

Learning unknown forces in nonlinear models with Gaussian processes and autoregressive flows

Wil O C Ward

`w.ward@sheffield.ac.uk`

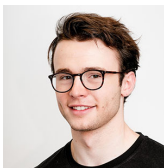
Department of Physics and Astronomy, The University of Sheffield

*GPSS Workshop: Structurally Constrained Gaussian
Processes
12 Sep 2019*

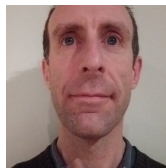
Collaborative Work



Mauricio Alvarez



Tom Ryder

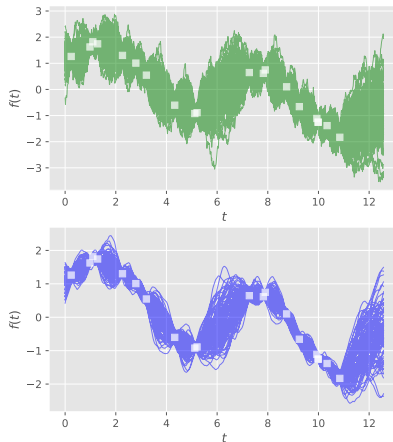


Dennis Prangle

Gaussian Processes

- GPs generalise Gaussian distribution
- Infinite dimension and non-parametric
- Defined in terms of mean and covariance function

$$f(t) \sim \text{GP}(m(t), k(t, t'))$$



Motivating Example

Consider the model,

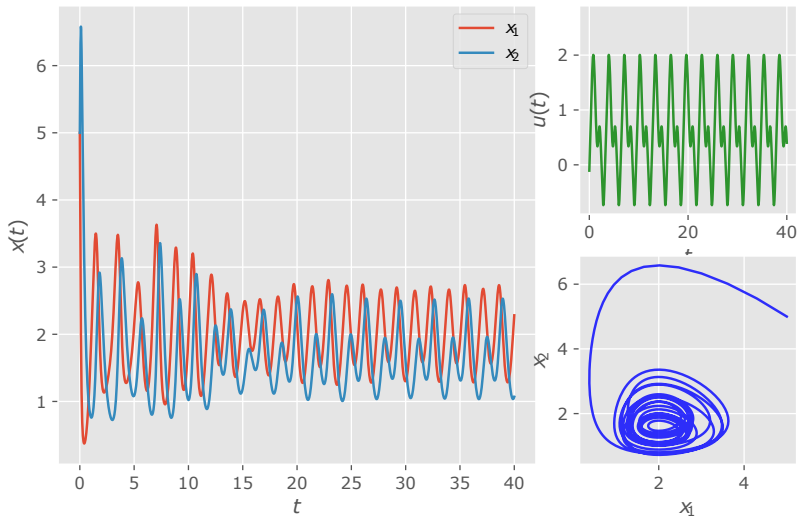
$$\frac{d}{dt}\mathbf{x} = \boldsymbol{\alpha}(\mathbf{x}(t), \boldsymbol{\theta}) + \begin{bmatrix} u(t) & 0 \end{bmatrix}^T$$

Where $\boldsymbol{\alpha} : \mathbb{R}^2 \times \Theta \rightarrow \mathbb{R}^2$ are known dynamics:

$$\boldsymbol{\alpha}(\mathbf{x}, \boldsymbol{\theta}) = \begin{bmatrix} \theta_1 x_1 - \theta_2 x_1 x_2 \\ \theta_2 x_1 x_2 - \theta_3 x_2 \end{bmatrix}$$

...but $\boldsymbol{\theta}$ and $u(t)$ are unknown. How can we infer $\mathbf{x}(t)$ and $u(t)$ given some noisy observations $\mathbf{y} = [\mathbf{x}(\tau_j) + \varepsilon_j]_{j=0}^N$?

Motivating Example



Contents

- 1 Stochastic Differential Equations and Gaussian Processes
- 2 Variational Solutions to Non-Linear Latent Force Models
- 3 Approximate Gaussian Processes
- 4 Some Results
- 5 Recap
- 6 Open Issues

Contents

- 1 Stochastic Differential Equations and Gaussian Processes
- 2 Variational Solutions to Non-Linear Latent Force Models
- 3 Approximate Gaussian Processes
- 4 Some Results
- 5 Recap
- 6 Open Issues

- Consider an ordinary differential equation describing the dynamics of some (vector-valued) function $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^d$
- The dynamics $\alpha_k : \mathbb{R}^d \rightarrow \mathbb{R}^d$ are known but it is driven by a white-noise process with covariance as function of \mathbf{x} ,
 $\Sigma : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$

Ordinary Differential Equation with White Noise

$$\sum_{k=0}^n \alpha_k(\mathbf{x}, t; \boldsymbol{\theta}) \frac{d^n}{dt^n} \mathbf{x}(t) = \Sigma^{1/2}(\mathbf{x}, t; \boldsymbol{\theta}) \mathbf{w}(t)$$

- Consider an ordinary differential equation describing the dynamics of some (vector-valued) function $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^d$
- The dynamics $\boldsymbol{\alpha}_k : \mathbb{R}^d \rightarrow \mathbb{R}^d$ are known but it is driven by a white-noise process with covariance as function of \mathbf{x} ,
 $\boldsymbol{\Sigma} : \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$

Stochastic Differential Equation

$$\sum_{k=0}^n \underbrace{\alpha_k(\mathbf{x}, t; \boldsymbol{\theta})}_{\text{drift terms}} \frac{d^n}{dt^n} \mathbf{x}(t) = \underbrace{\boldsymbol{\Sigma}^{1/2}(\mathbf{x}, t; \boldsymbol{\theta})}_{\text{diffusion}} \mathbf{w}(t)$$

Solutions to Itô Processes

- If system has linear dynamics, can solve exactly using Kalman filtering / Rauch-Tung-Streibel smoothing
- Assuming non-linearity, there are a number of approximation methods
- Stochastic extension to Euler method for iterative discrete-time estimation

Euler-Maruyama Discretisation

$$\mathbf{x}(t_{k+1}) - \mathbf{x}(t_k) \sim \mathcal{N}(\boldsymbol{\alpha}(\mathbf{x}(t_k))\Delta_t, \boldsymbol{\Sigma}\Delta_t)$$

Solutions to Itô Processes

- If system has linear dynamics, can solve exactly using Kalman filtering / Rauch-Tung-Streibel smoothing
- Assuming non-linearity, there are a number of approximation methods
- Stochastic extension to Euler method for iterative discrete-time estimation

Euler-Maruyama Discretisation as a Generative Prior

$$\mathbf{x}(t_{k+1}) \mid \mathbf{x}(t_k) \sim \mathcal{N}(\mathbf{x}(t_k) + \boldsymbol{\alpha}(\mathbf{x}(t_k))\Delta_t, \boldsymbol{\Sigma}\Delta_t)$$

Examples

- White noise process

$$w(t) \sim \text{GP}(0, \varsigma^2 \delta(t - t'))$$

- Half-integer ($\nu = p + 1/2$) Matérn models

$$f_\nu(t) \sim \text{GP}\left(0, \sigma^2 \exp(-\lambda|t - t'|) \frac{p!}{(2p)!} \sum_{i=0}^p \frac{(p+i)!}{i!(p-i)!} (2\lambda|t - t'|)^{p-i}\right)$$

- Gaussian Radial Basis / Exponentiated Quadratic ($\nu \rightarrow \infty$)

$$f(t) \sim \text{GP}(0, \sigma^2 \exp(-\lambda|t - t'|^2))$$

Examples

- White noise process

$$dw(t) = \varsigma d\beta$$

- Half-integer ($\nu = p + 1/2$) Matérn models

$$\sum_{i=1}^p \binom{p}{i-1} \lambda^{p+1-i} \frac{d^i}{dt^i} f(t) = -\lambda^{p+1} f(t) + w(t)$$

- Gaussian Radial Basis / Exponentiated Quadratic ($\nu \rightarrow \infty$)
infinitely differentiable so cannot represent as Itô process exactly

Examples

- White noise process

$$dw(t) = \varsigma d\beta$$

- Half-integer ($\nu = p + 1/2$) Matérn models

$$d\mathbf{f}(t) = \underbrace{\begin{bmatrix} 0 & 1 & & \\ \ddots & \ddots & \ddots & \\ -a_1\lambda^{p+1} & -a_2\lambda^p & \cdots & -a_p\lambda \end{bmatrix}}_{\mathbf{G}} \underbrace{\begin{bmatrix} f(t) \\ df/dt \\ \vdots \\ d^{p-1}f/dt^{p-1} \end{bmatrix}}_{\mathbf{f}(t)} dt + \varsigma \underbrace{\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}}_{\mathbf{w}(t)} d\beta$$

Stochastic Latent Force Models

- Recall our motivating example, a mixture of known dynamics with some hidden input function
- General form:

$$\alpha_0(\mathbf{x}, t; \boldsymbol{\theta})\mathbf{x}(t) + \alpha_1(\mathbf{x}, t; \boldsymbol{\theta})\frac{d}{dt}\mathbf{x}(t) + \dots = \mathbf{u}(t)$$

- Placing a GP prior over $\mathbf{u}(t)$
- Termed *latent force models*

M. A. Alvarez, D. Luengo, and N. D. Lawrence. Linear latent force models using Gaussian processes. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(11):2693–2705, 2013

Companion Form LFM's

- Easy enough to reframe n^{th} -order differential equation as first-order $d\mathbf{f}/dt = \mathbf{D}(\mathbf{f}(t), \boldsymbol{\theta}) + \mathbf{L}w(t)$

Companion Form LFM's

- Easy enough to reframe n^{th} -order differential equation as first-order $d\mathbf{f}/dt = \mathbf{D}(\mathbf{f}(t), \boldsymbol{\theta}) + \mathbf{L}w(t)$

Companion Form

$$\mathbf{f}(\tau) = \left[x(\tau) \quad \left. \frac{dx}{dt} \right|_{t=\tau} \quad \cdots \quad \left. \frac{d^{n-1}x}{dt^{n-1}} \right|_{t=\tau} \quad u(\tau) \quad \left. \frac{du}{dt} \right|_{t=\tau} \quad \cdots \quad \left. \frac{d^{m-1}u}{dt^{m-1}} \right|_{t=\tau} \right]^\top$$

$$D(\mathbf{f}(t), \boldsymbol{\theta}) = \begin{bmatrix} f_2 \\ f_3 \\ \vdots \\ \ddot{\alpha}_0 f_1 + \sum_{i=1}^{n-1} \ddot{\alpha}_i f_{i+1} + f_{n+1} \\ f_{n+2} \\ f_{n+3} \\ \vdots \\ a_0 f_{n+1} + \sum_{i=1}^{m-1} a_i f_{n+i+1} \end{bmatrix}, \quad \mathbf{L} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

Contents

- 1 Stochastic Differential Equations and Gaussian Processes
- 2 Variational Solutions to Non-Linear Latent Force Models
- 3 Approximate Gaussian Processes
- 4 Some Results
- 5 Recap
- 6 Open Issues

Inferring the Joint Posterior of a Non-Linear LFM

Problem: Infer \mathbf{f} and $\boldsymbol{\theta}$

$$\frac{d}{dt}\mathbf{f}(t) = \mathbf{D}(\mathbf{f}(t), \boldsymbol{\theta}) + \mathbf{L}w(t)$$

- We cannot infer \mathbf{f} exactly if \mathbf{D} is non-linear since the joint posterior is intractible
- Pseudo-chaos under some systems
- Non-linear versions of filters/smoothers, e.g. E/UKF, ADF, SMC
- Difficult to do joint parameter estimation, difficult to use autodifferentiation

J. Hartikainen, M. Seppänen, and S. Särkkä. State-space inference for non-linear latent force models with application to satellite orbit prediction. In *ICML*, pages 723–730, 2012.

Variational Bridge Constructs

We want to build variational approximation of conditional posterior: $p(\mathbf{x}, \mathbf{u}, \boldsymbol{\theta} \mid \mathbf{y})$.

Variational Bayes

Find $q^* \in \mathcal{Q}$, such that

$$q^* = \arg \min_{q \in \mathcal{Q}} \text{KL}[q(\mathbf{x}, \mathbf{u}, \boldsymbol{\theta}) \parallel p(\mathbf{x}, \mathbf{u}, \boldsymbol{\theta} \mid \mathbf{y})]$$

where \mathcal{Q} is a family of distributions parameterised by ϕ

Variational Bridge Constructs

We want to build variational approximation of conditional posterior: $p(\mathbf{f}, \boldsymbol{\theta} \mid \mathbf{y})$.

Variational Bayes

Find $q^* \in \mathcal{Q}$, such that

$$q^* = \arg \min_{q \in \mathcal{Q}} \text{KL}[q(\mathbf{f}, \boldsymbol{\theta}) \parallel p(\mathbf{f}, \boldsymbol{\theta} \mid \mathbf{y})]$$

where \mathcal{Q} is a family of distributions parameterised by ϕ

Variational Bridge Constructs

Evidence Lower Bound (ELBO)

$$\mathcal{L}(\phi) = \mathbb{E}_{\mathbf{f}, \boldsymbol{\theta} \sim q} [\log p(\mathbf{f}, \boldsymbol{\theta}, \mathbf{y}) - \log q(\mathbf{f}, \boldsymbol{\theta})]$$

Variational Bridge Constructs

Unbiased Evidence Lower Bound (ELBO)

$$\hat{\mathcal{L}}(\phi) = \frac{1}{n_s} \sum_{i=1}^{n_s} \log \frac{p(\boldsymbol{\theta}^{(i)}) p(\mathbf{f}^{(i)} | \boldsymbol{\theta}^{(i)}) p(\mathbf{y} | \mathbf{f}^{(i)}, \boldsymbol{\theta}^{(i)})}{q(\boldsymbol{\theta}^{(i)}) q(\mathbf{f}^{(i)} | \boldsymbol{\theta}^{(i)})}$$

where $\mathbf{f}^{(i)} \sim q(\mathbf{f} | \boldsymbol{\theta}^{(i)})$ and $\boldsymbol{\theta}^{(i)} \sim q(\boldsymbol{\theta})$ $i = 1, \dots, n_s$

Variational Bridge Constructs

Unbiased Evidence Lower Bound (ELBO)

$$\hat{\mathcal{L}}(\phi) = \frac{1}{n_s} \sum_{i=1}^{n_s} \log \frac{p(\boldsymbol{\theta}^{(i)}) p(\mathbf{f}^{(i)} | \boldsymbol{\theta}^{(i)}) p(\mathbf{y} | \mathbf{f}^{(i)}, \boldsymbol{\theta}^{(i)})}{q(\boldsymbol{\theta}^{(i)}) q(\mathbf{f}^{(i)} | \boldsymbol{\theta}^{(i)})}$$

where $\mathbf{f}^{(i)} \sim q(\mathbf{f} | \boldsymbol{\theta}^{(i)})$ and $\boldsymbol{\theta}^{(i)} \sim q(\boldsymbol{\theta})$ $i = 1, \dots, n_s$

Likelihood Agnostic

Valid for any (differentiable?) observation model $p(\mathbf{y} | \mathbf{f}, \boldsymbol{\theta})$

Black-box Variational Inference

- Black-box variational inference (BBVI) is predicated on the fact that the gradient of ELBO can be written as an unbiased average
- Straightforward since we have $\hat{\mathcal{L}}(\phi)$ as an unbiased average

Black-box Variational Inference

- Black-box variational inference (BBVI) is predicated on the fact that the gradient of ELBO can be written as an unbiased average
- Straightforward since we have $\hat{\mathcal{L}}(\phi)$ as an unbiased average

Monte Carlo approximation of ELBO gradient

$$\nabla_{\phi} \mathcal{L}(\phi) \approx \frac{1}{n_s} \sum_{i=1}^{n_s} \nabla_{\phi} \log q(\mathbf{f}^{(i)}, \boldsymbol{\theta}^{(i)}) \log \frac{p(\mathbf{f}^{(i)}, \boldsymbol{\theta}^{(i)}, \mathbf{y})}{q(\mathbf{f}^{(i)}, \boldsymbol{\theta}^{(i)})}$$

where $\mathbf{f}^{(i)} \sim q(\mathbf{f} | \boldsymbol{\theta}^{(i)})$ and $\boldsymbol{\theta}^{(i)} \sim q(\boldsymbol{\theta})$ $i = 1, \dots, n_s$

R. Ranganath, S. Gerrish, and D. Blei. Black box variational inference. In *Artificial Intelligence and Statistics*, 2014.

D. Duvenaud and R. P. Adams. Black-box stochastic variational inference in five lines of Python. In *NIPS Workshop on Black-box Learning and Inference*, 2015.

Algorithm 1 BBVI with gradient ascent

Initialise ϕ_0 (randomly)

$j \leftarrow 0$

while *not converged* **do**

 Calculate $\nabla_{\phi} \mathcal{L}(\phi_j)$

 Update ϕ w.r.t. ELBO gradient, e.g.:

$$\phi_{j+1} \leftarrow \phi_j + h \nabla_{\phi} \mathcal{L}(\phi_j)$$

$j \leftarrow j + 1$

end while

Variational approximation $q(\mathbf{f}, \boldsymbol{\theta} \mid \phi_j) \approx p(\mathbf{f}, \boldsymbol{\theta} \mid \mathbf{y})$

Parameter Estimation: $q(\boldsymbol{\theta})$

Commonly in system estimation, the model parameters, $\boldsymbol{\theta}$ are unknown.

We can also give these a Bayesian treatment by using variational representation of the posterior.

We can use any variational approach here, e.g. mean-field:

$$q(\boldsymbol{\theta}) = \prod \mathcal{N}(\theta_i | m_i, s_i)$$

Here, the free parameters are scalars, $\phi_{\theta} = \{(m_i, s_i)\}_{\forall i}$

Filtering Density: $p(\mathbf{f} \mid \boldsymbol{\theta})$

Represent stochastic process, \mathbf{f} as a filtering distribution with moments $\mathbf{m}(t)$ and $\mathbf{P}(t)$

$$\text{Mean term:} \quad \frac{d}{dt} \mathbf{m}(t) = \mathbf{D}(\mathbf{m}, t; \boldsymbol{\theta})$$

Filtering Density: $p(\mathbf{f} \mid \boldsymbol{\theta})$

Represent stochastic process, \mathbf{f} as a filtering distribution with moments $\mathbf{m}(t)$ and $\mathbf{P}(t)$

Mean term: $\frac{d}{dt}\mathbf{m}(t) = \mathbf{D}(\mathbf{m}, t; \boldsymbol{\theta})$

Covariance term: $??$

Extended Filtering Density: $p(\mathbf{f} \mid \boldsymbol{\theta})$

Represent stochastic process, \mathbf{f} as a filtering distribution with moments $\mathbf{m}(t)$ and $\mathbf{P}(t)$

Mean term:
$$\frac{d}{dt}\mathbf{m}(t) = \mathbf{D}(\mathbf{m}, t; \boldsymbol{\theta})$$

Covariance term:

$$\frac{d}{dt}\mathbf{P}(t) = \mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})\mathbf{P}(t) + \mathbf{P}(t)\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})^T + \mathbf{L}\Sigma^2\mathbf{L}^T$$

Extended Filtering Density: $p(\mathbf{f} | \boldsymbol{\theta})$

Represent stochastic process, \mathbf{f} as a filtering distribution with moments $\mathbf{m}(t)$ and $\mathbf{P}(t)$

Assume steady state: $d\mathbf{P}/dt = 0$.

Extended Filtering Density: $p(\mathbf{f} | \boldsymbol{\theta})$

Represent stochastic process, \mathbf{f} as a filtering distribution with moments $\mathbf{m}(t)$ and $\mathbf{P}(t)$

Assume steady state: $d\mathbf{P}/dt = 0$.

Denote covariance in steady state by $\tilde{\Sigma}$ and solve

$$\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})\tilde{\Sigma} + \tilde{\Sigma}\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})^T = -\mathbf{L}\zeta^2\mathbf{L}^T$$

Extended Filtering Density: $p(\mathbf{f} | \boldsymbol{\theta})$

Represent stochastic process, \mathbf{f} as a filtering distribution with moments $\mathbf{m}(t)$ and $\mathbf{P}(t)$

Assume steady state: $d\mathbf{P}/dt = 0$.

Denote covariance in steady state by $\tilde{\Sigma}$ and solve

$$\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})\tilde{\Sigma} + \tilde{\Sigma}\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})^T = -\mathbf{L}\zeta^2\mathbf{L}^T$$

Example of continuous Lyapunov equation; easy to solve numerically, but need an differentiable form of $\tilde{\Sigma}$

Extended Filtering Density: $p(\mathbf{f} | \boldsymbol{\theta})$

Represent stochastic process, \mathbf{f} as a filtering distribution with moments $\mathbf{m}(t)$ and $\mathbf{P}(t)$

Assume steady state: $d\mathbf{P}/dt = 0$.

Denote covariance in steady state by $\tilde{\Sigma}$ and solve

$$\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})\tilde{\Sigma} + \tilde{\Sigma}\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})^T = -\mathbf{L}\zeta^2\mathbf{L}^T$$

Example of continuous Lyapunov equation; easy to solve numerically, but need an differentiable form of $\tilde{\Sigma}$

$\tilde{\Sigma}$ may be a function of $\mathbf{m}(t)$ and $\boldsymbol{\theta}$ so is stochastic too.

Extended Filtering Density: $p(\mathbf{f} | \boldsymbol{\theta})$

Represent stochastic process, \mathbf{f} as a filtering distribution with moments $\mathbf{m}(t)$ and $\mathbf{P}(t)$

Assume steady state: $d\mathbf{P}/dt = 0$.

Extended Filtering Density: $p(\mathbf{f} | \boldsymbol{\theta})$

Represent stochastic process, \mathbf{f} as a filtering distribution with moments $\mathbf{m}(t)$ and $\mathbf{P}(t)$

Assume steady state: $d\mathbf{P}/dt = 0$.

Denote covariance in steady state by $\tilde{\Sigma}$ and solve

$$\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})\tilde{\Sigma} + \tilde{\Sigma}\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})^T = -\mathbf{L}\zeta^2\mathbf{L}^T$$

Extended Filtering Density: $p(\mathbf{f} | \boldsymbol{\theta})$

Represent stochastic process, \mathbf{f} as a filtering distribution with moments $\mathbf{m}(t)$ and $\mathbf{P}(t)$

Assume steady state: $d\mathbf{P}/dt = 0$.

Denote covariance in steady state by $\tilde{\Sigma}$ and solve

$$\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})\tilde{\Sigma} + \tilde{\Sigma}\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})^T = -\mathbf{L}\zeta^2\mathbf{L}^T$$

Example of continuous Lyapunov equation; easy to solve numerically, but need an differentiable form of $\tilde{\Sigma}$

Extended Filtering Density: $p(\mathbf{f} | \boldsymbol{\theta})$

Represent stochastic process, \mathbf{f} as a filtering distribution with moments $\mathbf{m}(t)$ and $\mathbf{P}(t)$

Assume steady state: $d\mathbf{P}/dt = 0$.

Denote covariance in steady state by $\tilde{\Sigma}$ and solve

$$\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})\tilde{\Sigma} + \tilde{\Sigma}\mathbf{J}_D(\mathbf{m}, t; \boldsymbol{\theta})^T = -\mathbf{L}\zeta^2\mathbf{L}^T$$

Example of continuous Lyapunov equation; easy to solve numerically, but need an differentiable form of $\tilde{\Sigma}$

$\tilde{\Sigma}$ may be a function of $\mathbf{m}(t)$ and $\boldsymbol{\theta}$ so is stochastic too.

Transition density: $p(f_k | f_{k-1} \boldsymbol{\theta})$

We construct a discrete-time transition density using Euler-Maruyama:

$$p(f_k | f_{k-1}, \boldsymbol{\theta}) = \mathcal{N}(f_k | \mu_\Delta, \Sigma_\Delta),$$

where

$$\mu_\Delta = f_{k-1} + \mathbf{D}(f_{k-1}, t_k; \boldsymbol{\theta}) \Delta_t$$

$$\Sigma_\Delta = \tilde{\Sigma}(f_{k-1}, t_k; \boldsymbol{\theta}) - \exp(\Delta_t \mathbf{J}_D) \tilde{\Sigma}(f_{k-1}, t_k; \boldsymbol{\theta}) \exp(\Delta_t \mathbf{J}_D)^T$$

Generative model: $p(\mathbf{f} \mid \boldsymbol{\theta})$

Marginal

$$p(\mathbf{f} \mid \boldsymbol{\theta}) = p(f_0 \mid \boldsymbol{\theta}) \prod_{k=1}^T p(f_k \mid f_{k-1}, \boldsymbol{\theta})$$

Generative model: $p(\mathbf{f} \mid \boldsymbol{\theta})$

Marginal

$$p(\mathbf{f} \mid \boldsymbol{\theta}) = p(f_0 \mid \boldsymbol{\theta}) \prod_{k=1}^T p(f_k \mid f_{k-1}, \boldsymbol{\theta})$$

Additional points

$$f_k = \mathbf{f}(t_k) \quad f_{k+1} = \mathbf{f}(t_{k+1}) = \mathbf{f}(t_k + \Delta_t)$$

$$p(f_k \mid f_{k-1}) \equiv p(x_k \mid x_{k-1}, u_k) p(u_k \mid u_{k-1})$$

Variational Approximation: $q(\boldsymbol{f} \mid \boldsymbol{\theta})$

- Family of distributions parameterised by ϕ
- Needs to be flexible, sampleable and invertible (for autodifferentiation)

Variational Approximation: $q(\boldsymbol{f} \mid \boldsymbol{\theta})$

- Family of distributions parameterised by ϕ
- Needs to be flexible, sampleable and invertible (for autodifferentiation)
- Look to (Bayesian) neural networks and other deep models

Variational Approximation: $q(f | \theta)$

- Family of distributions parameterised by ϕ
- Needs to be flexible, sampleable and invertible (for autodifferentiation)
- Look to (Bayesian) neural networks and other deep models
- Need to encode temporal (recurrent) structure

Variational Approximation: $q(f | \theta)$

- Family of distributions parameterised by ϕ
- Needs to be flexible, sampleable and invertible (for autodifferentiation)
- Look to (Bayesian) neural networks and other deep models
- Need to encode temporal (recurrent) structure
 - RNNs with priors on the weights

Variational Approximation: $q(f | \theta)$

- Family of distributions parameterised by ϕ
- Needs to be flexible, sampleable and invertible (for autodifferentiation)
- Look to (Bayesian) neural networks and other deep models
- Need to encode temporal (recurrent) structure
 - RNNs with priors on the weights
 - Normalising flows

Parametrising q with an RNN

Pros

- Can represent high dimensional recurrent structure
- Bi-directional RNN can represent first-order Markov properties of model
- Priors over weights and optimise in weight-space

Parametrising q with an RNN

Pros

- Can represent high dimensional recurrent structure
- Bi-directional RNN can represent first-order Markov properties of model
- Priors over weights and optimise in weight-space

Cons

- Need to sample sequentially
- BPTT inefficient for propagation of gradients
- Doesn't handle latent dimensions well

Inverse Autoregressive Flows

- Want to define a distribution for \mathbf{f} that is invertible and expressive
- Inverse autoregressive flows (IAFs) introduce a base random vector $\mathbf{z}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Layers of this random variable are shifted and scaled through 1-D convolutions to create autoregressive model
- Very flexible, and can sample in parallel

Inverse Autoregressive Flows

- Want to define a distribution for \mathbf{f} that is invertible and expressive
- Inverse autoregressive flows (IAFs) introduce a base random vector $\mathbf{z}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Layers of this random variable are shifted and scaled through 1-D convolutions to create autoregressive model
- Very flexible, and can sample in parallel

Autoregressive Flows

$$\mathbf{z}_j = \boldsymbol{\sigma}_j \odot \mathbf{z}_{j-1} + \boldsymbol{\mu}_j$$

where

$$[\boldsymbol{\mu}_j, \mathbf{s}_j] = \text{AUTOREGRESSIVENN}(\mathbf{z}_{j-1}, \mathbf{y}, \boldsymbol{\theta}) \text{ and } \boldsymbol{\sigma}_j = \log(1 + \exp \mathbf{s}_j)$$

$$\mathbf{f} = \text{BIJECTOR}(\mathbf{z}_N)$$

Autoregressive Neural Network Layers

Algorithm 2 j^{th} Autoregressive Neural Network Layer

$$\xi^{(0a)} \leftarrow \text{CONV1D}(\mathbf{z}_{j-1}, \mathbf{y}, \mathbf{t})$$

$$\xi^{(0b)} \leftarrow \text{DENSE}(\boldsymbol{\theta})$$

$$\xi^{(1)} \leftarrow \text{ELU}(\xi^{(0a)} + \xi^{(0b)})$$

for $i = 2 \dots n_\ell$ **do**

$$\xi^{(i)} \leftarrow \text{BATCHNORM}(\text{CONV1D}(\text{ELU}(\xi^{(i-1)})))$$

end for

$$[\boldsymbol{\mu}_j, \mathbf{s}_j] \leftarrow \text{CONV1D}(\xi^{(n_\ell)})$$

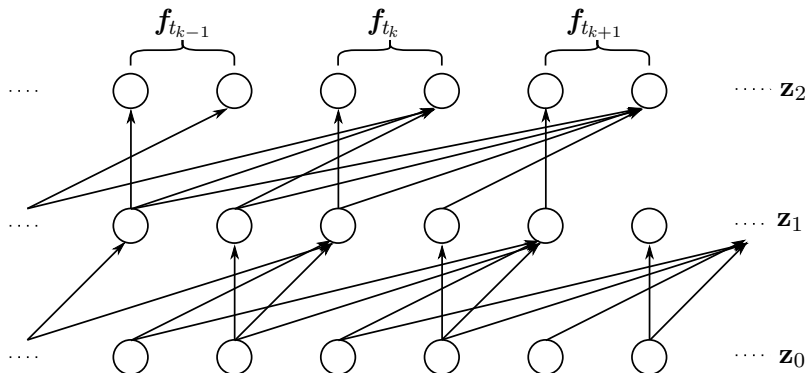
$$\boldsymbol{\sigma}_j \leftarrow \text{SOFTPLUS}(\mathbf{s}_j)$$

$$\mathbf{z}_j \leftarrow \boldsymbol{\sigma}_j \odot \mathbf{z}_{j-1} + \boldsymbol{\mu}_j$$

Locally Masked Multivariate Inverse Autoregressive Flows

- Passing the entire flow vector, \mathbf{z}_j can lead to complex (unrepresentative temporal dependencies)
- We use a local receptive field to update flow layers (similar to Wavenet)
- Rolled out multidimensional state (f_k) in sequence
- Hacks and tricks to approximate locally informed flow state

Locally Masked Multivariate Inverse Autoregressive Flows



Variational Log Density

$$\log q(\boldsymbol{f} \mid \boldsymbol{\theta}) = -\frac{1}{2} \boldsymbol{z}_0^T \boldsymbol{z}_0 + \frac{T}{2} \log 2\pi + T \sum_{i=1}^N \log \boldsymbol{\sigma}_j + \log |\boldsymbol{J}_{-1}(\boldsymbol{f})|$$

Variational Log Density

$$\log q(\mathbf{f}^{(i)} | \boldsymbol{\theta}^{(i)}) = -\frac{1}{2} \mathbf{z}_0^{(i)\text{T}} \mathbf{z}_0^{(i)} + \frac{T}{2} \log 2\pi + T \sum_{i=1}^N \log \boldsymbol{\sigma}_j^{(i)} + \log |\mathbf{J}_{-1}(\mathbf{f}^{(i)})|$$

$$\mathbf{z}_0^{(i)} \sim \text{N}(\mathbf{0}, \mathbf{I})$$

Unbiased Evidence Lower Bound

ELBO

$$\hat{\mathcal{L}}(\phi) = \frac{1}{n_s} \sum_{i=1}^{n_s} \log \frac{p(\boldsymbol{\theta}^{(i)}) p(\mathbf{f}^{(i)} | \boldsymbol{\theta}^{(i)}) p(\mathbf{y} | \mathbf{f}^{(i)}, \boldsymbol{\theta}^{(i)})}{q(\boldsymbol{\theta}^{(i)}) q(\mathbf{f}^{(i)} | \boldsymbol{\theta}^{(i)})}$$

where $\mathbf{f}^{(i)} \sim q(\mathbf{f} | \boldsymbol{\theta}^{(i)})$ and $\boldsymbol{\theta}^{(i)} \sim q(\boldsymbol{\theta})$ $i = 1, \dots, n_s$

Contents

- 1 Stochastic Differential Equations and Gaussian Processes
- 2 Variational Solutions to Non-Linear Latent Force Models
- 3 Approximate Gaussian Processes**
- 4 Some Results
- 5 Recap
- 6 Open Issues

Exponential Gaussian Process

$$f(t) \sim \text{GP}(0, \sigma_f^2 \exp(-\lambda|t - t'|))$$

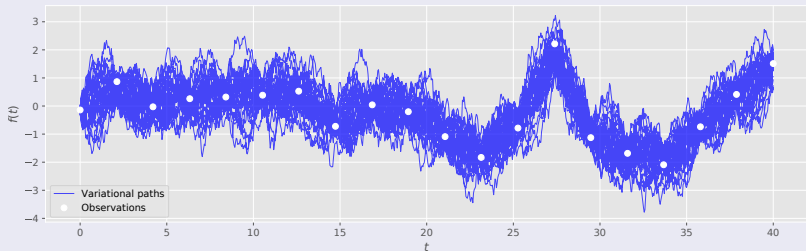
Exponential Gaussian Process

$$f(t) \sim \text{GP}(0, \sigma_f^2 \exp(-\lambda|t - t'|))$$

$$df(t) = -\lambda f(t) + 2\sigma_f^2 \lambda d\beta(t)$$

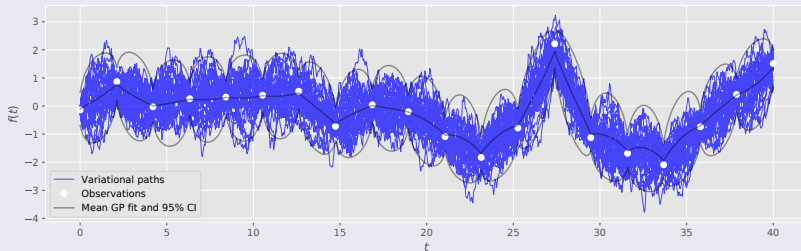
Exponential Gaussian Process

Samples from $q(f | \theta)$



Exponential Gaussian Process

Samples from $q(f | \theta)$ and mean and covariance of $p(f | \mathbf{y}, \theta)$



Model Criticism

- Visual confirmation fine, it *looks* like a good estimate
- Empirical evidence for reliability needed
- Map corresponding samples from p and q to RKHS
- Two-sample test with MMD to validate approximation

Maximum Mean Discrepancy (MMD)

- MMD is a measure of distance between two probabilities
- Samples are embedded in an RKHS
- Metric describes distance as some norm in the RKHS
- Two-sample testing for $H_0 : \hat{\text{MMD}}^2(\mu_p, \mu_q) = 0$
- Gretton, et al. A kernel two-sample test. *JMLR*, 2012

Model Criticism

MMD² values comparing samples from $q(f | \theta)$ and $p(f | \mathbf{y}, \theta)$

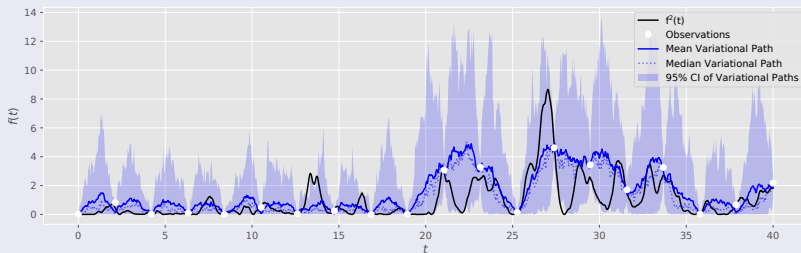
Fit on different number of observations, N

Epoch	10	100	500	1 000	2 500	25 000
$N = 6$	0.1111	0.1267	0.0596	0.0484	0.0556	–
$N = 20$	0.2731	0.1147	0.0654	0.0696	0.0471	0.0316

Thresholds for rejection at 95% confidence: 0.0371 ($N = 6$) and 0.0337 ($N = 20$)

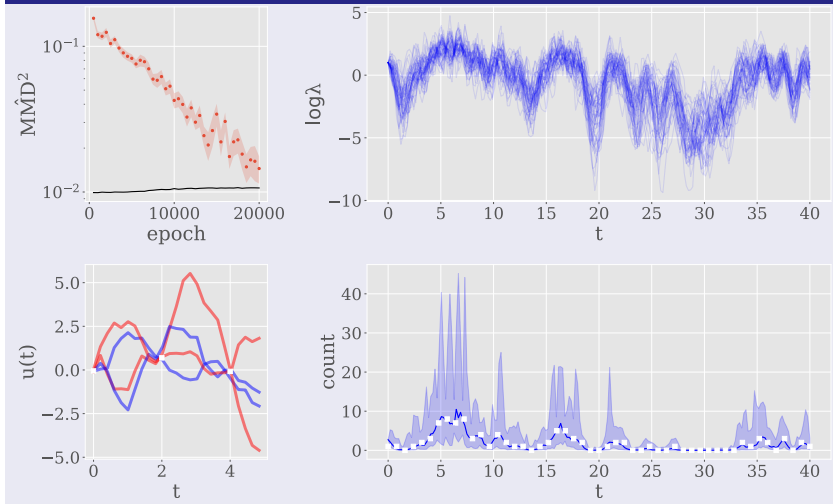
Matérn Covariances and Non-Gaussian Likelihoods

Summary statistics for $p(\mathbf{y} | f^{(i)}, \boldsymbol{\theta}), f^{(i)} \sim q(f | \boldsymbol{\theta})$, plotted against true latent function, $f^2(t)$



Matérn Covariances and Non-Gaussian Likelihoods

Approximating Matérn $3/2$ GPs



Contents

- 1 Stochastic Differential Equations and Gaussian Processes
- 2 Variational Solutions to Non-Linear Latent Force Models
- 3 Approximate Gaussian Processes
- 4 Some Results**
- 5 Recap
- 6 Open Issues

Toy Non-Linear ODE

$$\frac{d}{dt}x(t) = -\frac{2}{3}\sin(\omega x(t)) + u(t)$$

Toy Non-Linear ODE

$$\frac{d}{dt}x(t) = -\frac{2}{3}\sin(\omega x(t)) + u(t)$$

$$u(t) \sim \text{GP}(0, k_{\nu=1/2}(t, t'))$$

Toy Non-Linear ODE

$$\frac{d}{dt}x(t) = -\frac{2}{3}\sin(\omega x(t)) + u(t)$$

$$u(t) \sim \text{GP}(0, k_{\nu=1/2}(t, t'))$$

$$\underbrace{\frac{d}{dt} \begin{bmatrix} x(t) \\ u(t) \end{bmatrix}}_{\mathbf{f}(t)} = \underbrace{\begin{bmatrix} -2\cos(\omega f_1)/3 + f_2 \\ -\lambda f_2 \end{bmatrix}}_{D(\mathbf{f}(t), \boldsymbol{\theta})} + \underbrace{\begin{bmatrix} 0 \\ 1 \end{bmatrix}}_{\mathbf{L}} w(t),$$

Toy Non-Linear ODE

The Jacobian of $\mathbf{D}(\mathbf{f}(t), \boldsymbol{\theta})$ w.r.t \mathbf{f} is defined:

$$\mathbf{J}_D(\mathbf{f}(t)) = \begin{bmatrix} 2\omega \sin(\omega f_1)/3 & 1 \\ 0 & -\lambda \end{bmatrix},$$

Toy Non-Linear ODE

The Jacobian of $\mathbf{D}(\mathbf{f}(t), \boldsymbol{\theta})$ w.r.t \mathbf{f} is defined:

$$\mathbf{J}_D(\mathbf{f}(t)) = \begin{bmatrix} 2\omega \sin(\omega f_1)/3 & 1 \\ 0 & -\lambda \end{bmatrix},$$

and steady state covariance $\tilde{\boldsymbol{\Sigma}}$ such that

$$\mathbf{J}_D(\mathbf{f}(t))\tilde{\boldsymbol{\Sigma}} + \tilde{\boldsymbol{\Sigma}}[\mathbf{J}_D(\mathbf{f}(t))]^\top + 2\lambda\sigma^2\mathbf{L}\mathbf{L}^\top = 0$$

Toy Non-Linear ODE

The Jacobian of $\mathbf{D}(\mathbf{f}(t), \boldsymbol{\theta})$ w.r.t \mathbf{f} is defined:

$$\mathbf{J}_D(\mathbf{f}(t)) = \begin{bmatrix} 2\omega \sin(\omega f_1)/3 & 1 \\ 0 & -\lambda \end{bmatrix},$$

and steady state covariance $\tilde{\boldsymbol{\Sigma}}$ such that

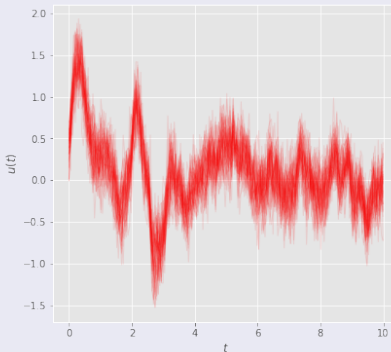
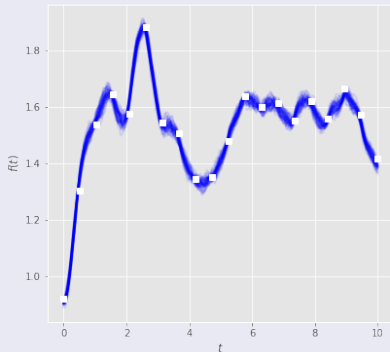
$$\mathbf{J}_D(\mathbf{f}(t))\tilde{\boldsymbol{\Sigma}} + \tilde{\boldsymbol{\Sigma}}[\mathbf{J}_D(\mathbf{f}(t))]^\top + 2\lambda\sigma^2\mathbf{L}\mathbf{L}^\top = 0$$

is

$$\tilde{\boldsymbol{\Sigma}} = \begin{bmatrix} \frac{\sigma^2\lambda}{2\lambda\omega \sin(\omega f_1)(2\omega \sin(\omega f_1)/3 - \lambda)/3} & \frac{\sigma^2\lambda}{\lambda^2 - 2\lambda\omega \sin(\omega f_1)/3} \\ \frac{\sigma^2\lambda}{\lambda^2 - 2\lambda\omega \sin(\omega f_1)/3} & \sigma^2 \end{bmatrix}$$

Toy Non-Linear ODE

Samples of joint posterior of x , u and θ



- Multi-output system with non-linear dependency on input

$$\frac{d}{dt}x_d(t) = a_d - b_dx_d(t) + s_d\frac{u(t)}{\gamma_d + u(t)}$$

- x_d is a model of gene expression, that's noisily observable
- u models the concentration of the transcription factor regulating the observed genes

- Multi-output system with non-linear dependency on input

$$\frac{d}{dt}x_d(t) = a_d - b_dx_d(t) + s_d\frac{u(t)}{\gamma_d + u(t)}$$

- $x_d(t), u(t) > 0, \theta_d = \{a_d, b_d, s_d, \gamma_d\}$
- $d = 1, \dots, ?$

Real World: Gene Expression Data

- Multi-output system with non-linear dependency on input

$$\frac{d}{dt}x_d(t) = a_d - b_d x_d(t) + s_d \frac{u(t)}{\gamma_d + u(t)}$$

- $x_d(t), u(t) > 0, \theta_d = \{a_d, b_d, s_d, \gamma_d\}$
- $d = 1, \dots, ?$
- Place GP prior over $\exp u(t)$

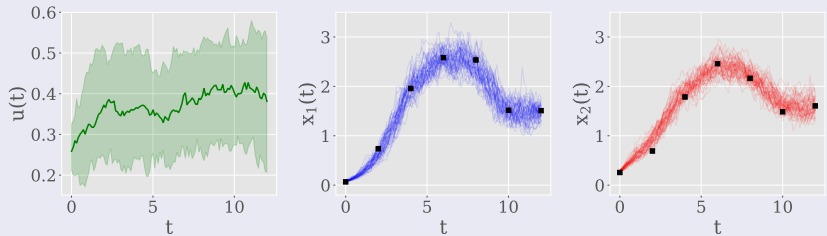
- Multi-output system with non-linear dependency on input

$$\frac{d}{dt}x_d(t) = a_d - b_d x_d(t) + s_d \frac{u(t)}{\gamma_d + u(t)}$$

- $x_d(t), u(t) > 0, \theta_d = \{a_d, b_d, s_d, \gamma_d\}$
- $d = 1, \dots, ?$
- Place GP prior over $\exp u(t)$
- Infer θ_d simultaneously

Real World: Gene Expression Data

Inferred TF concentration and predicted gene expressions for TNFRSF10b (blue) and p26 sesn1 (red)



Contents

- 1 Stochastic Differential Equations and Gaussian Processes
- 2 Variational Solutions to Non-Linear Latent Force Models
- 3 Approximate Gaussian Processes
- 4 Some Results
- 5 Recap**
- 6 Open Issues

- GP priors on non-linear forced models are non-linear SDEs

- GP priors on non-linear forced models are non-linear SDEs
- Filtering approaches struggle with joint parameter estimation
- Sequential inference slow for propagating gradients

- GP priors on non-linear forced models are non-linear SDEs
- Filtering approaches struggle with joint parameter estimation
- Sequential inference slow for propagating gradients
- With inverse autoregressive flows we can batch sample time-series

- GP priors on non-linear forced models are non-linear SDEs
- Filtering approaches struggle with joint parameter estimation
- Sequential inference slow for propagating gradients
- With inverse autoregressive flows we can batch sample time-series
- Can construct approximate model for joint posterior and infer state, input and parameters by optimising NN weights

- GP priors on non-linear forced models are non-linear SDEs
- Filtering approaches struggle with joint parameter estimation
- Sequential inference slow for propagating gradients
- With inverse autoregressive flows we can batch sample time-series
- Can construct approximate model for joint posterior and infer state, input and parameters by optimising NN weights
Approximation of GPs quantifiably good

Contents

- 1 Stochastic Differential Equations and Gaussian Processes
- 2 Variational Solutions to Non-Linear Latent Force Models
- 3 Approximate Gaussian Processes
- 4 Some Results
- 5 Recap
- 6 Open Issues**

Open Issue: Calculating Steady State Covariance

Solving continuous Lyapunov equation

$$\mathbf{J}_D(\mathbf{f}, t; \boldsymbol{\theta}) \tilde{\boldsymbol{\Sigma}} + \tilde{\boldsymbol{\Sigma}} \mathbf{J}_D(\mathbf{f}, t; \boldsymbol{\theta})^T = -\mathbf{L} \boldsymbol{\Sigma}^2 \mathbf{L}^T,$$

is possible to do for fixed values of \mathbf{f} , t , and $\boldsymbol{\theta}$ using numerical solvers, but hard to do online, so no gradients !

Manually solving is increasingly difficult with increase in dimension: solution is system of $d(d-1)/2$ equations

Open Issue: Dimensionality

- Smoother GPs have more orders of differentiation
- Approximations of infinitely-differentiable covariance functions, e.g. periodic, Gaussian RBF, require series approximation
- State dimension proportional to series threshold

References

- M. A. Alvarez, D. Luengo, and N. D. Lawrence. Linear latent force models using Gaussian processes. *IEEE transactions on pattern analysis and machine intelligence*, 35(11):2693–2705, 2013.
- D. Duvenaud and R. P. Adams. Black-box stochastic variational inference in five lines of Python. In *NIPS Workshop on Black-box Learning and Inference*, 2015.
- A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar): 723–773, 2012.
- J. Hartikainen, M. Seppänen, and S. Särkkä. State-space inference for non-linear latent force models with application to satellite orbit prediction. In *Proceedings of the 29th International Conference on Machine Learning*, pages 723–730, 2012.
- R. Ranganath, S. Gerrish, and D. Blei. Black box variational inference. In *Artificial Intelligence and Statistics*, pages 814–822, 2014.