m-1} \setminus \{0\}$ (we know that such a polynomial exists by Lemma 11). By Lemma 7 we have for every ω that

$$\omega(\sin(\frac{\sqrt{2}}{l}r\omega) - y_{2m-3}(\omega)) = (\frac{\sqrt{2}r}{l})^{2m-3} \frac{\cos(\xi)}{(2m-3)!} h_{m-1}^2(\omega) \quad (49)$$

so since $r \leq 1$ and $|\cos(\xi)| \leq 1$ we get

$$|\omega \sin(\frac{\sqrt{2}}{l}r\omega) - \omega y_{2m-3}(\omega)| \leq (\frac{\sqrt{2}r}{l})^{2m-3} \frac{h_{m-1}^2(\omega)}{(2m-3)!}. \quad (50)$$

Consequently, since the weights of a quadrature scheme W_i^m are positive and $Q_m(h_{m-1}^2(\omega)) = I(h_{m-1}^2(\omega))$ (because h_{m-1}^2 has degree less than $2m$), using the above relation we get

$$|Q_m(\omega y_{2m-3}(\omega)) - Q_m(\omega \sin(\omega r \frac{\sqrt{2}}{l}))| = |Q_m(\omega y_{2m-3}(\omega) - \omega \sin(\omega r \frac{\sqrt{2}}{l}))| \leq \quad (51)$$

$$Q_m(|\omega y_{2m-3}(\omega) - \omega \sin(\omega r \frac{\sqrt{2}}{l})|) \leq Q_m((\frac{\sqrt{2}}{l})^{2m-3} \frac{h_{m-1}^2(\omega)}{(2m-3)!}) = \quad (52)$$

$$I((\frac{\sqrt{2}}{l})^{2m-3} \frac{h_{m-1}^2(\omega)}{(2m-3)!}) := R \quad (53)$$

where

$$R = (\frac{\sqrt{2}}{l})^{2m-3} \frac{(m-1)! \sqrt{\pi}}{2^{m-1}(2m-3)!}. \quad (54)$$

Similarly,

$$|I(\omega y_{2m-3}(\omega)) - I(\omega \sin(\omega r \frac{\sqrt{2}}{l}))| = |I(\omega y_{2m-3}(\omega) - \omega \sin(\omega r \frac{\sqrt{2}}{l}))| \leq \quad (55)$$

$$I(|\omega y_{2m-3}(\omega) - \omega \sin(\omega r \frac{\sqrt{2}}{l})|) \leq I((\frac{\sqrt{2}}{l})^{2m-3} \frac{h_{m-1}^2(\omega)}{(2m-3)!}) = R. \quad (56)$$

Finally, using that $Q_m(\omega y_{2m-3}(\omega)) = I(\omega y_{2m-3}(\omega))$ (because $\omega y_{2m-3}(\omega)$ has degree less than $2m$) we get

$$\frac{\sqrt{2}}{l} |I(\omega \sin(\omega r \frac{\sqrt{2}}{l})) - Q_m(\omega \sin(\omega r \frac{\sqrt{2}}{l}))| = \quad (57)$$

$$\frac{\sqrt{2}}{l} |I(\omega \sin(\omega r \frac{\sqrt{2}}{l})) - I(\omega y_{2m-3}(\omega)) + Q_m(\omega y_{2m-3}(\omega)) - Q_m(\omega \sin(\omega r \frac{\sqrt{2}}{l}))| \leq \quad (58)$$

$$\frac{\sqrt{2}}{l} (|I(\omega \sin(\omega r \frac{\sqrt{2}}{l})) - I(\omega y_{2m-3}(\omega))| + |Q_m(\omega y_{2m-3}(\omega)) - Q_m(\omega \sin(\omega r \frac{\sqrt{2}}{l}))|) \leq \quad (59)$$

$$\frac{\sqrt{2}}{l} 2R = 2(\frac{\sqrt{2}}{l})^{2m-2} \frac{(m-1)!\sqrt{\pi}}{2^{m-1}(2m-3)!} \leq 4(m-1)E_{m-1}. \quad (60)$$

Suppose m is odd. Then $0 \in S^{m-2}$. Let $y_{2m-3}(\omega)$ be a polynomial of degree $2m-4$ that agrees with $\sin(\frac{\sqrt{2}}{l}r\omega)$ at S^{m-2} and its derivate also agrees with the derivative of $\sin(\frac{\sqrt{2}}{l}r\omega)$ at S^{m-2} and its second derivative agrees at 0 with with the second derivative of $\sin(\frac{\sqrt{2}}{l}r\omega)$ (we know that such a polynomial exists by Lemma 12). By Lemma 8 we have for every ω that

$$\omega(\sin(\frac{\sqrt{2}}{l}r\omega) - y_{2m-3}(\omega)) = (\frac{\sqrt{2}r}{l})^{2m-3} \frac{\cos(\xi)}{(2m-3)!} \omega^2 h_{m-2}^2(\omega) \quad (61)$$

so since $r \leq 1$ and $|\cos(\xi)| \leq 1$ we get

$$|\omega \sin(\frac{\sqrt{2}}{l}r\omega) - \omega y_{2m-3}(\omega)| \leq (\frac{\sqrt{2}}{l})^{2m-3} \frac{\omega^2 h_{m-2}^2(\omega)}{(2m-3)!}. \quad (62)$$

Consequently, since the weights of a quadrature scheme W_i^m are positive and $Q_m(\omega^2 h_{m-2}^2(\omega)) = I(\omega^2 h_{m-2}^2(\omega))$ (because $\omega^2 h_{m-2}^2$ has degree less than $2m$), using the above relation we get

$$|Q_m(\omega y_{2m-3}(\omega)) - Q_m(\omega \sin(\omega r \frac{\sqrt{2}}{l}))| = |Q_m(\omega y_{2m-3}(\omega) - \omega \sin(\omega r \frac{\sqrt{2}}{l}))| \leq \quad (63)$$

$$Q_m(|\omega y_{2m-3}(\omega) - \omega \sin(\omega r \frac{\sqrt{2}}{l})|) \leq Q_m((\frac{\sqrt{2}}{l})^{2m-3} \frac{\omega^2 h_{m-2}^2(\omega)}{(2m-3)!}) = \quad (64)$$

$$I((\frac{\sqrt{2}}{l})^{2m-3} \frac{\omega^2 h_{m-2}^2(\omega)}{(2m-3)!}) := R \quad (65)$$

where

$$R = (\frac{\sqrt{2}}{l})^{2m-3} \frac{I(\omega^2 h_{m-2}^2(\omega))}{(2m-3)!}.$$

We will now use the identity

$$xh_n(x) = h_{n+1}(x) + \frac{n}{2}h_{n-1}(x)$$

and that $I(h_{n+1}(\omega)h_{n-1}(\omega)) = 0$ (normality) so by squaring we have

$$I(\omega^2 h_{m-2}^2(\omega)) = I(h_{m-1}^2(\omega)) + (\frac{m-2}{2})^2 I(h_{m-3}^2(\omega)) \leq 2 \frac{(m-1)!\sqrt{\pi}}{2^{m-1}}. \quad (66)$$

So

$$R \leq 2\left(\frac{\sqrt{2}}{l}\right)^{2m-3} \frac{(m-1)!\sqrt{\pi}}{2^{m-1}(2m-3)!}. \quad (67)$$

Similarly,

$$|I(\omega y_{2m-3}(\omega)) - I(\omega \sin(\omega r \frac{\sqrt{2}}{l}))| = |I(\omega y_{2m-3}(\omega) - \omega \sin(\omega r \frac{\sqrt{2}}{l}))| \leq \quad (68)$$

$$I(|\omega y_{2m-3}(\omega) - \omega \sin(\omega r \frac{\sqrt{2}}{l})|) \leq I\left(\left(\frac{\sqrt{2}}{l}\right)^{2m-3} \frac{\omega^2 h_{m-2}^2(\omega)}{(2m-3)!}\right) = R. \quad (69)$$

Finally, using that $Q_m(\omega y_{2m-3}(\omega)) = I(\omega y_{2m-3}(\omega))$ (because $\omega y_{2m-3}(\omega)$ has degree less than $2m$) we get

$$\frac{\sqrt{2}}{l} |I(\omega \sin(\omega r \frac{\sqrt{2}}{l})) - Q_m(\omega \sin(\omega r \frac{\sqrt{2}}{l}))| = \quad (70)$$

$$\frac{\sqrt{2}}{l} |I(\omega \sin(\omega r \frac{\sqrt{2}}{l})) - I(\omega y_{2m-3}(\omega)) + Q_m(\omega y_{2m-3}(\omega)) - Q_m(\omega \sin(\omega r \frac{\sqrt{2}}{l}))| \leq \quad (71)$$

$$\frac{\sqrt{2}}{l} (|I(\omega \sin(\omega r \frac{\sqrt{2}}{l})) - I(\omega y_{2m-3}(\omega))| + |Q_m(\omega y_{2m-3}(\omega)) - Q_m(\omega \sin(\omega r \frac{\sqrt{2}}{l}))|) \leq \quad (72)$$

$$\frac{\sqrt{2}}{l} 2R = 4\left(\frac{\sqrt{2}}{l}\right)^{2m-2} \frac{(m-1)!\sqrt{\pi}}{2^{m-1}(2m-3)!} \leq 8(m-1)E_{m-1}. \quad (73)$$

Theorem 15 Consider the function $\frac{2}{l^2} \omega^2 \cos(\omega r \frac{\sqrt{2}}{l})$. If we approximate the integral $I(\frac{2}{l^2} \omega^2 \cos(\omega r \frac{\sqrt{2}}{l})) = \int_{-\infty}^{+\infty} e^{-\omega^2} \frac{2}{l^2} \omega^2 \cos(\omega r \frac{\sqrt{2}}{l}) d\omega$ by $Q_m(\frac{2}{l^2} \omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))$, the Gauss-Hermite quadrature scheme of order m , we have

$$\frac{2}{l^2} |I(\omega^2 \cos(\omega r \frac{\sqrt{2}}{l})) - Q_m(\omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))| \leq \frac{4}{l^2} (m-1)E_{m-2}. \quad (74)$$

Proof 5 Depending on whether m is odd or even we have:

Suppose m is even. Then $0 \in S^{m-1}$. Let $y_{2m-4}(\omega)$ be a polynomial of degree $2m-5$ that agrees with $\cos(\frac{\sqrt{2}}{l} r\omega)$ at $S^{m-1} \setminus \{0\}$ and its derivative also agrees with the derivative of $\cos(\frac{\sqrt{2}}{l} r\omega)$ at $S^{m-1} \setminus \{0\}$ (we know that such a polynomial exists by Lemma 11). By Lemma 7 we have for every ω that

$$\omega^2 (\cos(\frac{\sqrt{2}}{l} r\omega) - y_{2m-4}(\omega)) = \left(\frac{\sqrt{2}r}{l}\right)^{2m-4} \frac{\cos(\xi)}{(2m-4)!} h_{m-1}^2(\omega) \quad (75)$$

so since $r \leq 1$ and $|\cos(\xi)| \leq 1$ we get

$$|\omega^2 \cos(\frac{\sqrt{2}}{l} r\omega) - \omega^2 y_{2m-4}(\omega)| \leq \left(\frac{\sqrt{2}}{l}\right)^{2m-4} \frac{h_{m-1}^2(\omega)}{(2m-4)!}. \quad (76)$$

Consequently, since the weights of a quadrature scheme W_i^m are positive and $Q_m(h_{m-1}^2(\omega)) = I(h_{m-1}^2(\omega))$ (because h_{m-1}^2 has degree less than $2m$), using the above relation we get

$$|Q_m(\omega^2 y_{2m-4}(\omega)) - Q_m(\omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))| = |Q_m(\omega^2 y_{2m-4}(\omega) - \omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))| \leq \quad (77)$$

$$Q_m(|\omega^2 y_{2m-4}(\omega) - \omega^2 \cos(\omega r \frac{\sqrt{2}}{l})|) \leq Q_m((\frac{\sqrt{2}}{l})^{2m-4} \frac{h_{m-1}^2(\omega)}{(2m-4)!}) = \quad (78)$$

$$I((\frac{\sqrt{2}}{l})^{2m-4} \frac{h_{m-1}^2(\omega)}{(2m-4)!}) := R \quad (79)$$

where

$$R = (\frac{\sqrt{2}}{l})^{2m-4} \frac{(m-1)! \sqrt{\pi}}{2^{m-1} (2m-4)!}. \quad (80)$$

Similarly,

$$|I(\omega^2 y_{2m-4}(\omega)) - I(\omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))| = |I(\omega^2 y_{2m-4}(\omega) - \omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))| \leq \quad (81)$$

$$I(|\omega^2 y_{2m-4}(\omega) - \omega^2 \cos(\omega r \frac{\sqrt{2}}{l})|) \leq I((\frac{\sqrt{2}}{l})^{2m-4} \frac{h_{m-1}^2(\omega)}{(2m-4)!}) = R. \quad (82)$$

Finally using that $Q_m(\omega^2 y_{2m-4}(\omega)) = I(\omega^2 y_{2m-4}(\omega))$ (because $\omega^2 y_{2m-4}(\omega)$ has degree less than $2m$) and following exactly the same procedure as in the previous proof we get

$$\frac{2}{l^2} |I(\omega^2 \cos(\omega r \frac{\sqrt{2}}{l})) - Q_m(\omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))| \leq \frac{2}{l^2} 2R = 2(\frac{\sqrt{2}}{l})^{2m-2} \frac{(m-1)! \sqrt{\pi}}{2^{m-1} (2m-4)!} \leq \frac{2}{l^2} (m-1) E_{m-2}. \quad (83)$$

Suppose m is odd. Let $y_{2m-4}(\omega)$ be a polynomial of degree $2m-5$ that agrees with $\cos(\frac{\sqrt{2}}{l} r \omega)$ at S^{m-2} and its derivate also agrees with the derivative of $\cos(\frac{\sqrt{2}}{l} r \omega)$ at S^{m-2} (we know that such a polynomial exists by Lemma 10). By Lemma 6 we have for every ω that

$$\omega^2 (\cos(\frac{\sqrt{2}}{l} r \omega) - y_{2m-4}(\omega)) = (\frac{\sqrt{2} r}{l})^{2m-4} \frac{\cos(\xi)}{(2m-4)!} \omega^2 h_{m-2}^2(\omega) \quad (84)$$

so since $r \leq 1$ and $|\cos(\xi)| \leq 1$ we get

$$|\omega^2 \cos(\frac{\sqrt{2}}{l} r \omega) - \omega^2 y_{2m-4}(\omega)| \leq (\frac{\sqrt{2}}{l})^{2m-4} \frac{\omega^2 h_{m-2}^2(\omega)}{(2m-4)!}. \quad (85)$$

Consequently, since the weights of a quadrature scheme W_i^m are positive and $Q_m(\omega^2 h_{m-2}^2(\omega)) = I(\omega^2 h_{m-2}^2(\omega))$ (because $\omega^2 h_{m-2}^2$ has degree less than $2m$), using the above relation we get

$$|Q_m(\omega^2 y_{2m-4}(\omega)) - Q_m(\omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))| = |Q_m(\omega^2 y_{2m-4}(\omega) - \omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))| \leq \quad (86)$$

$$Q_m(|\omega^2 y_{2m-4}(\omega) - \omega^2 \cos(\omega r \frac{\sqrt{2}}{l})|) \leq Q_m((\frac{\sqrt{2}}{l})^{2m-4} \frac{\omega^2 h_{m-2}^2(\omega)}{(2m-4)!}) = \quad (87)$$

$$I((\frac{\sqrt{2}}{l})^{2m-4} \frac{\omega^2 h_{m-2}^2(\omega)}{(2m-4)!}) := R \quad (88)$$

where

$$R = \left(\frac{\sqrt{2}}{l}\right)^{2m-4} \frac{I(\omega^2 h_{m-2}^2(\omega))}{(2m-4)!}.$$

We will now use the identity

$$x h_n(x) = h_{n+1}(x) + \frac{n}{2} h_{n-1}(x)$$

and that $I(h_{n+1}(\omega)h_{n-1}(\omega)) = 0$ (normality) so by squaring we have

$$I(\omega^2 h_{m-2}^2(\omega)) = I(h_{m-1}^2(\omega)) + \left(\frac{m-2}{2}\right)^2 I(h_{m-3}^2(\omega)) \leq 2 \frac{(m-1)! \sqrt{\pi}}{2^{m-1}}. \quad (89)$$

So

$$R \leq 2 \left(\frac{\sqrt{2}}{l}\right)^{2m-4} \frac{(m-1)! \sqrt{\pi}}{2^{m-1} (2m-4)!}. \quad (90)$$

Similarly,

$$|I(\omega^2 y_{2m-4}(\omega)) - I(\omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))| = |I(\omega^2 y_{2m-4}(\omega) - \omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))| \leq \quad (91)$$

$$I(|\omega^2 y_{2m-4}(\omega) - \omega^2 \cos(\omega r \frac{\sqrt{2}}{l})|) \leq I\left(\left(\frac{\sqrt{2}}{l}\right)^{2m-4} \frac{\omega^2 h_{m-2}^2(\omega)}{(2m-4)!}\right) = R. \quad (92)$$

Finally, using that $Q_m(\omega^2 y_{2m-4}(\omega)) = I(\omega^2 y_{2m-4}(\omega))$ (because $\omega^2 y_{2m-4}(\omega)$ has degree less than $2m$) and following exactly the same procedure as in the previous proof we get

$$\frac{2}{l^2} |I(\omega^2 \cos(\omega r \frac{\sqrt{2}}{l})) - Q_m(\omega^2 \cos(\omega r \frac{\sqrt{2}}{l}))| \leq \frac{2}{l^2} 2R = 4 \left(\frac{\sqrt{2}}{l}\right)^{2m-2} \frac{(m-1)! \sqrt{\pi}}{2^{m-1} (2m-4)!} \leq \frac{4}{l^2} (m-1) E_{m-2}. \quad (93)$$

Appendix B. Kernel Approximation Additional Plots

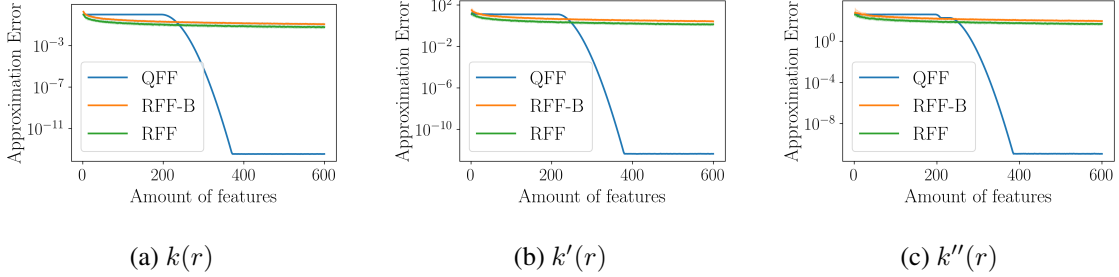


Figure 6: Comparing the maximum error of different feature expansions over $r \in [0, 1]$. For the random feature expansions, we show median as well as 12.5% and 87.5% quantiles over 100 random samples. Due to the exponential decay of the error of the QFF approximation, this stochasticity is barely visible. As given by the theoretical analysis, the error is a bit higher for the derivatives, but still decaying exponentially. In this plot, we set $l = 0.05$.

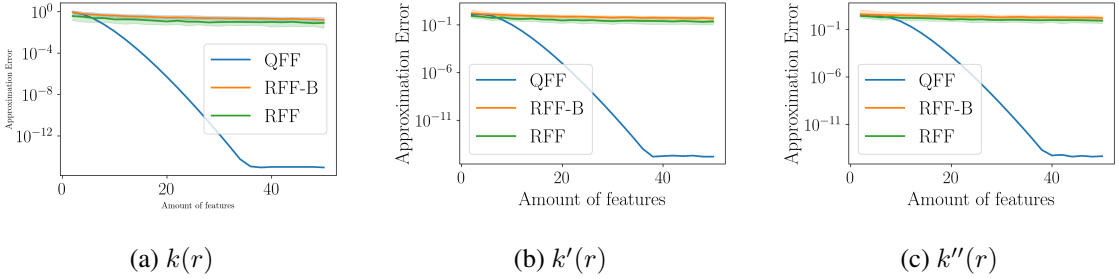


Figure 7: Comparing the maximum error of different feature expansions over $r \in [0, 1]$. For the random feature expansions, we show median as well as 12.5% and 87.5% quantiles over 100 random samples. Due to the exponential decay of the error of the QFF approximation, this stochasticity is barely visible. As given by the theoretical analysis, the error is a bit higher for the derivatives, but still decaying exponentially. In this plot, we set $l = 0.5$.

Appendix C. GP Regression with Derivatives

In this section, we will introduce the necessary notation and the proof for our theoretical results on GP regression with derivatives.

C.1 Notation

Consider the problem of Gaussian Process regression, using zero mean prior and the RBF kernel function $k_\phi(x, y) := \rho e^{-\frac{(x-y)^2}{2l^2}}$ for some fixed hyperparameters $\phi = (\rho, l)$, which denote the variance and the lengthscale. Suppose we are given at n observation points $\mathbf{t} := \mathbf{t}_n = (t_1, \dots, t_n)$ the n dimensional (column) vectors $\mathbf{y} := \mathbf{x} + \epsilon_{\sigma^2}$ and $\mathbf{F} := \dot{\mathbf{x}} + \epsilon_\gamma$ of noisy state and noisy state derivative observations. Our goal is to infer at the observation point T the values of state $x := x(T)$ and state derivative $\dot{x} := \dot{x}(T)$.

We repeat some of the definitions already given in this paper and this appendix as well as add a few new ones relevant for this section only:

Let $\mathbf{k}_\phi(\mathbf{t}, T)$ denote the n dimensional kernel (column) vector, i.e.

$$\mathbf{k}_\phi(\mathbf{t}, T)_i := k_\phi(t_i, T). \quad (94)$$

Let $'k_\phi(x, y)$ denote the partial derivative of k_ϕ w.r.t. its first argument, i.e.

$$'k_\phi(x, y) := \frac{\partial}{\partial a} k_\phi(a, b)|_{a=x, b=y}. \quad (95)$$

Let $'\mathbf{k}_\phi(\mathbf{t}, T)$ denote the n dimensional kernel derivative (column) vector, i.e.

$$' \mathbf{k}_\phi(\mathbf{t}, T)_i := 'k_\phi(t_i, T). \quad (96)$$

Let $k'_\phi(x, y)$ denote the partial derivative of k_ϕ w.r.t. its second argument, i.e.

$$k'_\phi(x, y) := \frac{\partial}{\partial b} k_\phi(a, b)|_{a=x, b=y}. \quad (97)$$

Let $\mathbf{k}'_\phi(\mathbf{t}, T)$ denote the n dimensional kernel derivative (column) vector, i.e.

$$\mathbf{k}'_\phi(\mathbf{t}, T)_i := k'_\phi(t_i, T). \quad (98)$$

Let $k''_\phi(x, y)$ denote the mixed partial derivative of k_ϕ , i.e.

$$k''_\phi(x, y) := \frac{\partial^2}{\partial a \partial b} k_\phi(a, b)|_{a=x, b=y}. \quad (99)$$

Let $\mathbf{k}''_\phi(\mathbf{t}, T)$ denote the n dimensional kernel derivative (column) vector, i.e.

$$\mathbf{k}''_\phi(\mathbf{t}, T)_i := k''_\phi(t_i, T). \quad (100)$$

Let \mathbf{C}_ϕ denote the $n \times n$ covariance kernel matrix, whose elements are given by

$$[\mathbf{C}_\phi]_{i,j} := k_\phi(t_i, t_j). \quad (101)$$

Let $'C_\phi$ denote the kernel derivative matrix, whose elements are given by

$$['C_\phi]_{i,j} := \frac{\partial}{\partial a} k_\phi(a, b)|_{a=t_i, b=t_j}. \quad (102)$$

Let C'_ϕ denote the kernel derivative matrix, whose elements are given by

$$[C'_\phi]_{i,j} := \frac{\partial}{\partial b} k_\phi(a, b)|_{a=t_i, b=t_j}. \quad (103)$$

Let C''_ϕ denote the mixed kernel derivative matrix, whose elements are given by

$$[C''_\phi]_{i,j} := \frac{\partial^2}{\partial a \partial b} k_\phi(a, b)|_{a=t_i, b=t_j}. \quad (104)$$

Let \hat{K}_ϕ denote the sum of the $2n \times 2n$ block matrix with the covariance matrix and its derivatives plus the diagonal noise matrix, i.e.

$$\hat{K}_\phi := \begin{pmatrix} C_\phi & C'_\phi \\ 'C_\phi & C''_\phi \end{pmatrix} + \begin{pmatrix} \sigma^2 \mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & \gamma \mathbb{I}_n \end{pmatrix}. \quad (105)$$

Let $\hat{\mathbf{k}}_\phi(\mathbf{t}, T)$ denote the $2n$ dimensional (column) vector, which is a concatenation of $\mathbf{k}_\phi(\mathbf{t}, T)$ and $'\mathbf{k}_\phi(\mathbf{t}, T)$, i.e.

$$\hat{\mathbf{k}}_\phi(\mathbf{t}, T) := \begin{pmatrix} \mathbf{k}_\phi(\mathbf{t}, T) \\ '\mathbf{k}_\phi(\mathbf{t}, T) \end{pmatrix}. \quad (106)$$

Let $\hat{\mathbf{k}}'_\phi(\mathbf{t}, T)$ denote the $2n$ dimensional (column) vector, which is a concatenation of $\mathbf{k}'_\phi(\mathbf{t}, T)$ and $\mathbf{k}''_\phi(\mathbf{t}, T)$, i.e.

$$\hat{\mathbf{k}}'_\phi(\mathbf{t}, T) := \begin{pmatrix} \mathbf{k}'_\phi(\mathbf{t}, T) \\ \mathbf{k}''_\phi(\mathbf{t}, T) \end{pmatrix}. \quad (107)$$

Finally, we are able to write down the formulas for the scalar predictive mean and covariance at a new point T . Here, we let μ denote the mean of the state, μ' denote the mean of the derivative, Σ denote the variance of the state and Σ' the variance of the derivative prediction. They are given by

$$\mu(T) = \hat{\mathbf{k}}_\phi(\mathbf{t}, T)^T \hat{K}_\phi^{-1} \begin{pmatrix} \mathbf{y} \\ \mathbf{F} \end{pmatrix}, \quad (108)$$

$$\mu'(T) = \hat{\mathbf{k}}'_\phi(\mathbf{t}, T)^T \hat{K}_\phi^{-1} \begin{pmatrix} \mathbf{y} \\ \mathbf{F} \end{pmatrix}, \quad (109)$$

$$\Sigma(T) = k(T, T) - \hat{\mathbf{k}}_\phi(\mathbf{t}, T)^T \hat{K}_\phi^{-1} \hat{\mathbf{k}}_\phi(\mathbf{t}, T), \quad (110)$$

$$\Sigma'(T) = k''(T, T) - \hat{\mathbf{k}}'_\phi(\mathbf{t}, T)^T \hat{K}_\phi^{-1} \hat{\mathbf{k}}'_\phi(\mathbf{t}, T). \quad (111)$$

C.2 Proof of Theorem 3

In this Section, we will prove the main theorem regarding the approximation error of GP regression with derivatives. For ease of reference, we will restate it as Theorem 16. Throughout this section, we will use \tilde{a} to denote the feature approximation of any scalar, vector or matrix a .

Theorem 16 *Let us consider an RBF kernel with hyperparameters (ρ, l) and domain $[0, 1]$. Define $e_{\tilde{\mu}}, e_{\tilde{\Sigma}}, e_{\tilde{\mu}'}, e_{\tilde{\Sigma}'}$ as the absolute error between the feature approximations and the corresponding accurate quantities of the means and covariances of Equations (4) and (5). For each $\tau \in [0, 1]$, define e_{tot} as the maximum of these four errors. Define $c := \min(\gamma, \sigma^2)$ and $R := \max(\|\mathbf{y}\|_\infty, \|\mathbf{F}\|_\infty)$. Let $C > 0$. Let us consider a QFF approximation scheme of order $m \geq 3 + \max\left(\frac{e}{2l^2}, \log\left(\frac{270n^2\rho^3R}{l^8c^2C}\right)\right)$. Then, it holds for all $\tau \in [0, 1]$ that $e_{\text{tot}} \leq C$.*

Proof 6 *Suppose that we apply a QFF approximation scheme of order m for the functions $k_\phi, {}^l k_\phi, k'_\phi, k''_\phi$, which gives us a deterministic and uniform (over their domain) approximation guarantee of absolute error less than $\epsilon := \epsilon_\phi(m)$ (for any of them). W.l.o.g we can assume that the domain of these functions is $[0, 1]^2$, $l \leq 1$ and $\rho \geq 1$ (so also $0 \leq T, t_1, \dots, t_n \leq 1$). Moreover, we assume that $\|\mathbf{y}\|_{\max}, \|\mathbf{F}\|_{\max} \leq R$, for some positive constant R .*

Let \mathbf{E}_1 be the error (matrix) when approximating $\hat{\mathbf{K}}_\phi$, i.e.

$$\mathbf{E}_1 := \hat{\mathbf{K}}_\phi - \tilde{\mathbf{K}}_\phi. \quad (112)$$

Let \mathbf{E}_2 be the error (matrix) when approximating $\hat{\mathbf{K}}_\phi^{-1}$, i.e.

$$\mathbf{E}_2 := \hat{\mathbf{K}}_\phi^{-1} - \tilde{\mathbf{K}}_{\phi^{-1}}. \quad (113)$$

Let \mathbf{e}_1 be the error (vector) when approximating $\hat{\mathbf{k}}_\phi(\mathbf{t}, T)$, i.e.

$$\mathbf{e}_1 := \hat{\mathbf{k}}_\phi(\mathbf{t}, T) - \tilde{\mathbf{k}}_\phi(\mathbf{t}, T). \quad (114)$$

Let \mathbf{e}_2 be the error (vector) when approximating $\hat{\mathbf{k}}_\phi(\mathbf{t}, T)$, i.e.

$$\mathbf{e}_2 := \hat{\mathbf{k}}'_\phi(\mathbf{t}, T) - \tilde{\mathbf{k}}'_\phi(\mathbf{t}, T). \quad (115)$$

Given these assumptions and definitions, we can introduce the following bounds:

$$\|\mathbf{e}_1\| \leq \sqrt{2n}\epsilon, \quad (116)$$

$$\|\mathbf{e}_2\| \leq \sqrt{2n}\epsilon, \quad (117)$$

$$\sigma_1(\mathbf{E}_1) \leq 2n\epsilon, \quad (118)$$

$$\left\| \begin{pmatrix} \mathbf{y} \\ \mathbf{F} \end{pmatrix} \right\| \leq \sqrt{2n}R. \quad (119)$$

Using the fact that $k_\phi \leq \rho$, ${}^l k_\phi \leq \frac{\rho}{l}$ and $k''_\phi \leq \frac{2\rho}{l^2}$, we can also get

$$\left\| \hat{\mathbf{k}}_\phi(\mathbf{t}, T) \right\| \leq \sqrt{2n} \frac{\rho}{l}, \quad (120)$$

$$\left\| \hat{\mathbf{k}}'_\phi(\mathbf{t}, T) \right\| \leq \sqrt{2n} \frac{2\rho}{l^2}. \quad (121)$$

We know that the matrix $\begin{pmatrix} \mathbf{C}_\phi & \mathbf{C}'_\phi \\ {}^l \mathbf{C}_\phi & \mathbf{C}''_\phi \end{pmatrix}$ and its QFF approximation are symmetric positive semi-definite and thus, their smallest singular value is non negative. Also, $\begin{pmatrix} \sigma^2 \mathbb{I}_n & \mathbf{0} \\ \mathbf{0} & \gamma \mathbb{I}_n \end{pmatrix}$ has smallest

singular value c . Thus, both $\hat{\mathbf{K}}$ and $\tilde{\mathbf{K}}$ are symmetric, positive definite matrices and we have

$$\sigma_1(\hat{\mathbf{K}}_\phi^{-1}) \leq \frac{1}{c}, \quad (122)$$

$$\sigma_1(\tilde{\mathbf{K}}_\phi^{-1}) \leq \frac{1}{c}. \quad (123)$$

Finally, we can use the Woodbury identity to get $\mathbf{E}_2 = -\hat{\mathbf{K}}_\phi^{-1} \mathbf{E}_1 \tilde{\mathbf{K}}_{\phi^{-1}}$, justifying

$$\sigma_1(\mathbf{E}_2) \leq \frac{2n}{c^2} \epsilon. \quad (124)$$

We now have all the building blocks needed to provide a bound for $e_{\tilde{\mu}}$: Using all the above notations, assumptions and results we get

$$|\epsilon_\mu| = \left| \hat{\mathbf{k}}_\phi(\mathbf{t}, T)^T \hat{\mathbf{K}}_\phi^{-1} \begin{pmatrix} \mathbf{y} \\ \mathbf{F} \end{pmatrix} - (\hat{\mathbf{k}}_\phi(\mathbf{t}, T) - \mathbf{e}_1)^T (\hat{\mathbf{K}}_\phi^{-1} - \mathbf{E}_1) \begin{pmatrix} \mathbf{y} \\ \mathbf{F} \end{pmatrix} \right| \quad (125)$$

$$= \left| (\mathbf{e}_1^T \hat{\mathbf{K}}_\phi^{-1} + \hat{\mathbf{k}}_\phi(\mathbf{t}, T)^T \mathbf{E}_1 - \mathbf{e}_1^T \mathbf{E}_1) \begin{pmatrix} \mathbf{y} \\ \mathbf{F} \end{pmatrix} \right| \quad (126)$$

$$\leq \sqrt{2n}R \left(\frac{\sqrt{2n}}{c} \epsilon + \frac{\rho}{l} \sqrt{2n} 2n\epsilon + \sqrt{2n} 2n\epsilon^2 \right) \quad (127)$$

$$\leq 10 \frac{n^2 R \rho}{lc} \epsilon. \quad (128)$$

Similarly, we get

$$|\epsilon_{\mu'}| = \left| (\mathbf{e}_2^T \hat{\mathbf{K}}_\phi^{-1} + \hat{\mathbf{k}}'_\phi(\mathbf{t}, T)^T \mathbf{E}_1 - \mathbf{e}_2^T \mathbf{E}_1) \begin{pmatrix} \mathbf{y} \\ \mathbf{F} \end{pmatrix} \right| \quad (129)$$

$$\leq \sqrt{2n}R \left(\frac{\sqrt{2n}}{c} \epsilon + \frac{\rho}{l^2} \sqrt{2n} 4n\epsilon + \sqrt{2n} 2n\epsilon^2 \right) \quad (130)$$

$$\leq 14 \frac{n^2 R \rho}{l^2 c} \epsilon. \quad (131)$$

Moreover,

$$|\epsilon_\Sigma| = |k(T, T) - \hat{\mathbf{k}}_\phi(\mathbf{t}, T)^T \hat{\mathbf{K}}_\phi^{-1} \hat{\mathbf{k}}_\phi(\mathbf{t}, T) - \tilde{k}(T, T)| \quad (132)$$

$$+ |(\hat{\mathbf{k}}_\phi(\mathbf{t}, T) - \mathbf{e}_1)^T (\hat{\mathbf{K}}_\phi^{-1} - \mathbf{E}_1) (\hat{\mathbf{k}}_\phi(\mathbf{t}, T) - \mathbf{e}_1)| \quad (133)$$

$$\leq \epsilon + |-2\mathbf{e}_1^T \tilde{\mathbf{K}}_{\phi^{-1}} \hat{\mathbf{k}}_\phi(\mathbf{t}, T) + \mathbf{e}_1^T \tilde{\mathbf{K}}_{\phi^{-1}} \mathbf{e}_1 - \hat{\mathbf{k}}_\phi(\mathbf{t}, T)^T \mathbf{E}_2 \hat{\mathbf{k}}_\phi(\mathbf{t}, T)| \quad (134)$$

$$\leq \epsilon + \frac{4n\rho}{lc} \epsilon + \frac{2n}{c} \epsilon^2 + \frac{4n^2 \rho^2}{l^2 c^2} \epsilon \quad (135)$$

$$\leq \frac{14n^2 \rho^2}{l^2 c^2} \epsilon \quad (136)$$

and

$$|\epsilon_{\Sigma'}| \leq \epsilon + |-2e_2^T \tilde{\mathbf{K}}_\phi^{-1} \hat{\mathbf{k}}'_\phi(t, T) + e_2^T \tilde{\mathbf{K}}_\phi^{-1} e_2 - \hat{\mathbf{k}}'_\phi(t, T)^T \mathbf{E}_2 \hat{\mathbf{k}}'_\phi(t, T)| \quad (137)$$

$$\leq \epsilon + \frac{8n\rho}{l^2 c} \epsilon + \frac{2n}{c} \epsilon^2 + \frac{16n^2 \rho^2}{l^4 c^2} \epsilon \quad (138)$$

$$\leq \frac{27n^2 \rho^2}{l^4 c^2} \epsilon. \quad (139)$$

To summarize all of these upper bounds, we can observe that

$$|e_{tot}| \leq \frac{27n^2 \rho^2 R}{l^4 c^2} \epsilon. \quad (140)$$

Let us now fix a constant $0 < C < 1$ and define $M := m - 3$. Let us choose M such that

$$M \geq \max\left(\frac{e}{2l^2}, \log\left(\frac{270n^2 \rho^3 R}{l^8 c^2 C}\right)\right). \quad (141)$$

With this choice, it holds that

$$\left(\frac{e}{4l^2 M}\right)^M \leq \frac{l^8 c^2}{270n^2 \rho^3 R} C \quad (142)$$

$$\implies \frac{2e\sqrt{\pi}\rho}{l^4} \left(\frac{e}{4l^2 M}\right)^M \leq \frac{l^4 c^2}{27n^2 \rho^2 R} C \quad (143)$$

$$\implies \epsilon \leq \frac{l^4 c^2}{27n^2 \rho^2 R} C \quad (144)$$

$$\iff \frac{27n^2 \rho^2 R}{l^4 c^2} \epsilon \leq C, \quad (145)$$

$$(146)$$

which concludes the proof of this theorem.

The above bound could be reformulated as $m \geq 12 + \max\left(\frac{e}{2l^2}, \log\left(\frac{n^2 \rho^3 R}{l^8 c^2 C}\right)\right)$. Moreover, we assumed that $R \geq 1$ and that $c \leq 1$. If any of these conditions is not met, then the bounds are still valid if we substitute these quantities by 1 (the same holds also for ρ and l). Finally, we implicitly assumed that ϵ , the uniform upper bound of the approximation for $k_\phi, k'_\phi, k''_\phi$ is smaller than 1. This happens when $m \geq 3 + \max\left(\frac{e}{2l^2}, \log\left(\frac{10\rho}{l^4}\right)\right)$, a condition which is met if $m \geq 3 + \max\left(\frac{e}{2l^2}, \log\left(\frac{270n^2 \rho^3 R}{l^8 c^2 C}\right)\right)$.

Appendix D. Additional Empirical Evaluation GPR

D.1 Lotka Volterra

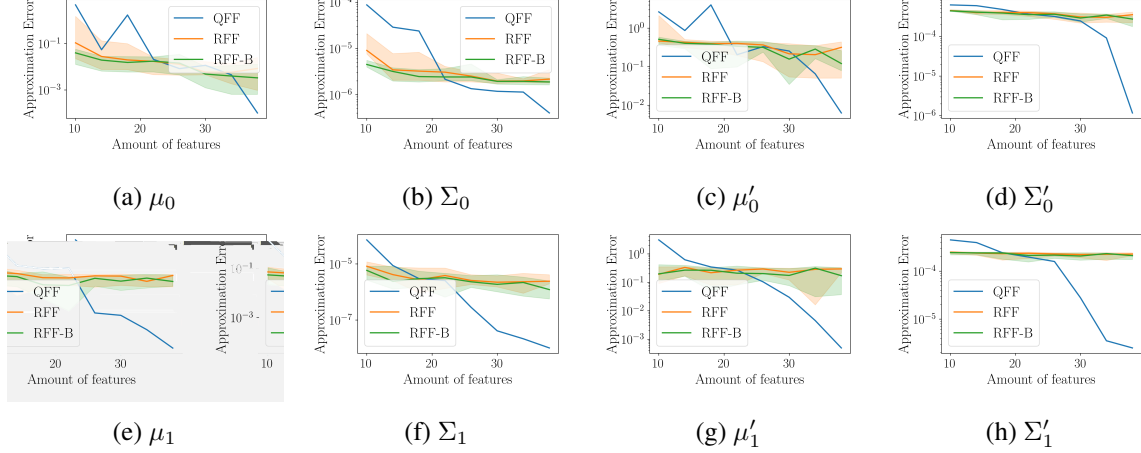


Figure 8: Approximation error of the different feature approximations compared to the accurate GP, evaluated at $t = 1.75$ for the Lotka Volterra system with 1000 observations and $\sigma^2 = 0.1$. For each feature, we show the median as well as the 12.5% and 87.5% quantiles over 10 independent noise realizations, separately for each state dimension.

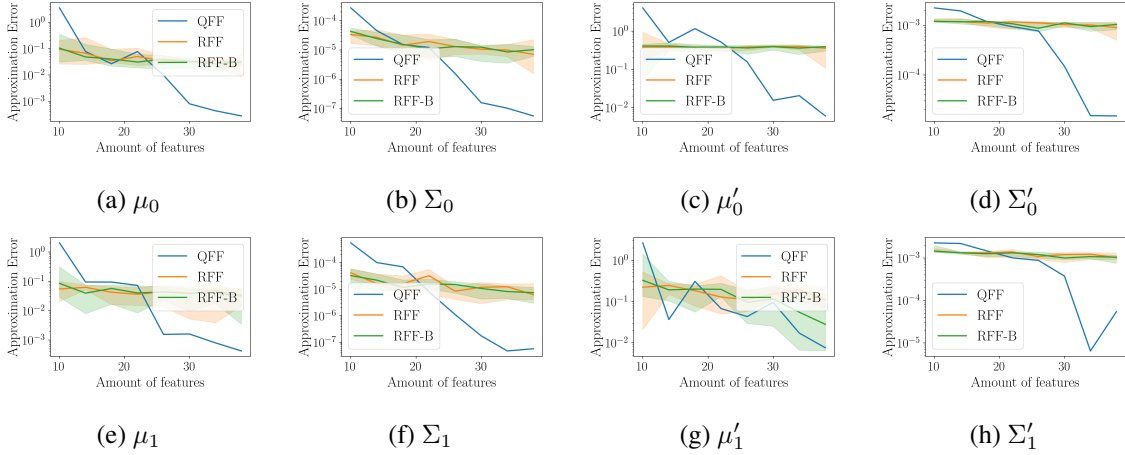


Figure 9: Approximation error of the different feature approximations compared to the accurate GP, evaluated at $t = 1.75$ for the Lotka Volterra system with 1000 observations and $\sigma^2 = 0.5$. For each feature, we show the median as well as the 12.5% and 87.5% quantiles over 10 independent noise realizations, separately for each state dimension.

D.2 Protein Transduction

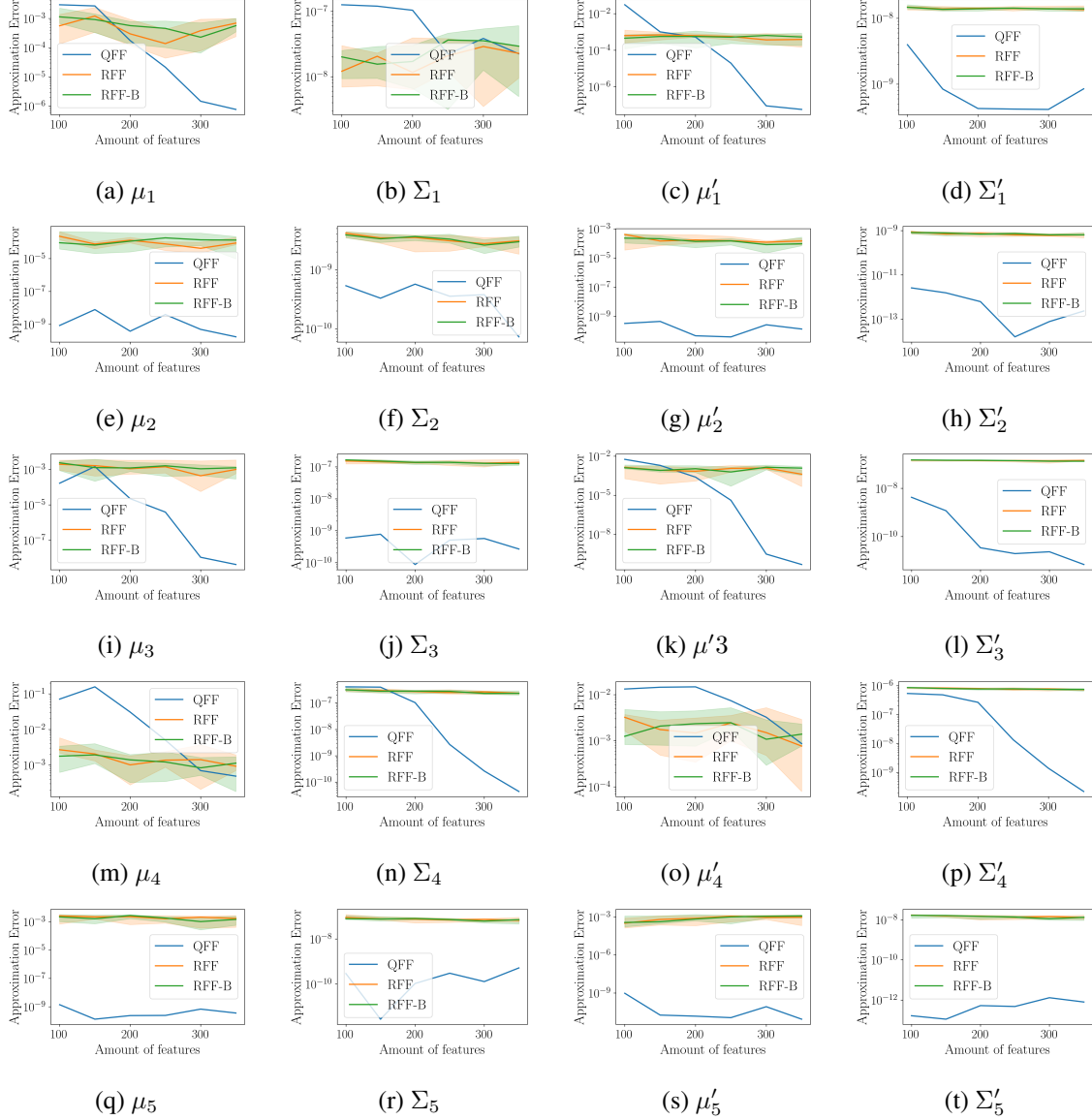


Figure 10: Approximation error of the different feature approximations compared to the accurate GP, evaluated at $t = 30$ for the Protein Transduction system with 1000 observations and $\sigma^2 = 0.0001$. For each feature, we show the median as well as the 12.5% and 87.5% quantiles over 10 independent noise realizations, separately for each state dimension.

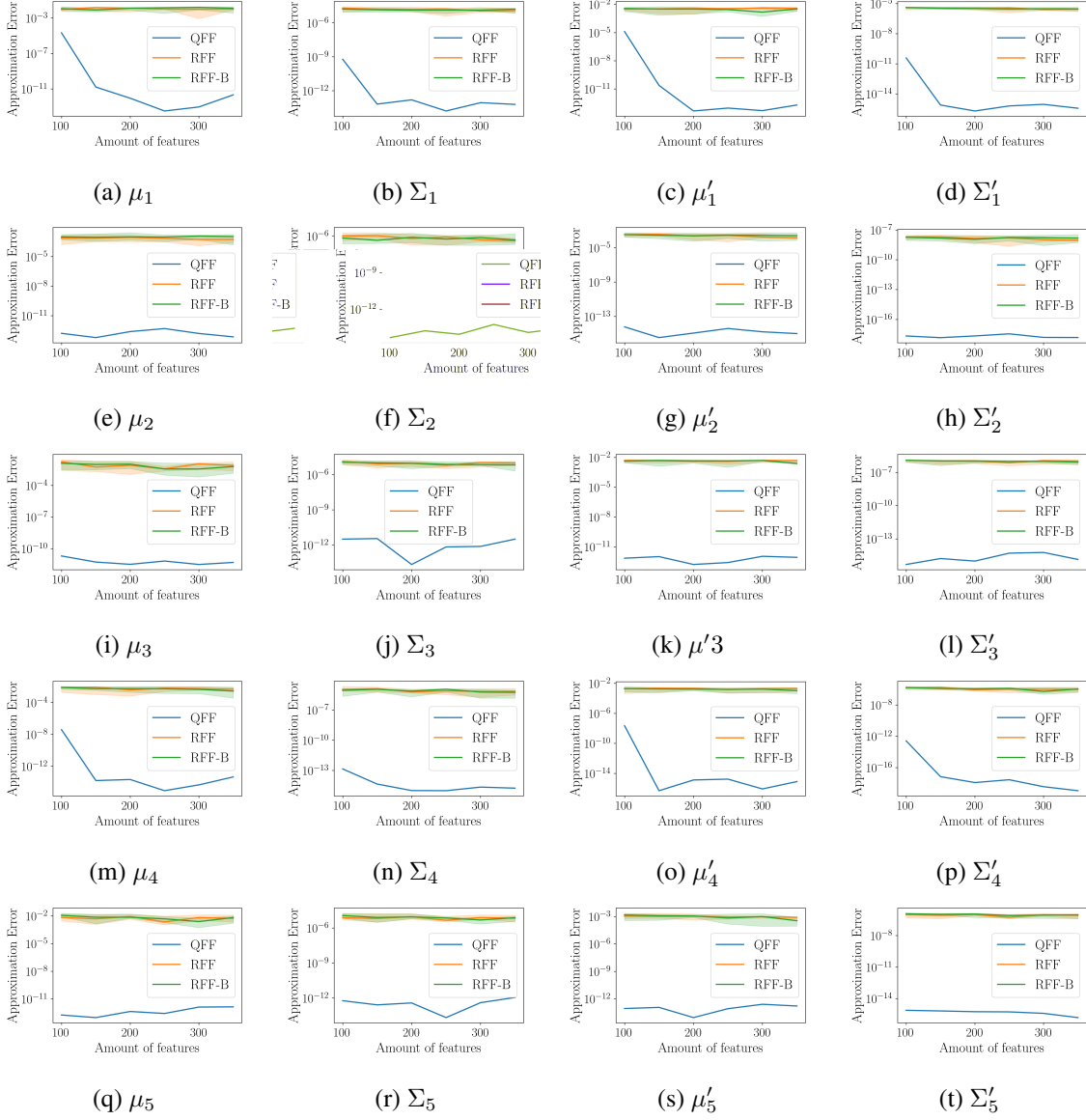


Figure 11: Approximation error of the different feature approximations compared to the accurate GP, evaluated at $t = 30$ for the Protein Transduction system with 1000 observations and $\sigma^2 = 0.01$. For each feature, we show the median as well as the 12.5% and 87.5% quantiles over 10 independent noise realizations, separately for each state dimension.

D.3 Lorenz

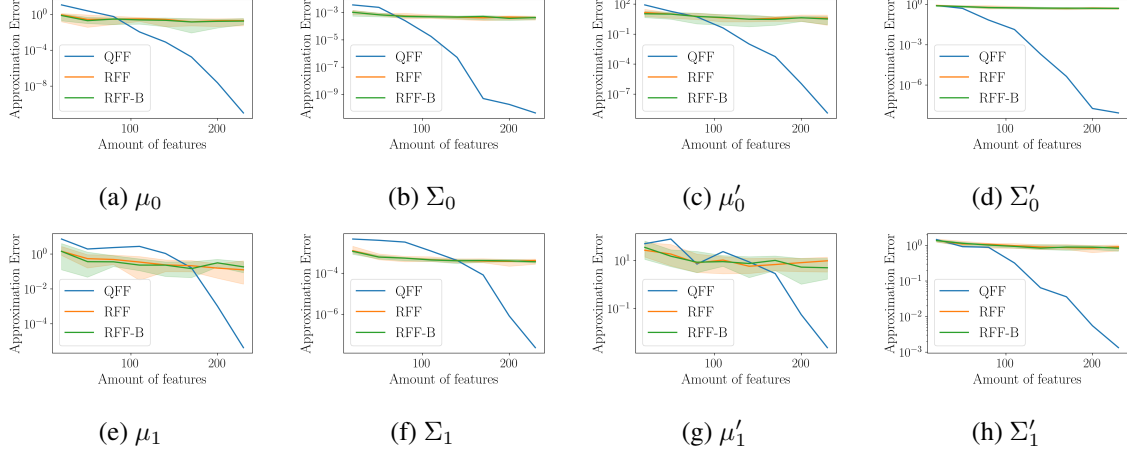


Figure 12: Approximation error of the different feature approximations compared to the accurate GP, evaluated at $t = 0.8$ for the Lorenz system with 1000 observations and an SNR of 100. For each feature, we show the median as well as the 12.5% and 87.5% quantiles over 10 independent noise realizations, separately for each state dimension.

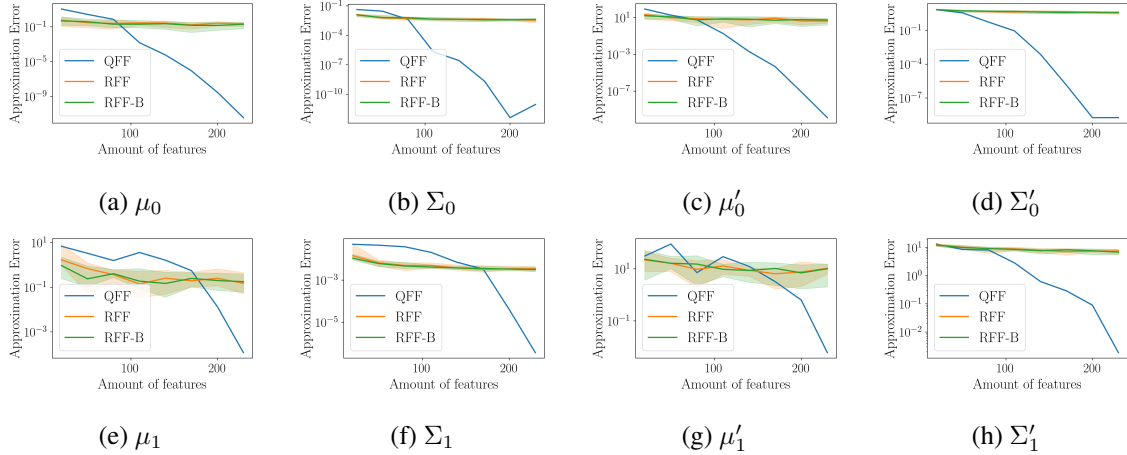


Figure 13: Approximation error of the different feature approximations compared to the accurate GP, evaluated at $t = 0.8$ for the Lorenz system with 1000 observations and an SNR of 10. For each feature, we show the median as well as the 12.5% and 87.5% quantiles over 10 independent noise realizations, separately for each state dimension.

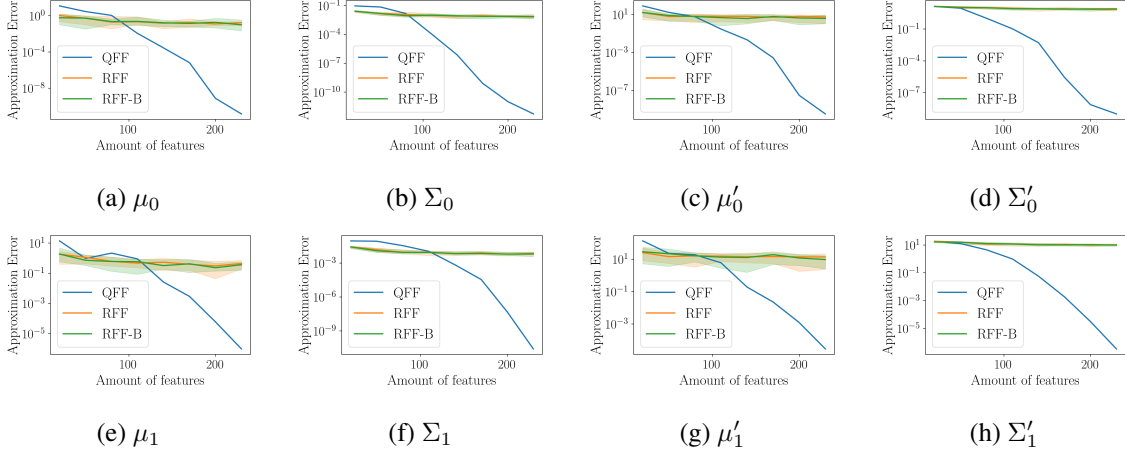


Figure 14: Approximation error of the different feature approximations compared to the accurate GP, evaluated at $t = 0.8$ for the Lorenz system with 1000 observations and an SNR of 5. For each feature, we show the median as well as the 12.5% and 87.5% quantiles over 10 independent noise realizations, separately for each state dimension.

Appendix E. Risk Approximation Error Bounds

Let

$$\mathcal{R}_{\lambda\gamma\phi}(\mathbf{x}, \boldsymbol{\theta}, \mathbf{y}) = \mathbf{x}^T (\mathbf{C}_\phi + \lambda \mathbf{I})^{-1} \mathbf{x} \quad (147)$$

$$+ (\mathbf{x} - \mathbf{y})^T \sigma^{-2} (\mathbf{x} - \mathbf{y}) \quad (148)$$

$$+ (\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D}\mathbf{x})^T (\mathbf{A} + \gamma \mathbf{I})^{-1} (\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D}\mathbf{x}) \quad (149)$$

where $\mathbf{C}_\phi[i, j] = \rho e^{-\frac{r_{ij}^2}{2l^2}}$ for some fixed hyperparameters $\phi = (\rho, l)$, which denote the variance and the lengthscale. Let n be the number of data points and m be the order of the Quadrature scheme used to approximate the kernel. By writing $\tilde{\mathbf{C}}_\phi$, $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{D}}$ for the approximated quantities as described in Section ??, we get the approximate risk function

$$\tilde{\mathcal{R}}_{\lambda\gamma\phi}(\mathbf{x}, \boldsymbol{\theta}) := \mathbf{x}^T (\tilde{\mathbf{C}}_\phi + \lambda \mathbf{I})^{-1} \mathbf{x} \quad (150)$$

$$+ (\mathbf{x} - \mathbf{y})^T \sigma^{-2} (\mathbf{x} - \mathbf{y}) \quad (151)$$

$$+ (\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \tilde{\mathbf{D}}\mathbf{x})^T (\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-1} (\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \tilde{\mathbf{D}}\mathbf{x}). \quad (152)$$

To prove Theorem 4, we show that $\tilde{\mathcal{R}}$ converges to \mathcal{R} in the relative error sense: As m increases, the relative error $\frac{|\mathcal{R} - \tilde{\mathcal{R}}|}{\mathcal{R}}$ will become arbitrarily small. Theorem 4 is restated here as Theorem 17.

Theorem 17 *Let \mathcal{R} and $\tilde{\mathcal{R}}$ be defined as above. The parameters λ and γ , the kernel hyperparameters $\phi = (\rho, l)$ and the number of data points n are considered fixed. We assume $n \geq 60$. Consider $1 > \epsilon > 0$. If m , the order of the quadrature scheme, is at least*

$$m \geq 10 + \max\left\{\frac{e}{2l^2}, \log_2\left(\frac{\rho^2 n^3}{\lambda^2 \gamma l^4 \epsilon}\right)\right\} \quad (153)$$

then we have

$$\frac{|\mathcal{R}_{\lambda\gamma\phi}(\mathbf{x}, \boldsymbol{\theta}) - \tilde{\mathcal{R}}_{\lambda\gamma\phi}(\mathbf{x}, \boldsymbol{\theta})|}{\mathcal{R}_{\lambda\gamma\phi}(\mathbf{x}, \boldsymbol{\theta})} \leq \epsilon \quad (154)$$

for any configuration of the variables \mathbf{x} and $\boldsymbol{\theta}$.

The lengthscale corresponds to time observations normalized in the $[0, 1]$ interval. In order to make things a bit less complicated in the calculations, we shall assume that $\gamma \leq 1, \lambda \leq 1, \rho \geq 1, l \leq 1$ (and specifically $l \leq \frac{\epsilon}{4}$). If any of these assumptions is violated, we can just substitute the corresponding parameter with 1 in the previous bound, and the resulting bound will still be valid. Moreover, the logarithm is the binary one and we will simply use \log from now on. Before we prove the theorem, we will introduce some notations and some preliminary results.

Let $\|\mathbf{K}\|_F = \sqrt{\text{tr}(\mathbf{K}\mathbf{K}^T)}$ be the Frobenius norm of a matrix \mathbf{K} , $\|\mathbf{K}\|_{\max} = \max_{i,j} |\mathbf{K}_{ij}|$ the max norm of \mathbf{K} and $\sigma_1(\mathbf{K}) := \|\mathbf{K}\|_2$ the spectral norm of \mathbf{K} (which is by definition the largest singular value). It holds that $\|\mathbf{K}\mathbf{x}\| \leq \sigma_1(\mathbf{K})\|\mathbf{x}\|$. Moreover, for a $n \times n$ matrix \mathbf{K} we have

$$\sigma_1(\mathbf{K}) \leq \sqrt{\sum \sigma_i^2(\mathbf{K})} = \sqrt{\|\mathbf{K}\|_F^2} \leq \sqrt{n^2 \|\mathbf{K}\|_{\max}^2} = n\|\mathbf{K}\|_{\max} \quad (155)$$

Specifically, if $C_\phi[i, j] = \rho e^{-\frac{r_{ij}^2}{2l^2}}$, then

$$\sigma_1(\mathbf{C}_\phi) \leq \rho n \quad (156)$$

while

$$\sigma_1(\mathbf{C}'_\phi) \leq \frac{\rho}{l} n \quad (157)$$

since $\mathbf{C}'_\phi[i, j] = -\frac{\rho}{l} \frac{r_{ij}}{l} e^{-\frac{r_{ij}^2}{2l^2}}$ as $xe^{-\frac{x^2}{2}} \leq \frac{1}{\sqrt{e}} \leq 1$.

From Woodbury's identity for matrix inversion, for invertible \mathbf{K} and $\mathbf{K} + \mathbf{E}$ we have:

$$(\mathbf{K} + \mathbf{E})^{-1} = \mathbf{K}^{-1} - \mathbf{K}^{-1}(\mathbf{I} + \mathbf{E}\mathbf{K}^{-1})^{-1}\mathbf{E}\mathbf{K}^{-1} \quad (158)$$

Let $\mathbf{E}_1 := \mathbf{E}_1(m) = \tilde{\mathbf{C}}_\phi - \mathbf{C}_\phi$. We know that $\|\mathbf{E}_1\|_{\max} \leq \sqrt{\frac{\pi}{2}} \frac{1}{m^m} (\frac{e}{4l^2})^m$ so for $a := \frac{e}{4l^2}$ we have

$$\sigma_1(\mathbf{E}_1) \leq 2\left(\frac{a}{m}\right)^m n \quad (159)$$

Now let us bound the first term of the risk

Lemma 18 Consider the term $\mathbf{x}^T(\mathbf{C}_\phi + \lambda\mathbf{I})^{-1}\mathbf{x}$ approximated by $\mathbf{x}^T(\tilde{\mathbf{C}}_\phi + \lambda\mathbf{I})^{-1}\mathbf{x}$. Then for $m \geq M := 7 + \max\{\frac{e}{2l^2}, \log(\frac{\rho^2 n^3}{\lambda^2 \gamma l^4})\}$ we get

$$|\mathbf{x}^T(\mathbf{C}_\phi + \lambda\mathbf{I})^{-1}\mathbf{x} - \mathbf{x}^T(\tilde{\mathbf{C}}_\phi + \lambda\mathbf{I})^{-1}\mathbf{x}| \leq \epsilon_1 \mathbf{x}^T(\mathbf{C}_\phi + \lambda\mathbf{I})^{-1}\mathbf{x} \quad (160)$$

where $\epsilon_1 = \frac{4n}{\lambda} \left(\frac{a}{m}\right)^m$.

Proof 7 From Woodbury's inversion formula (158), if we set $\mathbf{K} = \mathbf{C}_\phi + \lambda\mathbf{I}$ and $\mathbf{E} = \mathbf{E}_1$ we have $\sigma_1(\mathbf{K}^{-1}) < \frac{1}{\lambda}$. Regarding the matrix $\mathbf{I} + \mathbf{E}\mathbf{K}^{-1}$, we would like to find a lower bound on the smallest singular value $\sigma_n(\mathbf{I} + \mathbf{E}\mathbf{K}^{-1})$. For any vector u with $|u| = 1$ we have

$$|(\mathbf{I} + \mathbf{E}\mathbf{K}^{-1})u| \geq |u| - |\mathbf{E}\mathbf{K}^{-1}u| \geq 1 - \sigma_1(\mathbf{E})\sigma_1(\mathbf{K}^{-1}) \geq 1 - \frac{2n}{\lambda} \left(\frac{a}{m}\right)^m \geq \frac{1}{2} \quad (161)$$

for $m \geq \max\{\frac{e}{2l^2}, \log(\frac{4n}{\lambda})\}$, which is fulfilled if $m \geq M$. So $\sigma_n(\mathbf{I} + \mathbf{E}\mathbf{K}^{-1}) \geq \frac{1}{2}$ and consequently $\sigma_1((\mathbf{I} + \mathbf{E}\mathbf{K}^{-1})^{-1}) \leq 2$. Moreover, $\sigma_1(\mathbf{K}^{-1}(\mathbf{I} + \mathbf{E}\mathbf{K}^{-1})^{-1}\mathbf{E}\mathbf{K}^{-1}) \leq \sigma_1(\mathbf{K}^{-1})^2\sigma_1((\mathbf{I} + \mathbf{E}\mathbf{K}^{-1})^{-1})\sigma_1(\mathbf{E}) \leq \frac{4n}{\lambda^2}(\frac{a}{m})^m$. So

$$\sigma_1((\mathbf{C}_\phi + \lambda\mathbf{I})^{-1} - (\tilde{\mathbf{C}}_\phi + \lambda\mathbf{I})^{-1}) \leq \frac{4n}{\lambda^2}(\frac{a}{m})^m \quad (162)$$

In total we have

$$|\mathbf{x}^T(\mathbf{C}_\phi + \lambda\mathbf{I})^{-1}\mathbf{x} - \mathbf{x}^T(\tilde{\mathbf{C}}_\phi + \lambda\mathbf{I})^{-1}\mathbf{x}| = |\mathbf{x}^T\mathbf{K}^{-1}(\mathbf{I} + \mathbf{E}\mathbf{K}^{-1})^{-1}\mathbf{E}\mathbf{K}^{-1}\mathbf{x}| \leq \quad (163)$$

$$|\mathbf{x}^T\mathbf{K}^{-\frac{1}{2}}|\sigma_1(\mathbf{K}^{-\frac{1}{2}}(\mathbf{I} + \mathbf{E}\mathbf{K}^{-1})^{-1}\mathbf{E}\mathbf{K}^{-\frac{1}{2}})|\mathbf{x}^T\mathbf{K}^{-\frac{1}{2}}| \leq \quad (164)$$

$$|\mathbf{x}^T\mathbf{K}^{-\frac{1}{2}}||\mathbf{x}^T\mathbf{K}^{-\frac{1}{2}}|\sigma_1(\mathbf{K}^{-\frac{1}{2}})^2\sigma_1((\mathbf{I} + \mathbf{E}\mathbf{K}^{-1})^{-1})\sigma_1(\mathbf{E}) \leq \quad (165)$$

$$|\mathbf{x}^T(\mathbf{C}_\phi + \lambda\mathbf{I})^{-1}\mathbf{x}|\frac{4n}{\lambda}(\frac{a}{m})^m \quad (166)$$

Now we define $\mathbf{E}_2 = \mathbf{E}_2(m) := \tilde{\mathbf{A}} - \mathbf{A}$ and $\mathbf{E}_3 = \mathbf{E}_3(m) := \tilde{\mathbf{D}} - \mathbf{D}$. We need to find an upper bound for the spectral norm of these two error terms.

Lemma 19 Consider the error terms \mathbf{E}_2 and \mathbf{E}_3 as defined above. Then for $m \geq M := 7 + \max\{\frac{e}{2l^2}, \log(\frac{\rho^2 n^3}{\lambda^2 \gamma l^4})\}$ we get

$$\sigma_1(\mathbf{E}_2(m+3)) \leq 20\frac{\rho^2 n^3 a^2}{\lambda^2}(\frac{a}{m})^m \quad \sigma_1(\mathbf{E}_3(m+3)) \leq 10\frac{n^2 \rho a}{\lambda^2}(\frac{a}{m})^m \quad (167)$$

Proof 8 We have

$$\mathbf{E}_2 = \tilde{\mathbf{C}}_\phi'' - \mathbf{C}_\phi'' + {}'\tilde{\mathbf{C}}_\phi(\tilde{\mathbf{C}}_\phi + \lambda\mathbf{I})^{-1}\tilde{\mathbf{C}}_\phi' - {}'\mathbf{C}_\phi(\mathbf{C}_\phi + \lambda\mathbf{I})^{-1}\mathbf{C}_\phi' \quad (168)$$

and we define $\mathbf{E}_{21} := \tilde{\mathbf{C}}_\phi' - \mathbf{C}_\phi'$, $\mathbf{E}_{22} := (\tilde{\mathbf{C}}_\phi + \lambda\mathbf{I})^{-1} - (\mathbf{C}_\phi + \lambda\mathbf{I})^{-1}$ and $\mathbf{E}_{23} := \tilde{\mathbf{C}}_\phi'' - \mathbf{C}_\phi''$. We know that $|\mathbf{E}_{21}(m+2)|_{\max} \leq 16a(\frac{a}{m})^m$ so $\sigma_1(\mathbf{E}_{21}(m+3)) \leq 16a(\frac{a}{m})^m n$. Moreover, we know $|\mathbf{E}_{23}(m+3)|_{\max} \leq 32a^2(\frac{a}{m})^m$ so $\sigma_1(\mathbf{E}_{23}(m+3)) = 32a^2(\frac{a}{m})^m n$. Finally, $\sigma_1(\mathbf{E}_{22}(m)) \leq \frac{4n}{\lambda^2}(\frac{a}{m})^m$ (by (162)). So

$$\sigma_1(\mathbf{E}_2(m+3)) = \quad (169)$$

$$\sigma_1(\mathbf{E}_{23} + ({}'\mathbf{C}_\phi + \mathbf{E}_{21}^T)((\mathbf{C}_\phi + \lambda\mathbf{I})^{-1} + \mathbf{E}_{22})(\mathbf{C}_\phi' + \mathbf{E}_{21}) - {}'\mathbf{C}_\phi(\mathbf{C}_\phi + \lambda\mathbf{I})^{-1}\mathbf{C}_\phi') \leq \quad (170)$$

$$\sigma_1(\mathbf{E}_{23}) + 2\sigma_1({}'\mathbf{C}_\phi)\sigma_1((\mathbf{C}_\phi + \lambda\mathbf{I})^{-1})\sigma_1(\mathbf{E}_{21}) + \sigma_1(\mathbf{E}_{21})^2\sigma_1((\mathbf{C}_\phi + \lambda\mathbf{I})^{-1}) + \quad (171)$$

$$\sigma_1({}'\mathbf{C}_\phi)^2\sigma_1(\mathbf{E}_{22}) + 2\sigma_1(\mathbf{E}_{21})\sigma_1(\mathbf{E}_{22})\sigma_1({}'\mathbf{C}_\phi) + \sigma_1(\mathbf{E}_{21})^2\sigma_1(\mathbf{E}_{22}) \leq \quad (172)$$

$$20\frac{\rho^2 n^3 a^2}{\lambda^2}(\frac{a}{m})^m \quad (173)$$

Now we will prove for each summand in the last sum that it is at most $\frac{\rho^2 n^3 a^2}{\lambda^2} (\frac{a}{m})^m$ times a constant. We will use the fact that $m \geq M$. Indeed:

$$\sigma_1(\mathbf{E}_{23}) \leq 32a^2 (\frac{a}{m})^m n \leq \frac{1}{2} \frac{\rho^2 n^3 a^2}{\lambda^2} (\frac{a}{m})^m \quad (174)$$

$$2\sigma_1('C_\phi)\sigma_1((C_\phi + \lambda I)^{-1})\sigma_1(\mathbf{E}_{21}) \leq 32 \frac{\rho a}{l\lambda} (\frac{a}{m})^m n^2 \leq \frac{\rho^2 n^3 a^2}{\lambda^2} (\frac{a}{m})^m \quad (175)$$

$$\sigma_1(\mathbf{E}_{21})^2 \sigma_1((C_\phi + \lambda I)^{-1}) \leq 16^2 \frac{a^2 n^2}{\lambda} (\frac{a}{m})^{2m} \leq \frac{1}{2} \frac{\rho^2 n^3 a^2}{\lambda^2} (\frac{a}{m})^m \quad (176)$$

$$\sigma_1('C_\phi)^2 \sigma_1(\mathbf{E}_{22}) \leq 16 \frac{\rho^2 n^3 a}{\lambda^2} (\frac{a}{m})^m \leq 16 \frac{\rho^2 n^3 a^2}{\lambda^2} (\frac{a}{m})^m \quad (177)$$

$$2\sigma_1(\mathbf{E}_{21})\sigma_1(\mathbf{E}_{22})\sigma_1('C_\phi) \leq 16a (\frac{a}{m})^m n \frac{4n}{\lambda^2} (\frac{a}{m})^m \frac{n\rho}{l} \leq 2^9 \frac{a^2 n^2 \rho^2}{\lambda^2} (\frac{a}{m})^{2m} \leq \frac{\rho^2 n^3 a^2}{\lambda^2} (\frac{a}{m})^m \quad (178)$$

$$\sigma_1(\mathbf{E}_{21})^2 \sigma_1(\mathbf{E}_{22}) \leq 2^{10} \frac{a^2 n^3}{\lambda^2} (\frac{a}{m})^{3m} \leq \frac{\rho^2 n^3 a^2}{\lambda^2} (\frac{a}{m})^m \quad (179)$$

Similarly,

$$\mathbf{E}_3 = ' \tilde{C}_\phi (\tilde{C}_\phi + \lambda I)^{-1} - ' C_\phi (C_\phi + \lambda I)^{-1} = ' C_\phi \mathbf{E}_{22} + \mathbf{E}_{21}^T (C_\phi + \lambda I)^{-1} + \mathbf{E}_{21}^T \mathbf{E}_{22} \quad (180)$$

Using again that $m \geq M$ we have:

$$\sigma_1('C_\phi \mathbf{E}_{22}) \leq \frac{4n^2 \rho}{\lambda^2 l} (\frac{a}{m})^m \leq 8 \frac{n^2 \rho a}{\lambda^2} (\frac{a}{m})^m \quad (181)$$

$$\sigma_1(\mathbf{E}_{21}^T (C_\phi + \lambda I)^{-1}) \leq 16 \frac{n}{\lambda} a (\frac{a}{m})^m \leq \frac{n^2 \rho a}{\lambda^2} (\frac{a}{m})^m \quad (182)$$

$$\sigma(\mathbf{E}_{21}^T \mathbf{E}_{22}) \leq 16na (\frac{a}{m})^m \frac{4n}{\lambda^2} (\frac{a}{m})^m = 64 \frac{n^2 a}{\lambda^2} (\frac{a}{m})^{2m} \leq \frac{n^2 \rho a}{\lambda^2} (\frac{a}{m})^m \quad (183)$$

$$(184)$$

Thus, we have shown that $\sigma_1(\mathbf{E}_3) \leq 10 \frac{n^2 \rho a}{\lambda^2} (\frac{a}{m})^m$

For the third term, $(\mathbf{f}(x, \theta) - \mathbf{D}x)^T (\mathbf{A} + \gamma \mathbf{I})^{-1} (\mathbf{f}(x, \theta) - \mathbf{D}x)$ we have

$$|(\mathbf{f}(x, \theta) - \tilde{\mathbf{D}}x)^T (\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-1} (\mathbf{f}(x, \theta) - \tilde{\mathbf{D}}x) - (\mathbf{f}(x, \theta) - \mathbf{D}x)^T (\mathbf{A} + \gamma \mathbf{I})^{-1} (\mathbf{f}(x, \theta) - \mathbf{D}x)| = \quad (185)$$

$$|(\mathbf{f}(x, \theta) - \tilde{\mathbf{D}}x)^T (\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-1} (\mathbf{f}(x, \theta) - \tilde{\mathbf{D}}x) - (\mathbf{f}(x, \theta) - \mathbf{D}x)^T (\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-1} (\mathbf{f}(x, \theta) - \mathbf{D}x) + \quad (186)$$

$$|(\mathbf{f}(x, \theta) - \mathbf{D}x)^T (\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-1} (\mathbf{f}(x, \theta) - \mathbf{D}x) - (\mathbf{f}(x, \theta) - \mathbf{D}x)^T (\mathbf{A} + \gamma \mathbf{I})^{-1} (\mathbf{f}(x, \theta) - \mathbf{D}x)| \leq \quad (187)$$

$$|(\mathbf{f}(x, \theta) - \tilde{\mathbf{D}}x)^T (\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-1} (\mathbf{f}(x, \theta) - \tilde{\mathbf{D}}x) - (\mathbf{f}(x, \theta) - \mathbf{D}x)^T (\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-1} (\mathbf{f}(x, \theta) - \mathbf{D}x)| + \quad (188)$$

$$|(\mathbf{f}(x, \theta) - \mathbf{D}x)^T (\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-1} (\mathbf{f}(x, \theta) - \mathbf{D}x) - (\mathbf{f}(x, \theta) - \mathbf{D}x)^T (\mathbf{A} + \gamma \mathbf{I})^{-1} (\mathbf{f}(x, \theta) - \mathbf{D}x)| \quad (189)$$

$$:= T_1 + T_2 \quad (190)$$

This sum can be bounded by bounding each summand separately.

We start with T_2 :

Lemma 20 Consider the term T_2 as defined above. Then for $m \geq M := 7 + \max\{\frac{e}{2l^2}, \log(\frac{\rho^2 n^3}{\lambda^2 \gamma l^4})\}$, if we use for the approximations quadrature schemes of order $m + 3$ we have:

$$T_2 \leq \epsilon_{32}(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D}\mathbf{x})^T(\mathbf{A} + \gamma\mathbf{I})^{-1}(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D}\mathbf{x}) \quad (191)$$

where

$$\epsilon_{32} = 40 \frac{\rho^2 n^3 a^2}{\lambda^2 \gamma} \left(\frac{a}{m}\right)^m \quad (192)$$

Proof 9 As in the proof of Lemma 18, we can use Woodbury's identity. Now, we use $\mathbf{A} + \gamma\mathbf{I}$ instead of $\mathbf{C}_\phi + \lambda\mathbf{I}$, \mathbf{E}_2 instead of \mathbf{E}_1 and $(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D}\mathbf{x})$ instead of \mathbf{x} . This directly yields

$$\sigma_1(\mathbf{K}^{-\frac{1}{2}})^2 \sigma_1((\mathbf{I} + \mathbf{E}\mathbf{K}^{-1})^{-1}) \sigma_1(\mathbf{E}_2) \leq 40 \frac{\rho^2 n^3 a^2}{\lambda^2 \gamma} \left(\frac{a}{m}\right)^m. \quad (193)$$

Lemma 21 Consider the term T_1 as defined above. Then for $m \geq M := 7 + \max\{\frac{e}{2l^2}, \log(\frac{\rho^2 n^3}{\lambda^2 \gamma l^4})\}$, if we use for the approximations quadrature schemes of order $m + 3$ we have:

$$T_1 \leq \epsilon_{31}(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D}\mathbf{x})^T(\mathbf{A} + \gamma\mathbf{I})^{-1}(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D}\mathbf{x}) + \epsilon_{31}\mathbf{x}^T(\mathbf{C}_\phi + \lambda\mathbf{I})^{-1}\mathbf{x} \quad (194)$$

where $\epsilon_{31} = 30 \frac{n^{\frac{5}{2}} \rho^{\frac{3}{2}} a}{\lambda^2 \gamma^{\frac{1}{2}}} \left(\frac{a}{m}\right)^m$.

Proof 10 It holds that

$$\sigma_1((\tilde{\mathbf{A}} + \gamma\mathbf{I})^{-\frac{1}{2}}) \sigma_1(\mathbf{E}_3) \sigma_1((\mathbf{C}_\phi + \lambda\mathbf{I})^{\frac{1}{2}}) \leq 10 \frac{n^2 \rho a}{\lambda^2} \left(\frac{a}{m}\right)^m \left(\frac{\rho n + \lambda}{\gamma}\right)^{\frac{1}{2}} \leq \quad (195)$$

$$15 \frac{n^{\frac{5}{2}} \rho^{\frac{3}{2}} a}{\lambda^2 \gamma^{\frac{1}{2}}} \left(\frac{a}{m}\right)^m = \frac{\epsilon_{31}}{2} \quad (196)$$

Simply using that $\mathbf{E}_3 = \tilde{\mathbf{D}} - \mathbf{D}$ we get:

$$T_1 = \quad (197)$$

$$|(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \tilde{\mathbf{D}}\mathbf{x})^T(\tilde{\mathbf{A}} + \gamma\mathbf{I})^{-1}(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \tilde{\mathbf{D}}\mathbf{x}) - (\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D}\mathbf{x})^T(\tilde{\mathbf{A}} + \gamma\mathbf{I})^{-1}(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D}\mathbf{x})| = \quad (198)$$

$$|2(\mathbf{E}_3\mathbf{x})^T(\tilde{\mathbf{A}} + \gamma\mathbf{I})^{-1}(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D}\mathbf{x}) + (\mathbf{E}_3\mathbf{x})^T(\tilde{\mathbf{A}} + \gamma\mathbf{I})^{-1}(\mathbf{E}_3\mathbf{x})| \leq \quad (199)$$

$$2|\mathbf{E}_3\mathbf{x}| \sigma_1((\tilde{\mathbf{A}} + \gamma\mathbf{I})^{-\frac{1}{2}}) |(\tilde{\mathbf{A}} + \gamma\mathbf{I})^{-\frac{1}{2}}(\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D}\mathbf{x})| + \sigma_1((\tilde{\mathbf{A}} + \gamma\mathbf{I})^{-1}) |\mathbf{E}_3\mathbf{x}|^2 \quad (200)$$

Since $|\mathbf{E}_3\mathbf{x}| = |\mathbf{E}_3(\mathbf{C}_\phi + \lambda\mathbf{I})^{\frac{1}{2}}(\mathbf{C}_\phi + \lambda\mathbf{I})^{-\frac{1}{2}}\mathbf{x}| \leq \sigma_1(\mathbf{E}_3) \sigma_1((\mathbf{C}_\phi + \lambda\mathbf{I})^{\frac{1}{2}}) |(\mathbf{C}_\phi + \lambda\mathbf{I})^{-\frac{1}{2}}\mathbf{x}|$ we have:

$$\sigma_1((\tilde{\mathbf{A}} + \gamma\mathbf{I})^{-1}) |\mathbf{E}_3\mathbf{x}|^2 \leq \sigma_1((\tilde{\mathbf{A}} + \gamma\mathbf{I})^{-1}) \sigma_1(\mathbf{E}_3)^2 \sigma_1((\mathbf{C}_\phi + \lambda\mathbf{I})^{\frac{1}{2}})^2 |(\mathbf{C}_\phi + \lambda\mathbf{I})^{-\frac{1}{2}}\mathbf{x}|^2 \leq \quad (201)$$

$$\frac{\epsilon_{31}^2}{4} |(\mathbf{C}_\phi + \lambda\mathbf{I})^{-\frac{1}{2}}\mathbf{x}|^2 \leq \frac{\epsilon_{31}}{2} |(\mathbf{C}_\phi + \lambda\mathbf{I})^{-\frac{1}{2}}\mathbf{x}|^2 \quad (202)$$

Note that $|(\mathbf{C}_\phi + \lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{x}|^2$ is the first term of $\mathcal{R}_{\lambda\gamma\phi}(\mathbf{x}, \boldsymbol{\theta})$

Moreover

$$2|\mathbf{E}_3 \mathbf{x}| \sigma_1((\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-\frac{1}{2}}) |(\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-\frac{1}{2}} (\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D} \mathbf{x})| = \quad (203)$$

$$\sigma_1((\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-\frac{1}{2}}) 2|\mathbf{E}_3 (\mathbf{C}_\phi + \lambda \mathbf{I})^{\frac{1}{2}} (\mathbf{C}_\phi + \lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{x}| |(\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-\frac{1}{2}} (\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D} \mathbf{x})| \leq \quad (204)$$

$$\sigma_1((\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-\frac{1}{2}}) \sigma_1(\mathbf{E}_3) \sigma_1((\mathbf{C}_\phi + \lambda \mathbf{I})^{\frac{1}{2}}) 2|(\mathbf{C}_\phi + \lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{x}| |(\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-\frac{1}{2}} (\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D} \mathbf{x})| \leq \quad (205)$$

$$\frac{\epsilon_{31}}{2} (|(\mathbf{C}_\phi + \lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{x}|^2 + |(\tilde{\mathbf{A}} + \gamma \mathbf{I})^{-\frac{1}{2}} (\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D} \mathbf{x})|^2) \leq \quad (206)$$

$$\frac{\epsilon_{31}}{2} (|(\mathbf{C}_\phi + \lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{x}|^2 + (1 + \epsilon_{32}) |(\mathbf{A} + \gamma \mathbf{I})^{-\frac{1}{2}} (\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D} \mathbf{x})|^2) \text{ from (191)} \leq \quad (207)$$

$$\frac{\epsilon_{31}}{2} |(\mathbf{C}_\phi + \lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{x}|^2 + \epsilon_{31} |(\mathbf{A} + \gamma \mathbf{I})^{-\frac{1}{2}} (\mathbf{f}(\mathbf{x}, \boldsymbol{\theta}) - \mathbf{D} \mathbf{x})|^2 \quad (\epsilon_{32} \leq 1 \text{ for } m \geq M) \quad (208)$$

Combining all of these pieces, we can now proof the original Theorem.

Proof 11 (Proof of Theorem 17) Consider $m \geq M := 7 + \max\{\frac{e}{2l^2}, \log(\frac{\rho^2 n^3}{\lambda^2 \gamma l^4})\}$ and let $m' = m + 3$. Then, if we apply quadrature schemes of order m' for $\tilde{\mathcal{R}}$, using the results from Lemmata ((18), (21) and (20)) and that $\epsilon_1 \leq \epsilon_{32}$ and $\epsilon_{31} \leq \epsilon_{32}$ we get

$$\frac{|\mathcal{R}_{\lambda\gamma\phi}(\mathbf{x}, \boldsymbol{\theta}) - \tilde{\mathcal{R}}_{\lambda\gamma\phi}(\mathbf{x}, \boldsymbol{\theta})|}{\mathcal{R}_{\lambda\gamma\phi}(\mathbf{x}, \boldsymbol{\theta})} \leq 2\epsilon_{32} = 80 \frac{\rho^2 n^3 a^2}{\lambda^2 \gamma} \left(\frac{a}{m}\right)^m \leq 50 \frac{\rho^2 n^3}{\lambda^2 \gamma l^4} \left(\frac{a}{m}\right)^m \quad (209)$$

In order to make that smaller than ϵ it suffices $m \geq \max\{\frac{e}{2l^2}, \log(50 \frac{\rho^2 n^3}{\lambda^2 \gamma l^4 \epsilon})\}$. So we can choose for the order of the quadrature scheme m' to be

$$m' = 10 + \max\{\frac{e}{2l^2}, \log(\frac{\rho^2 n^3}{\lambda^2 \gamma l^4 \epsilon})\} \quad (210)$$

Appendix F. Experimental Setups

Here, we will give a brief overview of the experimental setups we used to evaluate ODIN-S.

F.1 Basic Definitions

The trajectory RMSE has proven to be an efficient metric to evaluate the quality of a parameter inference scheme, especially in the context of non-identifiable systems. Here, we restate its definition exactly as provided by Wenk et al. (2020).

Definition 22 (Trajectory RMSE) *Let $\hat{\theta}$ be the parameters estimated by an algorithm. Let \mathbf{t} be the vector collecting the observation times. Define $\tilde{\mathbf{x}}(t)$ as the trajectory one obtains by integrating the ODEs using the estimated parameters, but the true initial value, i.e.*

$$\tilde{\mathbf{x}}(0) = \mathbf{x}^*(0) \quad (211)$$

$$\tilde{\mathbf{x}}(t) = \int_0^t f(\tilde{\mathbf{x}}(s), \hat{\theta}) ds \quad (212)$$

and define $\tilde{\mathbf{x}}$ element-wise as its evaluation at observation times \mathbf{t} , i.e. $\tilde{\mathbf{x}}_i = \tilde{\mathbf{x}}(t_i)$. The trajectory RMSE is then defined as

$$tRMSE := \frac{1}{N} \|\tilde{\mathbf{x}} - \mathbf{x}\|_2 \quad (213)$$

where $\|\cdot\|_2$ denotes the standard Euclidean 2-norm.

Additionally, we would like to restate the definition of the signal-to-noise ratio, as it was used in our work to create the observation noise for the Quadcopter system.

Definition 23 (Signal to Noise Ratio) *Let $x(t)$ be a time-continuous signal for a closed time interval. Let σ_x^2 denote its variance across time. Furthermore, let σ^2 be the variance of an additive, Gaussian noise signal. Then we define the SNR as the ration of these two variances, i.e.*

$$SNR = \frac{\sigma_x^2}{\sigma^2}. \quad (214)$$

F.2 Lotka Volterra

$$\begin{aligned} \dot{x}_1(t) &= \theta_1 x_1(t) - \theta_2 x_1(t) x_2(t) \\ \dot{x}_2(t) &= -\theta_3 x_2(t) + \theta_4 x_1(t) x_2(t). \end{aligned} \quad (215)$$

The Lotka Volterra system (Lotka, 1932) has become a widely used benchmarking system. Due to its locally linear dynamics (Gorbach et al., 2017) and relatively tame trajectories, it is a system many algorithms can solve. We follow the standard setting in the literature and use $\theta = [2, 1, 4, 1]$ and $\mathbf{x}(0) = [5, 3]$ to generate trajectories over the time interval $[0, 2]$. The dynamics are shown in Equation (215).

F.3 Protein Transduction

$$\begin{aligned}
 \dot{S} &= -\theta_1 S - \theta_2 SR + \theta_3 R_S \\
 dS &= \theta_1 S \\
 \dot{R} &= -\theta_2 SR + \theta_3 R_S + \theta_5 \frac{R_{pp}}{\theta_6 + R_{pp}} \\
 \dot{R}_S &= \theta_2 SR - \theta_3 R_S - \theta_4 R_S \\
 \dot{R}_{pp} &= \theta_4 R_S - \theta_5 \frac{R_{pp}}{\theta_6 + R_{pp}}
 \end{aligned} \tag{216}$$

A more challenging system was introduced by Vyshemirsky and Girolami (2007). Its nonlinear terms and non-stationarity introduce interesting challenges for many collocation methods. We follow the standard setting in the literature and use $\theta = [0.07, 0.6, 0.05, 0.3, 0.017, 0.3]$ and $\mathbf{x}(0) = [1, 0, 1, 0, 0]$, but change the time interval to generate trajectories over the time interval $[0, 50]$, since they stay pretty much constant for $t > 50$. The dynamics are shown in Equation (216).

F.4 Lorenz 63

$$\dot{x} = \theta_0(y - x) \tag{217}$$

$$\dot{y} = x(\theta_1 - z) - y \tag{218}$$

$$\dot{z} = xy - \theta_2 z \tag{219}$$

The Lorenz 63 system was introduced by Lorenz (1963) to model atmospheric flows. It is an optimal test bed for parameter inference algorithms, as it exhibits chaotic behavior for the parameter settings we chose. Working with chaotic dynamics is notoriously challenging due to high sensitivity to parameter changes and the presence of many local optima. We follow standard literature and use $\theta = [10, 28, 8/3]$ and $\mathbf{x}(0) = [1, 1, 1]$ to generate trajectories over the time interval $[0, 1]$. The dynamics are shown in Equation (219).

F.5 Quadcopter

$$\begin{aligned}
\dot{x}_0 &= -g \sin(x_7) + x_5 x_1 - x_4 x_2 \\
\dot{x}_1 &= g \sin(x_6) \cos(x_7) - x_0 x_5 + x_2 x_3 \\
\dot{x}_2 &= -\frac{u_0 + u_1 + u_2 + u_3}{\theta_0} + g \cos(x_6) \cos(x_7) + x_0 x_4 - \theta_4 x_1 \\
\dot{x}_3 &= \frac{1}{\theta_1} (\theta_5 (-u_0 + u_1 + u_2 - u_3)) + (\theta_2 - \theta_3 (\theta_2 + \theta_1)) x_4 x_5 \\
\dot{x}_4 &= \frac{1}{\theta_2} (\theta_4 (u_0 - u_1 + u_3 - u_4) + (\theta_3 (\theta_2 + \theta_1) - \theta_1) x_3 x_5) \\
\dot{x}_5 &= \frac{(\theta_1 - \theta_2) x_3 x_4}{\theta_3 (\theta_2 + \theta_1)} \\
\dot{x}_6 &= x_3 + (x_4 \sin(x_6) + \frac{x_5 \cos(x_6) \sin(x_7)}{\cos(x_7)}) \\
\dot{x}_7 &= x_4 \cos(x_6) - x_5 \sin(x_6) \\
\dot{x}_8 &= \frac{x_4 \sin(x_6) + x_5 \cos(x_6)}{\cos(x_7)} \\
\dot{x}_9 &= \cos(x_7) \cos(x_8) x_0 + (-\cos(x_6) \sin(x_8) + \sin(x_6) \sin(x_7) \cos(x_8)) x_1 \\
&\quad + (\sin(x_6) \sin(x_8) + \cos(x_6) \sin(x_7) \cos(x_8)) x_2 \\
\dot{x}_{10} &= \cos(x_7) \sin(x_8) x_0 + (\cos(x_6) \cos(x_8) + \sin(x_6) \sin(x_7) \sin(x_8)) x_1 \\
&\quad + (\cos(x_6) \sin(x_7) \sin(x_8) - \sin(x_6) \cos(x_8)) x_2 \\
\dot{x}_{11} &= \sin(x_7) x_0 - \sin(x_6) \cos(x_7) x_1 - \cos(x_6) \cos(x_7) x_2
\end{aligned} \tag{220}$$

As an ultimate benchmark, we introduce a parametric model describing the dynamics of a 6DOF quadcopter, shown in Equation (220). Its strongly nonlinear dynamics and the presence of inputs make it a formidable challenge. The states of this system are representing the linear velocities (x_0, x_1, x_2) , the angular velocities (x_3, x_4, x_5) , the angles (x_6, x_7, x_8) and the position (x_9, x_{10}, x_{11}) of the quadcopter. The four inputs represent the forces applied at the four different propellers. While in principle any input commands could be incorporated, we keep the inputs constant at $u = [0.248, 0.2475, 0.24775, 0.24775]$. This input leads to interesting nonstationary climbing, pitching and rolling behavior. We use $\theta = [0.1, 0.00062, 0.00113, 0.9, 0.114, 0.0825, 9.85]$ and $x_i(0) = 0$ for $i = 0 \dots 11$ to generate trajectories over the time interval $[0, 15]$.

Appendix G. Additional Empirical Evaluation SLEIPNIR

G.1 tRMSE vs Features

G.1.1 LOTKA VOLTERRA

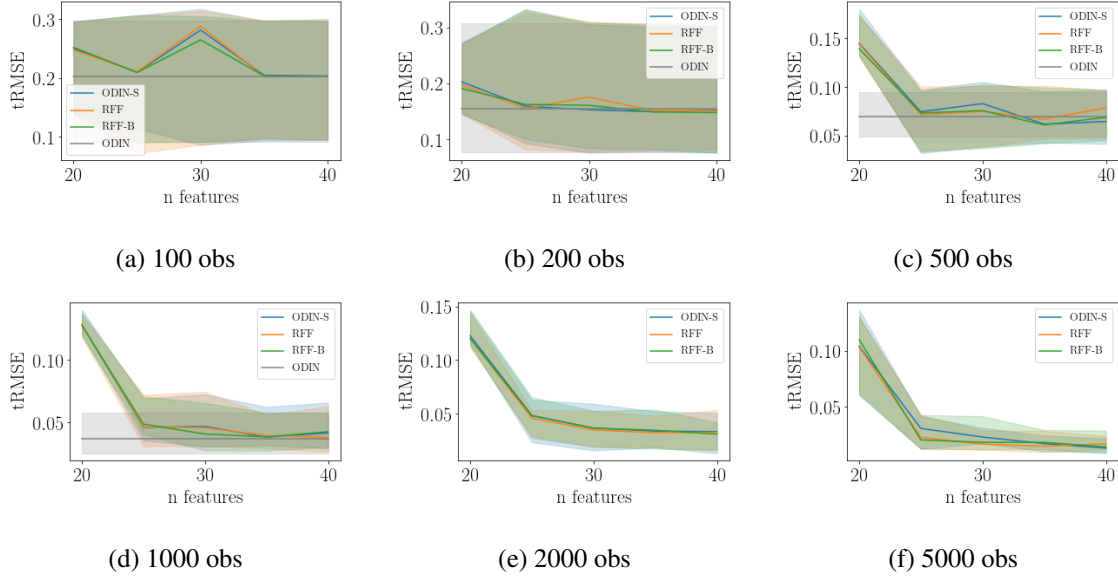


Figure 15: tRMSE vs features for the Lotka Volterra system using additive observation noise with $\sigma^2 = 0.1$.

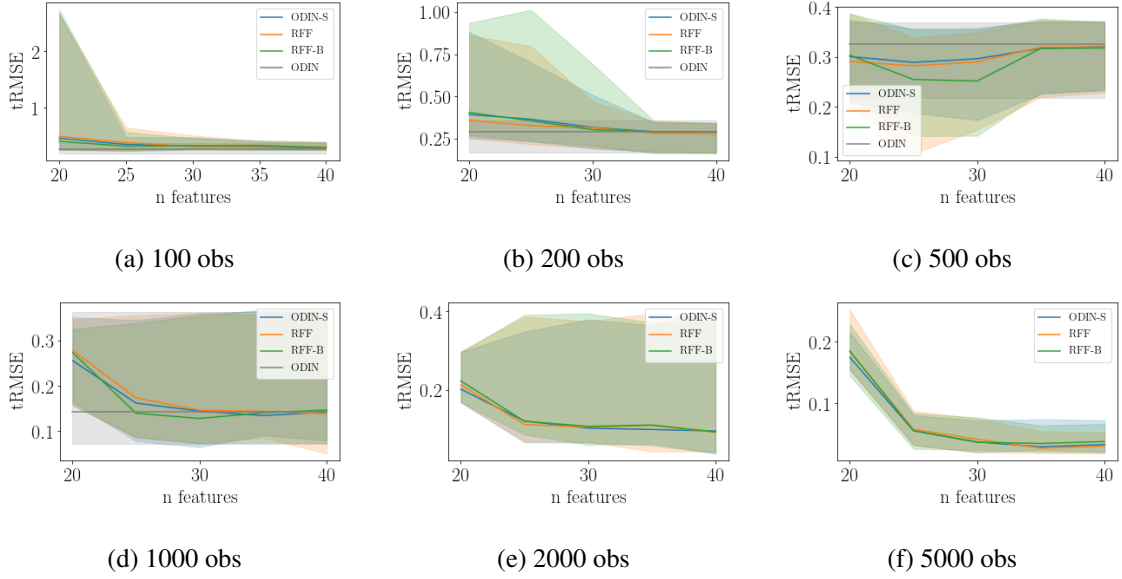


Figure 16: tRMSE vs features for the Lotka Volterra system using additive observation noise with $\sigma^2 = 0.5$.

G.1.2 PROTEIN TRANSDUCTION

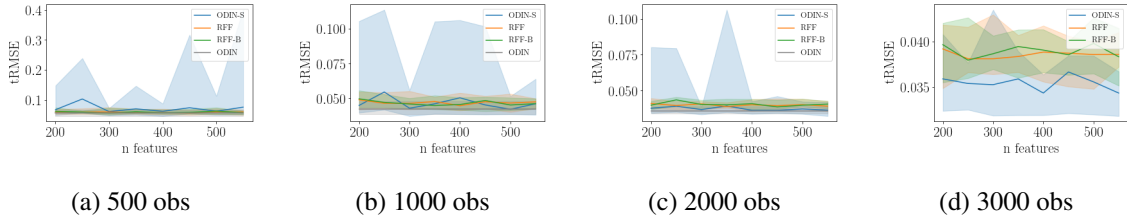


Figure 17: tRMSE vs features for the Protein Transduction system using additive observation noise with $\sigma^2 = 0.01$.

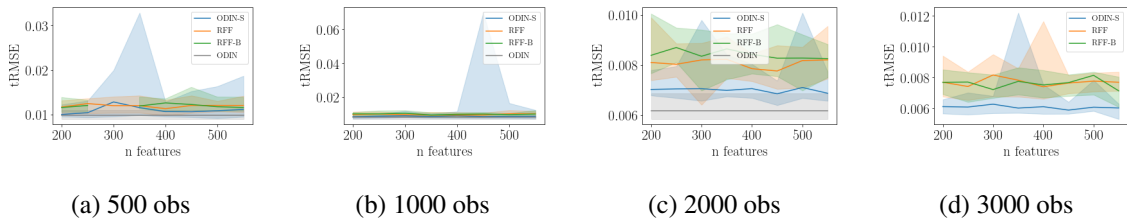


Figure 18: tRMSE vs features for the Protein Transduction system using additive observation noise with $\sigma^2 = 0.0001$.

G.1.3 LORENZ

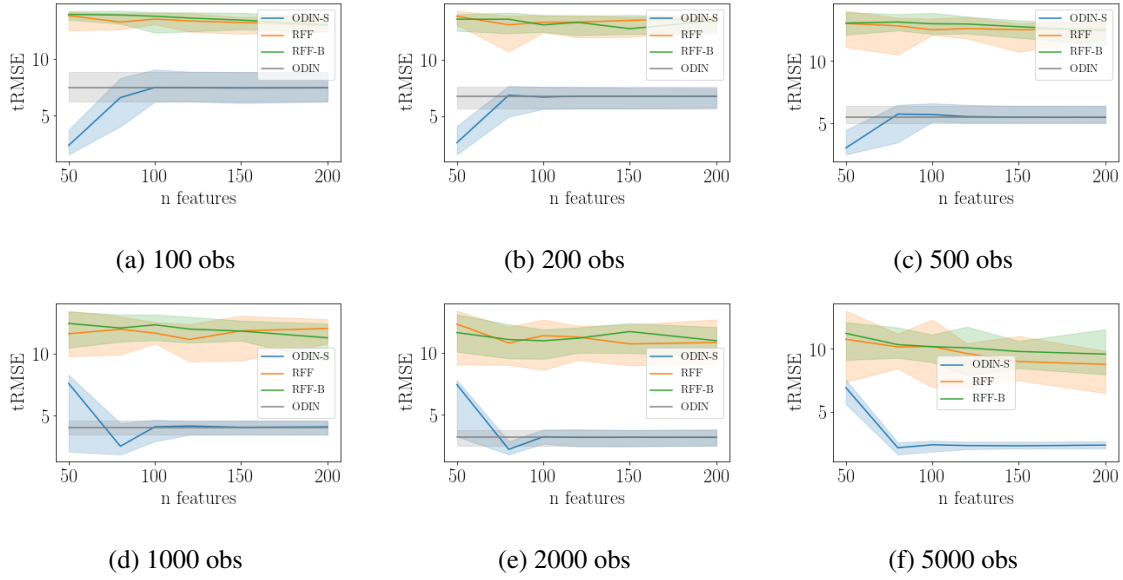


Figure 19: tRMSE vs features for the Lorenz system with noise created using a signal-to-noise ratio of 5.

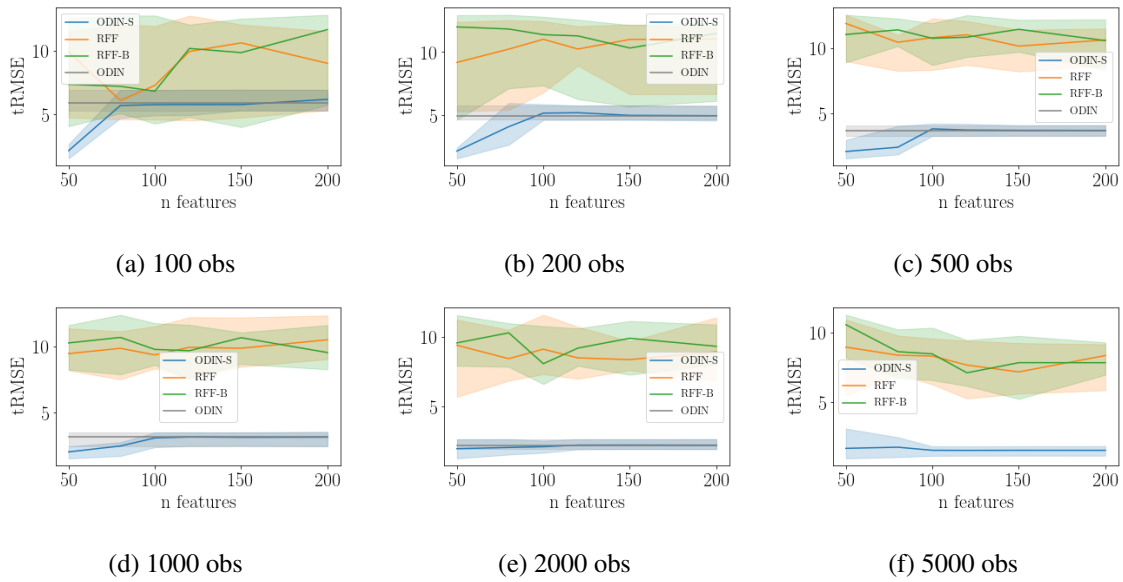


Figure 20: tRMSE vs features for the Lorenz system with noise created using a signal-to-noise ratio of 10.

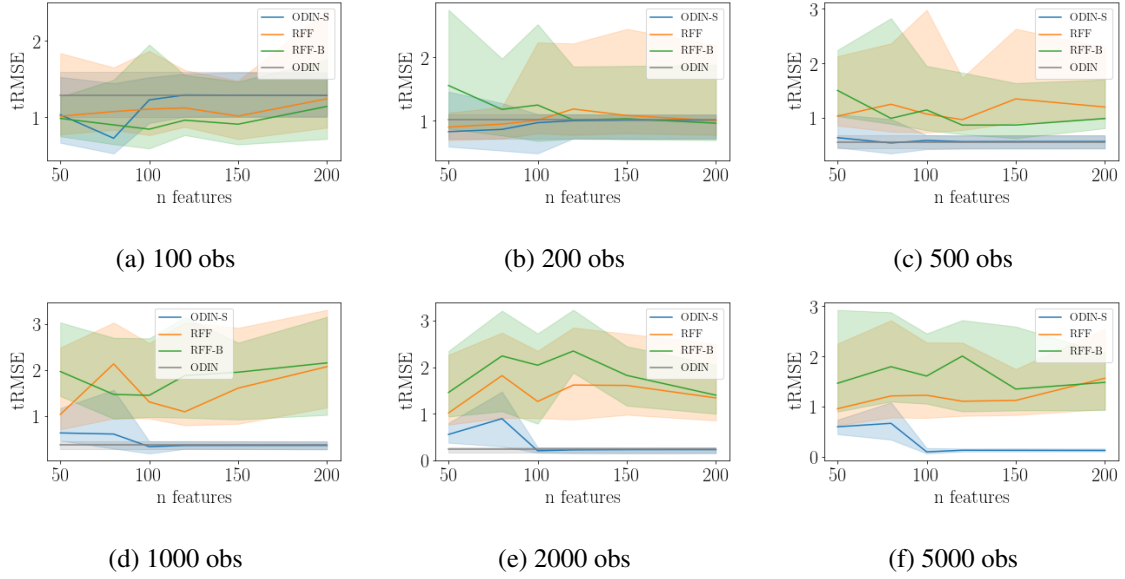


Figure 21: tRMSE vs features for the Lorenz system with noise created using a signal-to-noise ratio of 100.

G.2 Learning Curves

G.2.1 LOTKA VOLTERRA

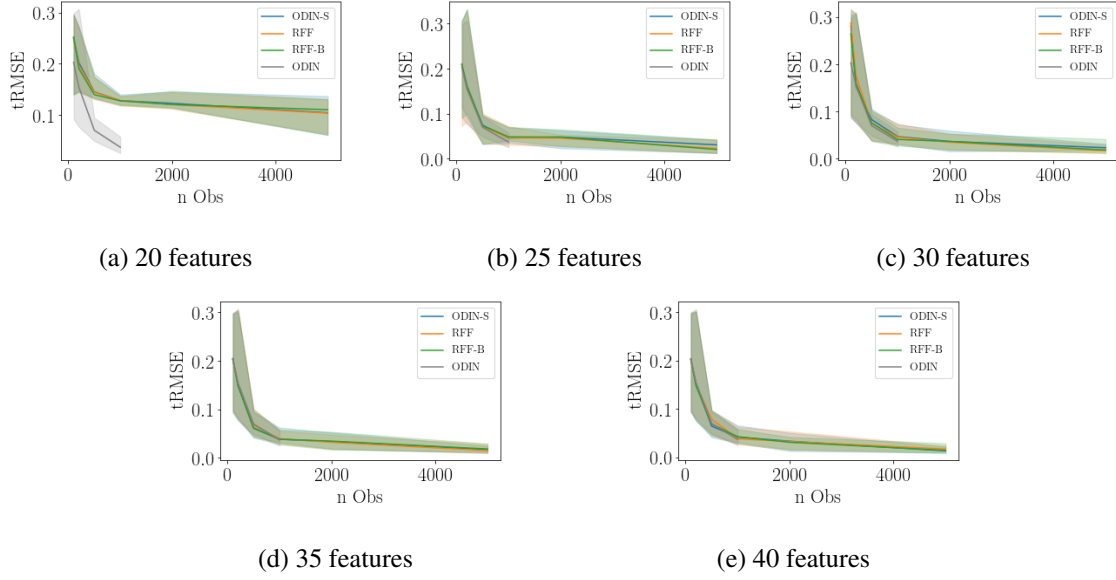


Figure 22: tRMSE vs amount of observations for the Lotka Volterra system with additive noise with $\sigma^2 = 0.1$.

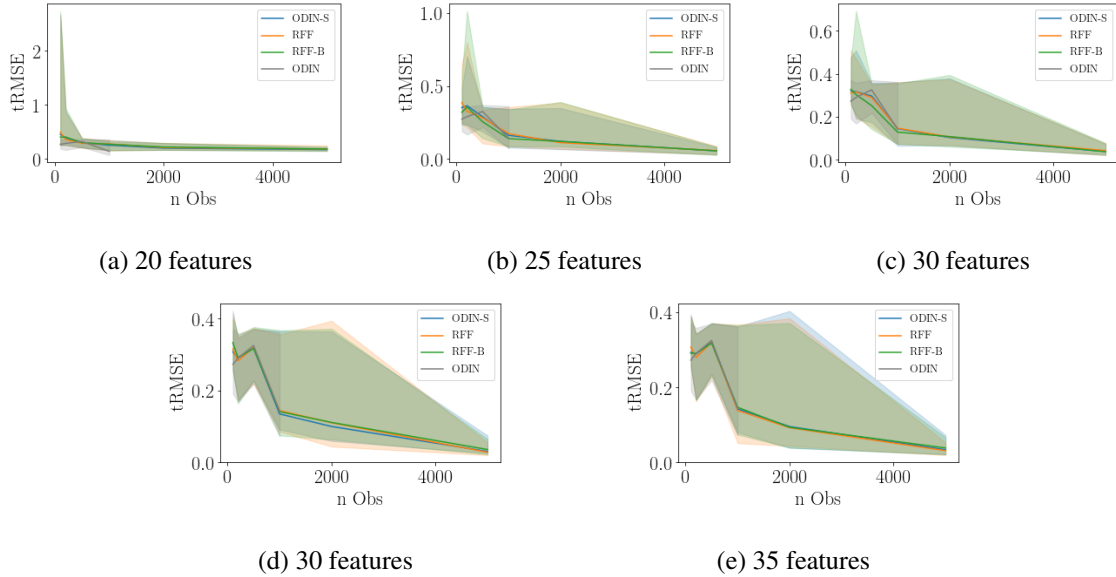


Figure 23: tRMSE vs amount of observations for the Lotka Volterra system with additive noise with $\sigma^2 = 0.5$.

G.2.2 PROTEIN TRANSDUCTION

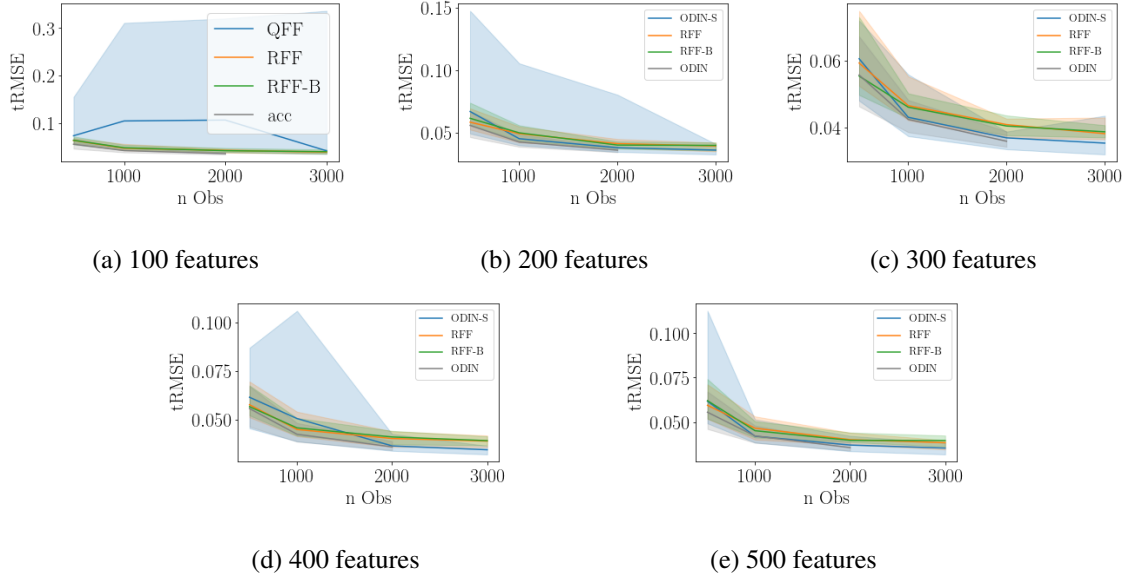


Figure 24: tRMSE vs amount of observations for the Protein Transduction system using additive Gaussian noise with $\sigma^2 = 0.01$.

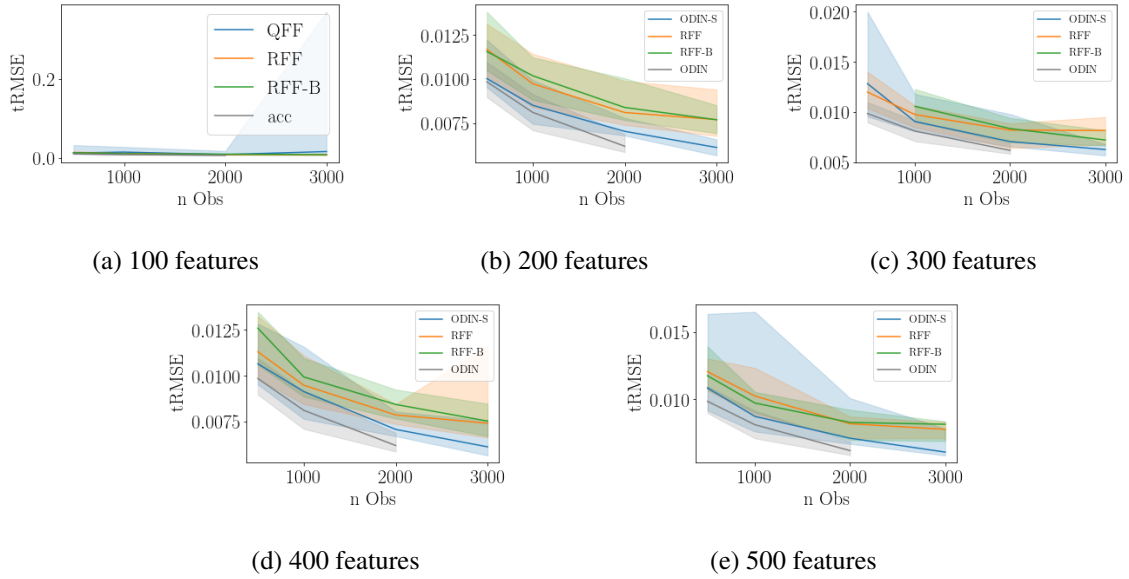


Figure 25: tRMSE vs amount of observations for the Protein Transduction system using additive Gaussian noise with $\sigma^2 = 0.0001$.

G.2.3 LORENZ

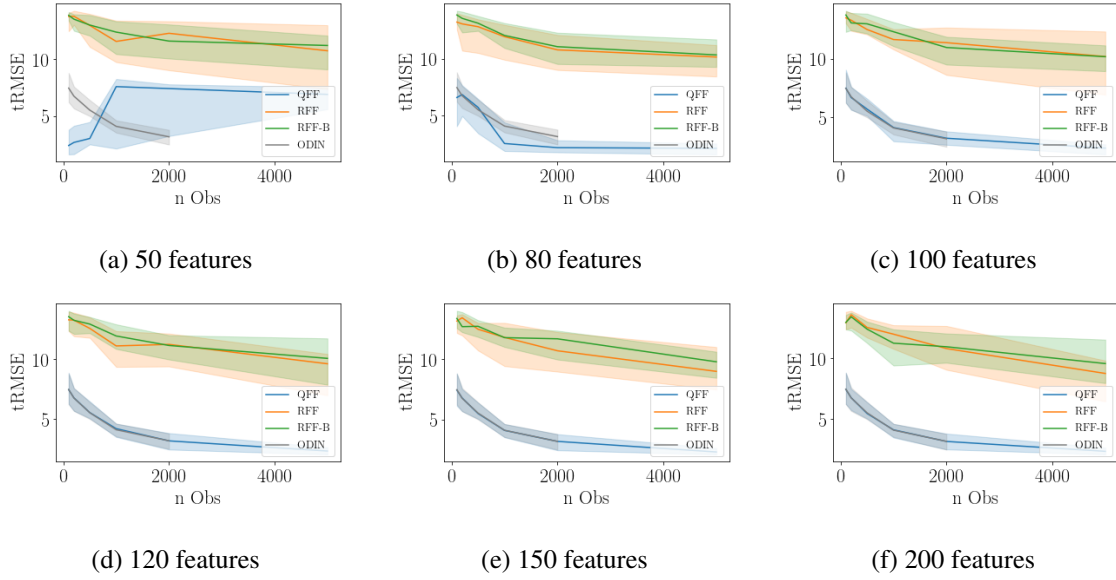


Figure 26: tRMSE vs amount of observations for the Lorenz system with noise created using a signal-to-noise ratio of 5.

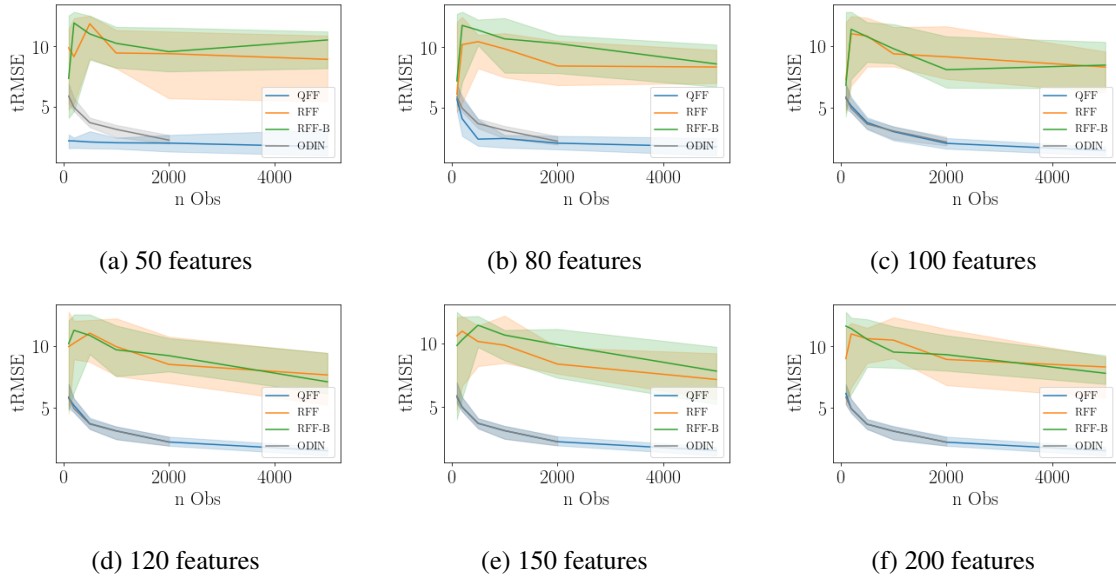


Figure 27: tRMSE vs amount of observations for the Lorenz system with noise created using a signal-to-noise ratio of 10.

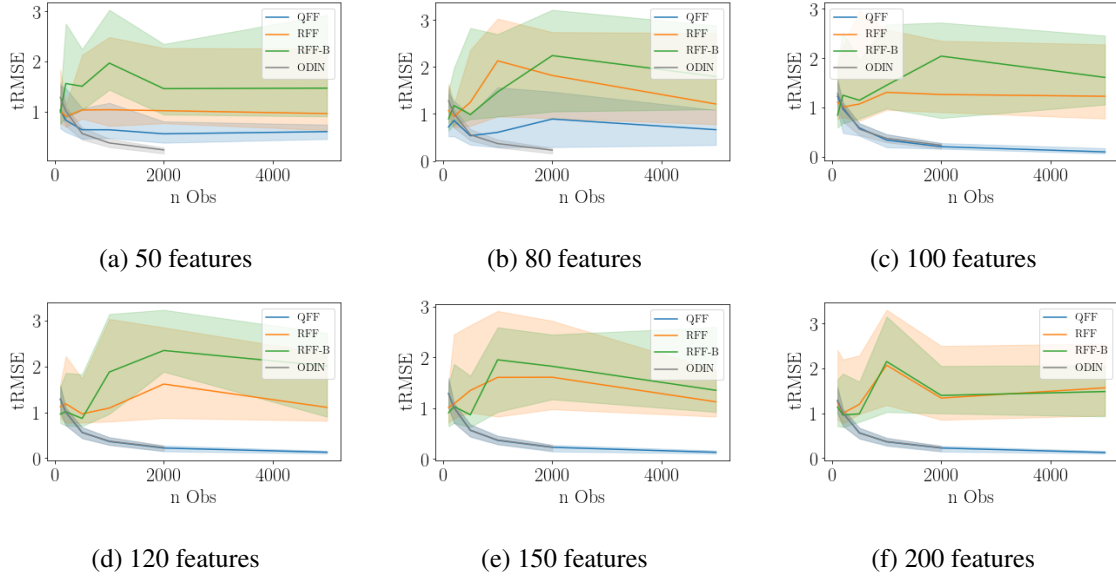


Figure 28: tRMSE vs amount of observations for the Lorenz system with noise created using a signal-to-noise ratio of 100.

G.2.4 RUN TIME

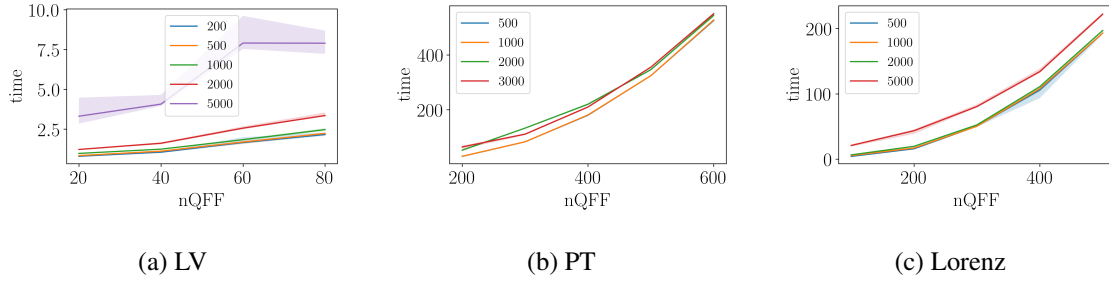


Figure 29: Run time per iteration in ms vs amount of features for different amounts of observations. As expected from theoretical analysis, the run time per iteration scales approximately cubic.

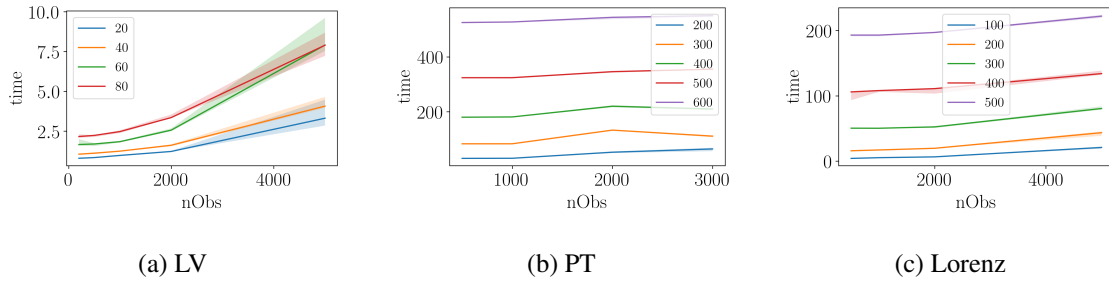


Figure 30: Run time per iteration in ms vs amount of observations for different amounts of Fourier features. As expected from theoretical analysis, the run time per iteration scales approximately linear, even though there is a strong bias term.

G.3 Quadcopter State Inference

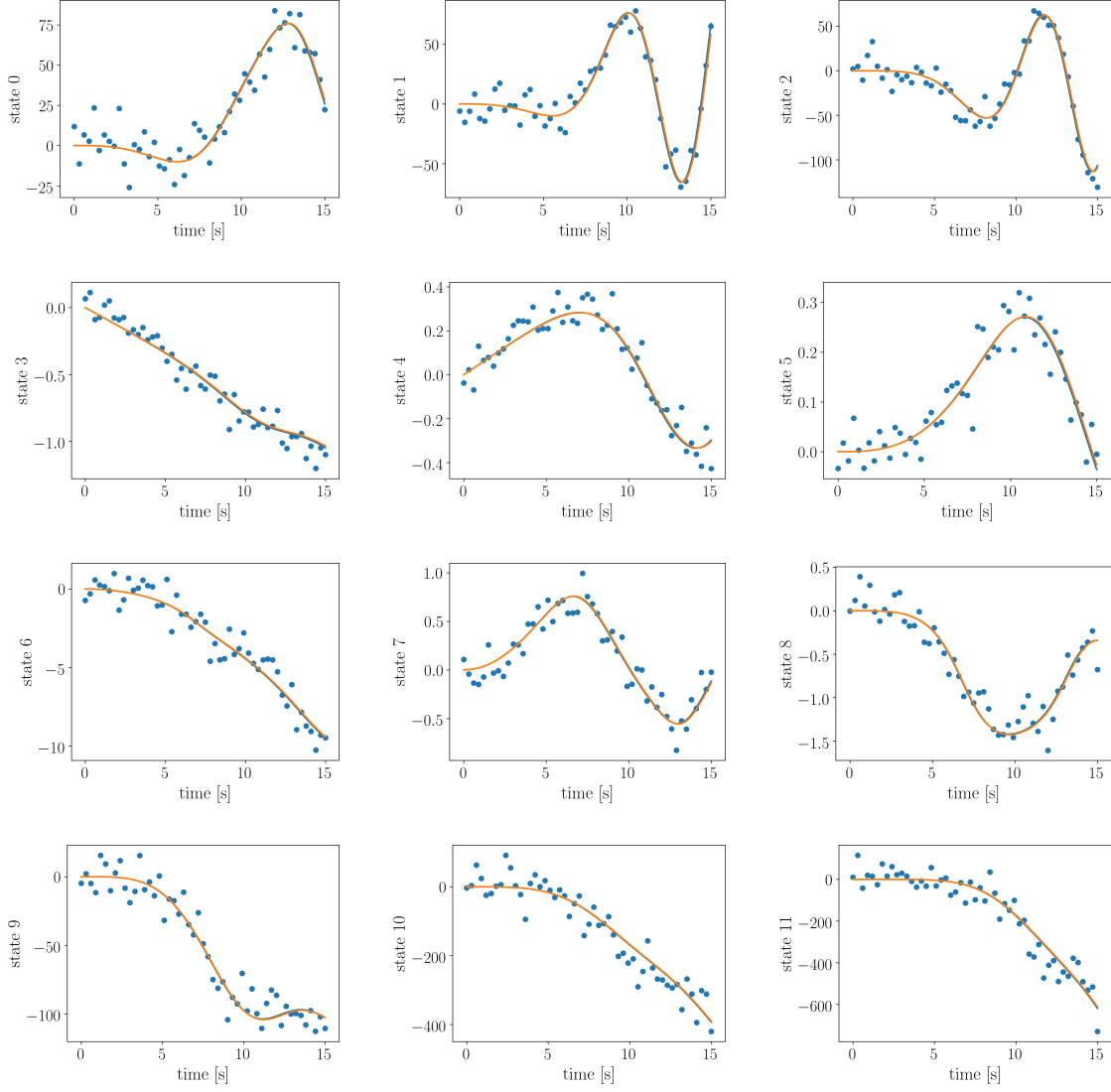


Figure 31: State trajectories obtained by integrating the parameters inferred by ODIN-S (orange). The blue line represents the ground truth, while the blue dots show every 300-th observation for a signal-to-noise ratio of 10.