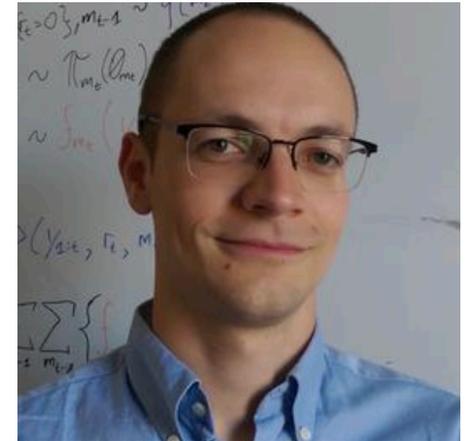




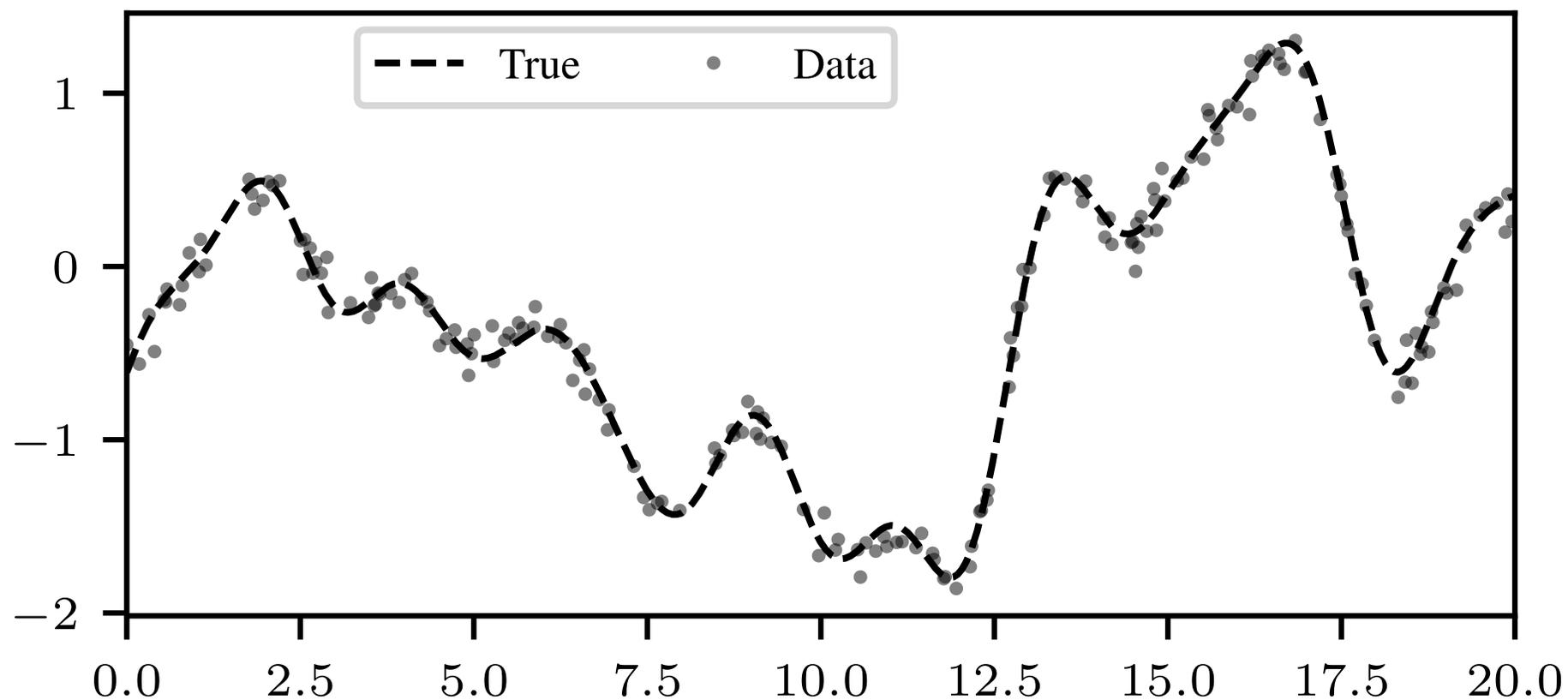
UCL

Robust and Conjugate Gaussian Process Regression

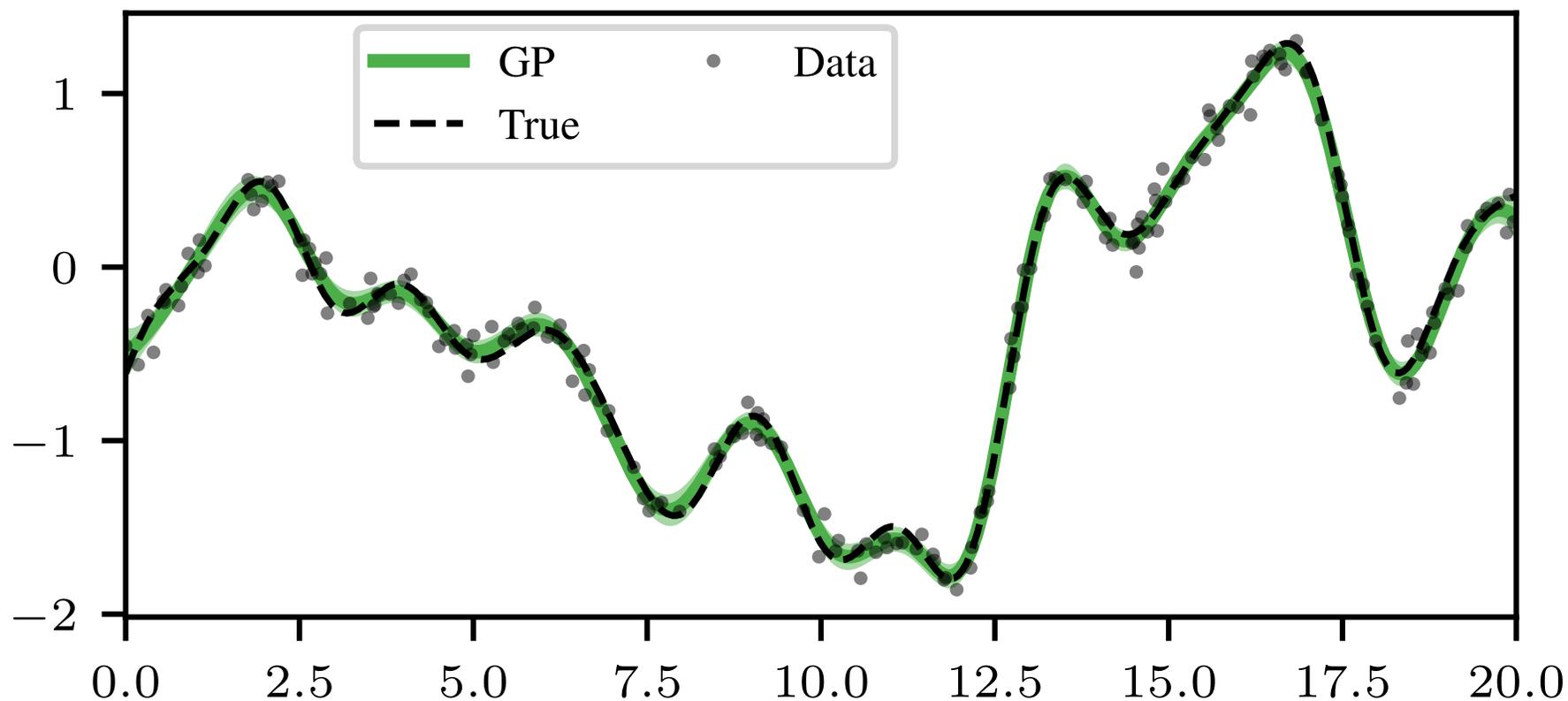
Dr François-Xavier Briol
Department of Statistical Science
University College London



A synthetic problem



GP regression on the synthetic problem



[I am being a bad Bayesian by plotting only the mean... sorry....]

Gaussian process regression

- **Regression problem:** Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be some unknown function of interest. we have access to data $\{x_i, y_i\}_{i=1}^n$ where:

$$y_i = f(x_i) + \epsilon_i$$

- Two main assumptions:

$$f \sim GP(m, k) \quad \leftarrow \quad \text{“Prior”}$$

$$\epsilon_i \sim N(0, \sigma^2) \quad \leftarrow \quad \text{“Likelihood/
Observation
Model”}$$

Gaussian process regression

$$f \sim GP(m, k)$$

- A GP is a stochastic process often used as prior in Bayesian (non-parametric) inference.
- It is fully determined by its mean function $m : \mathcal{X} \rightarrow \mathbb{R}$ and covariance function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$
- The GP posterior can then be obtained in closed form as follows:

$$p(\mathbf{f} | \mathbf{y}, \mathbf{x}) = N(\mathbf{f}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

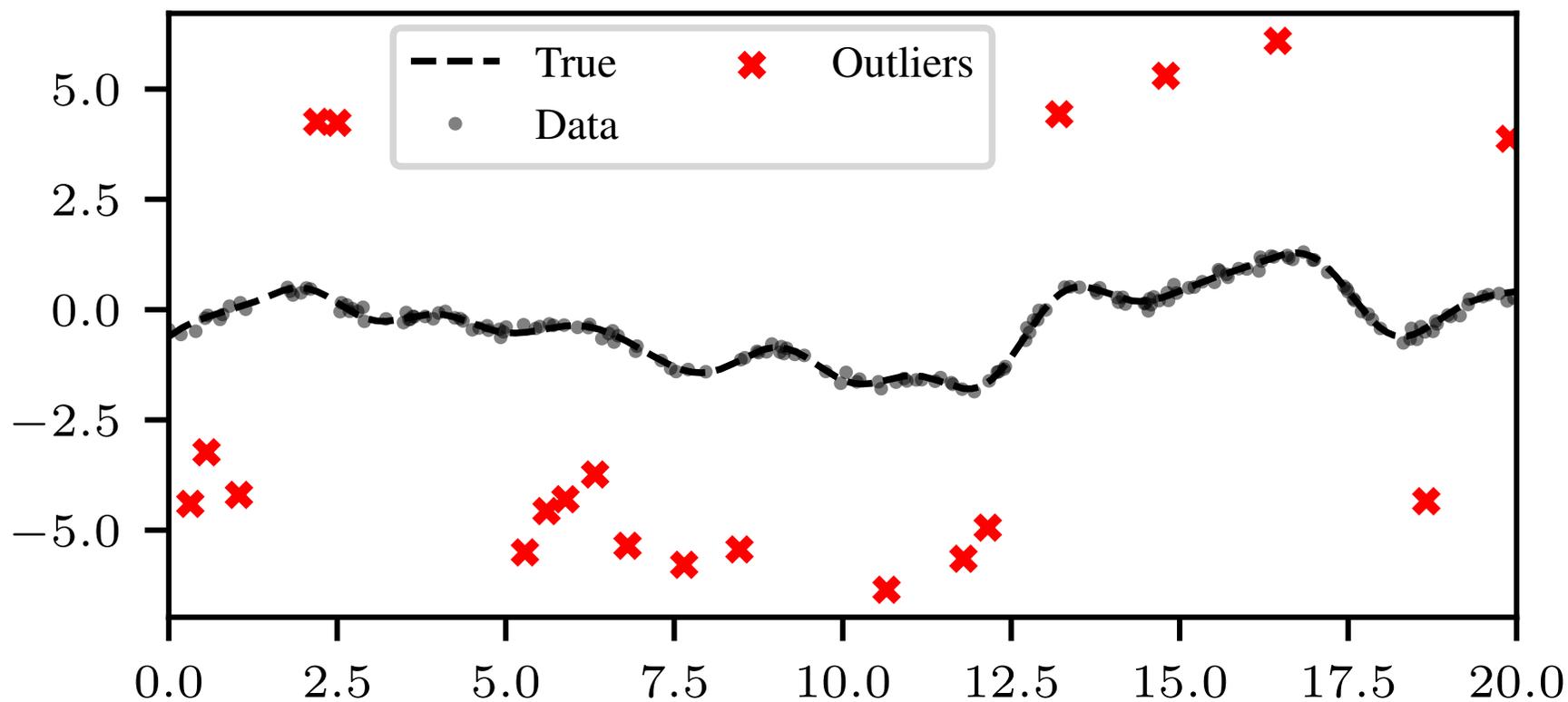
$$\boldsymbol{\mu} = \mathbf{m} + K(K + \sigma^2 I_n)^{-1}(\mathbf{y} - \mathbf{m})$$

$$\boldsymbol{\Sigma} = K(K + \sigma^2 I_n)^{-1} \sigma^2 I_n$$

Why Gaussian processes?

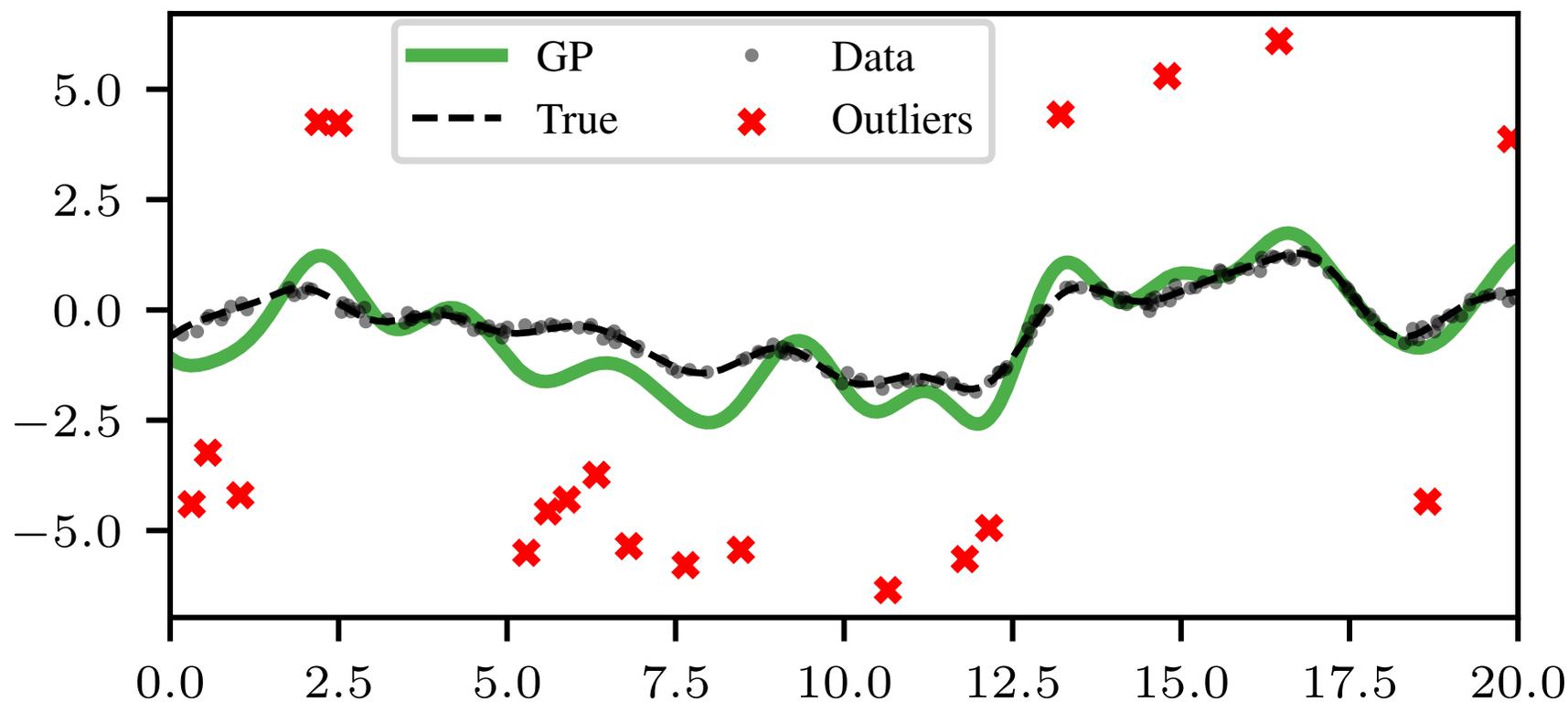
1. A very **flexible and interpretable model** through the choice of prior mean function m and covariance k function (e.g. smoothness, periodicity, sparsity, etc...).
2. We get a posterior on f which quantifies **epistemic uncertainty**.
3. We can do **exact conditioning** through Gaussian conjugacy! We therefore don't need to do any approximation of the posterior!

Regression in the “real world”



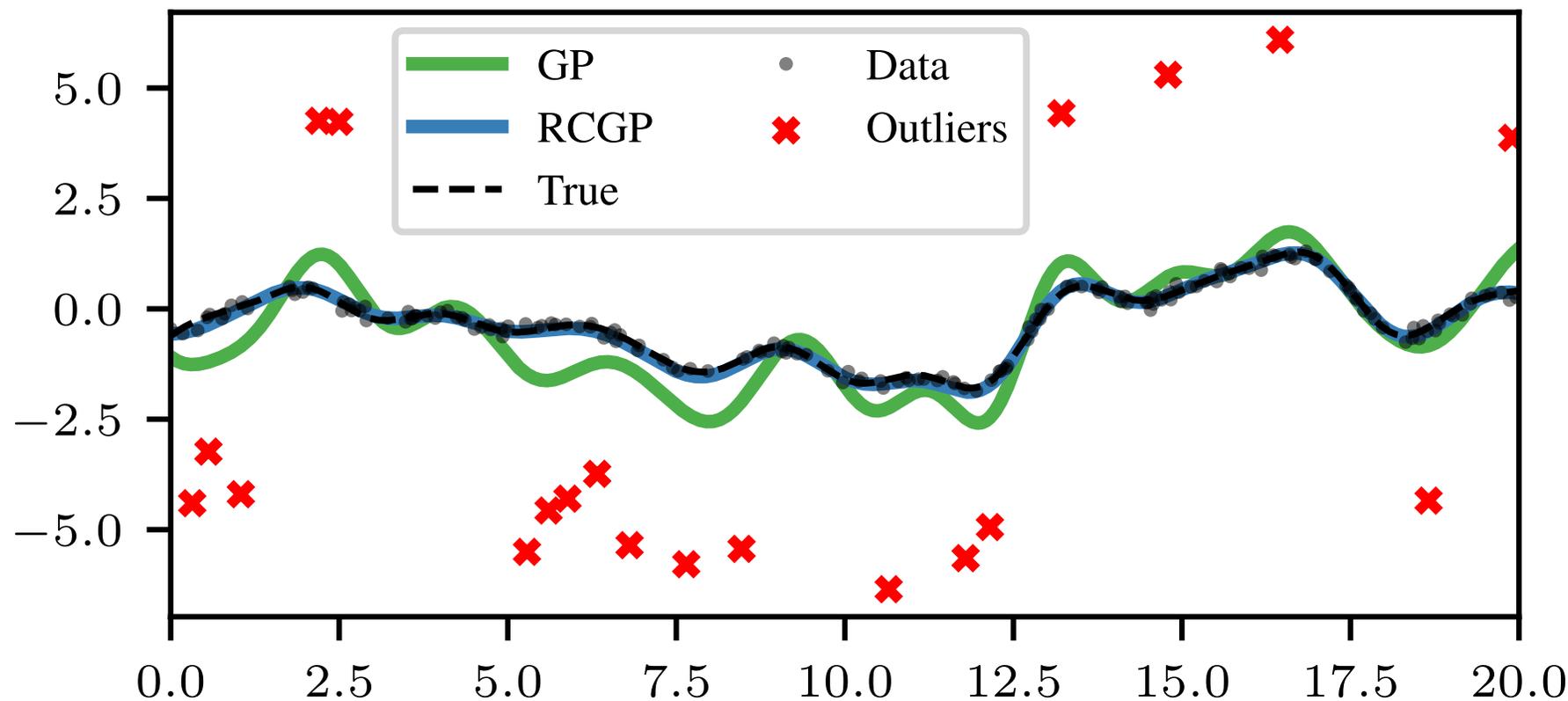
~~$$\epsilon_i \sim \mathcal{N}(0, \sigma^2)$$~~

GP regression in the “real world”



We assumed $\epsilon_i \sim N(0, \sigma^2)$ but its wrong...

Our goal: robust GP regression



We assumed $\epsilon_i \sim N(0, \sigma^2)$ but its wrong...

Existing work

IEEE TRANSACTIONS ON BIOMEDICAL ENGINEERING, VOL. 55, NO. 9, SEPTEMBER 2008

Gaussian process regression with Student-*t* likelihood

Jarno Vanhatalo
Department of Biomedical Engineering
and Computational Science
Helsinki University of Technology
Finland
jarno.vanhatalo@tkk.fi

Pasi Jylänki
Department of Biomedical Engineering
and Computational Science
Helsinki University of Technology
Finland
pasi.jylanki@tkk.fi

Aki Vehtari
Department of Biomedical Engineering
and Computational Science
Finland
Helsinki University of Technology
aki.vehtari@tkk.fi

Gaussian Process Robust Regression for Noisy Heart Rate Data

Oliver Stegle*, Sebastian V. Fallert, David J. C. MacKay, and Søren Brage

Corruption-Tolerant Gaussian Process Bandit Optimization

Ilija Bogunovic
ETH Zürich

Andreas Krause
ETH Zürich

Jonathan Scarlett
National University of Singapore

Robust Gaussian Process Regression with a Bias Model

Chiwoo Park

Department of Industrial and Manufacturing Engineering
Florida State University
Tallahassee, FL 32310, USA

CPARK5@FSU.EDU

Robust and Scalable Gaussian Process Regression and Its Applications

Yifan Lu¹, Jiayi Ma^{1*}, Leyuan Fang², Xin Tian¹, and Junjun Jiang³

¹ Wuhan University, China ² Hunan University, China ³ Harbin Institute of Technology, China
{lyf048, xin.tian}@whu.edu.cn, {jyma2010, fangleyuan}@gmail.com, jiangjunjun@hit.edu.

ROBUST GAUSSIAN PROCESS REGRESSION WITH HUBER LIKELIHOOD

*

BY POOJA ALGIKAR^{1,a}, LAMINE MILI^{2,b}

**Robust Gaussian process regression with
G-confluent likelihood**
Martin Lindfors^{*,**} Tianshi Chen^{**} Christian A. Naesseth^{***}

Identification of robust Gaussian Process Regression with noisy input using EM algorithm

Atefeh Daemi, Yousef Alipouri, Biao Huang^{*}

Department of Chemical and Materials Engineering, University of Alberta, Edmonton, Alberta, T6G 1H9, Canada

Robust Gaussian process modeling using EM algorithm

Rishik Ranjan^a, Biao Huang^{a,*}, Alireza Fatehi^{a,b}

^a Department of Chemical and Materials Engineering, University of Alberta, Edmonton, Alberta, Canada T6G 2G6

^b APAC Research Group, Industrial Control Center of Excellence, Faculty of Electrical Engineering, K.N. Toosi University of Technology, Tehran 16317-14191, Iran

Robust Bayesian Optimization with Student-*t* Likelihood

Ruben Martinez-Cantin
SigOpt Inc.

Centro Universitario de la Defensa, Zaragoza

RUBEN@SIGOPT.COM

Michael McCourt
Kevin Tee

SigOpt Inc.

MCCOURT@SIGOPT.COM

KEVIN@SIGOPT.COM

Robust Regression with Twinned Gaussian Processes

Andrew Naish-Guzman & Sean Holden
Computer Laboratory
University of Cambridge
Cambridge, CB3 0FD, United Kingdom
{agpn2, sbh11}@cl.cam.ac.uk

Robust Gaussian Process Regression with the Trimmed Marginal Likelihood

Daniel Andrade¹

Akiko Takeda^{2,3}

Robust Gaussian process regression based on iterative trimming

Zhao-Zhou Li^{a,*}, Lu Li^{b,c}, Zhengyi Shao^{b,d}

Existing work

- There are two main categories:
 1. **Extended models:** i.e. use more flexible likelihood model to ensure that the outliers are well modelled. Examples include Student-t, mixtures, Laplace, etc...

$$\epsilon \sim P \neq N(0, \sigma^2)$$

2. **Outlier detection/removal:** i.e. find the outliers, remove them, then fit a standard GP model (with Gaussian observations) to the rest of the data.

Issues with existing work

- The main issue with all of the methods above is that they are **very slow!**
- This is because they all **break Gaussian conjugacy** and so we must resort to approximate methods such as MCMC, Laplace or Variational Bayes.



Issues with existing work

- The main issue with all of the methods above is that they are **very slow!**
- This is because they all **break Gaussian conjugacy** and so we must resort to approximate methods such as MCMC, Laplace or Variational Bayes.

	GP	t-GP	m-GP	
Synthetic	1.5 (0.1)	2.2 (0.0)	3.0 (0.0)	$n = 300, d = 1$
Boston	1.9 (0.5)	30.7 (6.1)	16.7 (1.7)	$n = 506, d = 13$
Energy	3.8 (0.9)	34.0 (11)	33.8 (0.3)	$n = 768, d = 8$
Yacht	1.6 (0.3)	5.6 (0.7)	4.5 (0.4)	$n = 308, d = 6$

Table: Fitting time in second, including time for hyper parameter optimisation.

Goal of this project

- **Robust Gaussian Process regression without the additional computational cost!**

arXiv > stat > arXiv:2311.00463

Statistics > Machine Learning

[Submitted on 1 Nov 2023]

Robust and Conjugate Gaussian Process Regression

Matias Altamirano, François-Xavier Briol, Jeremias Knoblauch

Bayesian inference for regression

- In standard GP regression, we do:

Posterior Likelihood Prior


$$p(\mathbf{f} | \mathbf{y}, \mathbf{x}) \propto p(\mathbf{y} | \mathbf{f}, \mathbf{x}) \times p(\mathbf{f} | \mathbf{x})$$

$$\mathbf{x} = (x_1, \dots, x_n)^\top$$

$$\mathbf{f} = (f(x_1), \dots, f(x_n))^\top$$

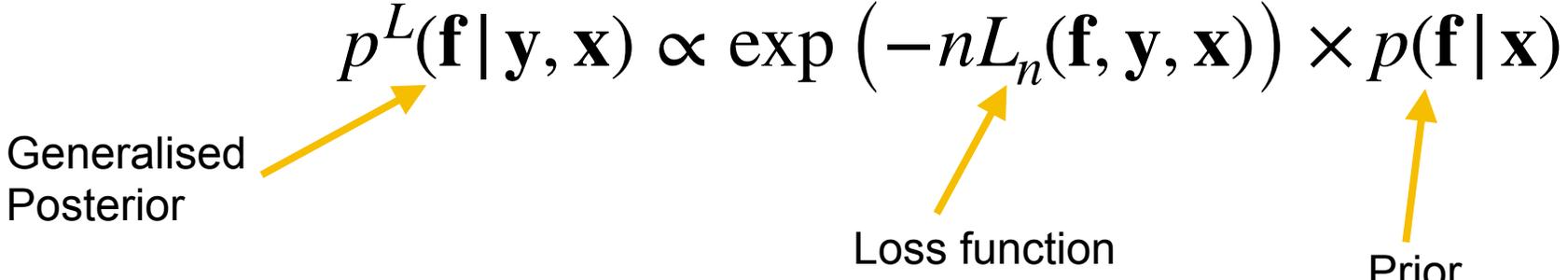
$$\mathbf{y} = (y_1, \dots, y_n)^\top$$

Generalised Bayesian inference for regression

- In standard GP regression, we do:


$$p(\mathbf{f} | \mathbf{y}, \mathbf{x}) \propto p(\mathbf{y} | \mathbf{f}, \mathbf{x}) \times p(\mathbf{f} | \mathbf{x})$$

- We take a generalised Bayesian approach and do:


$$p^L(\mathbf{f} | \mathbf{y}, \mathbf{x}) \propto \exp(-nL_n(\mathbf{f}, \mathbf{y}, \mathbf{x})) \times p(\mathbf{f} | \mathbf{x})$$

Standard vs Generalised Bayesian inference

$$p^L(\mathbf{f} | \mathbf{y}, \mathbf{x}) \propto \exp(-nL_n(\mathbf{f}, \mathbf{y}, \mathbf{x})) \times p(\mathbf{f} | \mathbf{x})$$

- Standard Bayes is recovered by taking

$$L_n(\mathbf{f}, \mathbf{y}, \mathbf{x}) = -\frac{1}{n} \log p(\mathbf{y} | \mathbf{f}, \mathbf{x})$$

- This is **optimal**, but **only when the model is well-specified**; i.e. when $\epsilon \sim N(0, \sigma^2)$!



Key Question: What should we do when this is not the case??

Generalised Bayesian inference

$$p^L(\mathbf{f} | \mathbf{y}, \mathbf{x}) \propto \exp(-nL_n(\mathbf{f}, \mathbf{y}, \mathbf{x})) \times p(\mathbf{f} | \mathbf{x})$$

- We can choose the loss function to induce **robustness to mild model misspecification**.
- Common choice is a loss based on the Beta divergence. But we have already seen other examples (i.e. MMD) this week.
- In this talk, we will also choose the loss function for computational convenience!

Bissiri, P., Holmes, C., & Walker, S. (2016). A general framework for updating belief distributions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 78, 1103–1130.

Knoblauch, J., Jewson, J., & Damoulas, T. (2022). An optimization-centric view on Bayes' rule: reviewing and generalizing variational inference. *Journal of Machine Learning Research*, 23(132), 1–109.

Score-matching and generalisations

- The score-matching divergence is given by:

$$D(p || q) := \mathbb{E}_{X \sim q} [\|(\nabla \log p - \nabla \log q)(X)\|_2^2]$$

- We consider a weighted generalisation based on $w : \mathcal{X} \rightarrow \mathbb{R}$:

$$D(p || q) := \mathbb{E}_{X \sim q} [\|w(\nabla \log p - \nabla \log q)(X)\|_2^2]$$

Hyvärinen, A. (2006). Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6, 695–708.

Barp, A., Briol, F.-X., Duncan, A. B., Girolami, M., & Mackey, L. (2019). Minimum Stein discrepancy estimators. *Neural Information Processing Systems*, 12964–12976.

Score-matching and generalisations

- For regression setting, we need to extend this divergence (now $w : \mathcal{X} \times \mathbb{R} \rightarrow \mathbb{R}$):

$$D(p || q) := \mathbb{E}_{X \sim q_x} \left[\mathbb{E}_{Y \sim q(\cdot | X)} \left[\|(w(\nabla \log p - \nabla \log q))(X, Y)\|_2^2 \right] \right]$$

- With integration by part and replacing q by our samples, we get that D is equal to the following loss up to some additive constant which does not depend on f:

$$L_n^w(\mathbf{f}, \mathbf{y}, \mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \left((w \nabla \log p_f)^2 + 2 \nabla_y (w^2 \nabla \log p_f) \right) (x_i, y_i)$$

Likelihood



RCGPs are conjugate!

- Suppose $f \sim GP(m, k)$ and $\epsilon \sim N(0, \sigma^2 I_n)$, then the GP and RCGP posteriors are:

Standard GP

$$p(\mathbf{f} | \mathbf{y}, \mathbf{x}) = N(\mathbf{f}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\mu} = \mathbf{m} + K(K + \sigma^2 I_n)^{-1}(\mathbf{y} - \mathbf{m})$$

$$\boldsymbol{\Sigma} = K(K + \sigma^2 I_n)^{-1} \sigma^2 I_n$$


$$K_{ij} = k(x_i, x_j)$$



Identity matrix

RCGPs are conjugate!

- Suppose $f \sim GP(m, k)$ and $\epsilon \sim N(0, \sigma^2 I_n)$, then the GP and RCGP posteriors are:

Standard GP

$$p(\mathbf{f} | \mathbf{y}, \mathbf{x}) = N(\mathbf{f}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\mu} = \mathbf{m} + K(K + \sigma^2 I_n)^{-1}(\mathbf{y} - \mathbf{m})$$

$$\boldsymbol{\Sigma} = K(K + \sigma^2 I_n)^{-1} \sigma^2 I_n$$

RCGP

$$p^w(\mathbf{f} | \mathbf{y}, \mathbf{x}) = N(\mathbf{f}; \boldsymbol{\mu}^R, \boldsymbol{\Sigma}^R)$$

$$\boldsymbol{\mu}^R = \mathbf{m} + K(K + \sigma^2 J_w)^{-1}(\mathbf{y} - \mathbf{m}_w)$$

$$\boldsymbol{\Sigma}^R = K(K + \sigma^2 J_w)^{-1} \sigma^2 J_w$$

RCGPs are conjugate!

- Suppose $f \sim GP(m, k)$ and $\epsilon \sim N(0, \sigma^2 I_n)$, then the GP and RCGP posteriors are:

Standard GP

$$p(\mathbf{f} | \mathbf{y}, \mathbf{x}) = N(\mathbf{f}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\mu} = \mathbf{m} + K(K + \sigma^2 I_n)^{-1}(\mathbf{y} - \mathbf{m})$$

$$\boldsymbol{\Sigma} = K(K + \sigma^2 I_n)^{-1} \sigma^2 I_n$$

RCGP

$$p^w(\mathbf{f} | \mathbf{y}, \mathbf{x}) = N(\mathbf{f}; \boldsymbol{\mu}^R, \boldsymbol{\Sigma}^R)$$

$$\boldsymbol{\mu}^R = \mathbf{m} + K(K + \sigma^2 J_w)^{-1}(\mathbf{y} - \mathbf{m}_w)$$

$$\boldsymbol{\Sigma}^R = K(K + \sigma^2 J_w)^{-1} \sigma^2 J_w$$

$$J_w = \text{diag}(w^{-2})$$

$$\mathbf{m}_w = \mathbf{m} + \sigma^2 \nabla_y \log(w^2)$$

Measuring outlier-robustness

- The posterior influence function measures the impact of a single outlier on the posterior:

$$\text{PIF}(y_m^c, D) = \text{KL} (p(f | D), p(f | D_m^c))$$

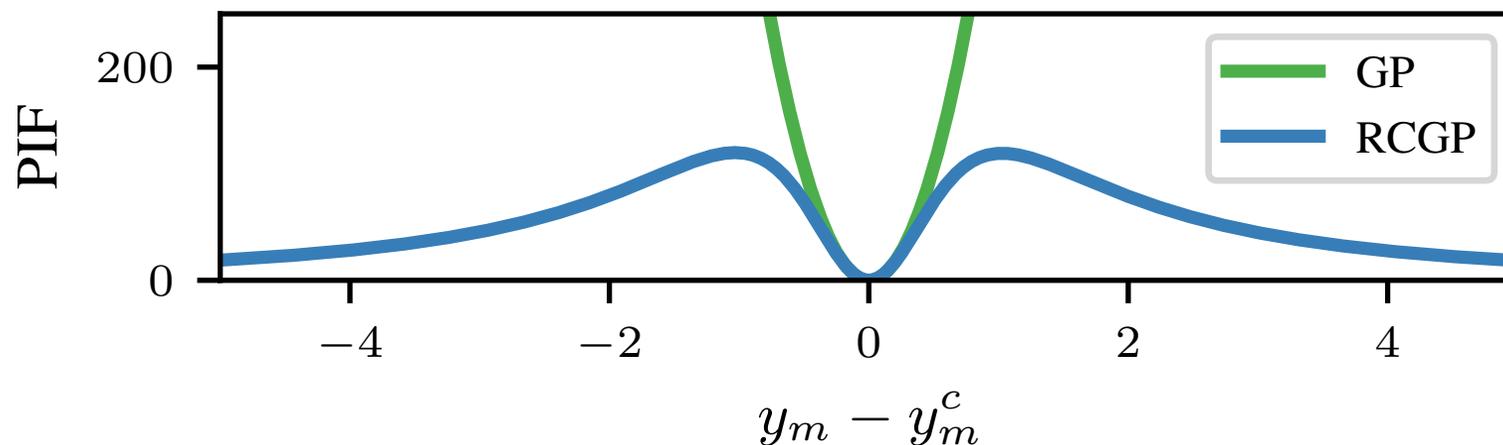
$$D = \{x_i, y_i\}_{i=1}^n$$

$$D_m^c = (D \setminus \{x_m, y_m\}) \cup \{x_m, y_m^c\}$$

RCGPs are provably outlier-robust

- **Theorem (informal):** Suppose $w(x, y) = (1 + (y - m(x))^2/c^2)^{-\frac{1}{2}}$ for some $c > 0$, then RCGPs are robust since:

$$\sup_{y_m^c} \text{PIF}_{\text{RCGP}}(y_m^c, D) < \infty$$



Hyperparameter selection

- The standard approach for selecting hyper parameters is to do empirical Bayes and **maximise the marginal likelihood**.
- This of course does not make sense when the likelihood is wrong!
- Our alternative is to do **leave-one-out cross-validation**

$$\hat{\sigma}^2, \hat{\theta} = \arg \max_{\sigma^2, \theta} \left\{ \sum_{i=1}^n \log p^w(y_i | \mathbf{x}, \mathbf{y}_{-i}, \theta, \sigma^2) \right\},$$

- This can be done efficiently through clever linear algebra tricks and gradient-based optimisation.

Performance when well-specified

	GP	RCGP	t-GP	m-GP
		No Outliers		
Synthetic	0.08 (0.00)	0.08 (0.00)	0.09 (0.00)	0.33 (0.0)
Boston	0.20 (0.01)	0.20 (0.00)	0.20 (0.00)	0.28 (0.0)
Energy	0.02 (0.00)	0.02 (0.00)	0.03 (0.00)	0.61 (0.0)
Yacht	0.01 (0.00)	0.02 (0.00)	0.02 (0.00)	0.33 (0.0)

GPs and RCGPs are comparable when the model is well-specified!

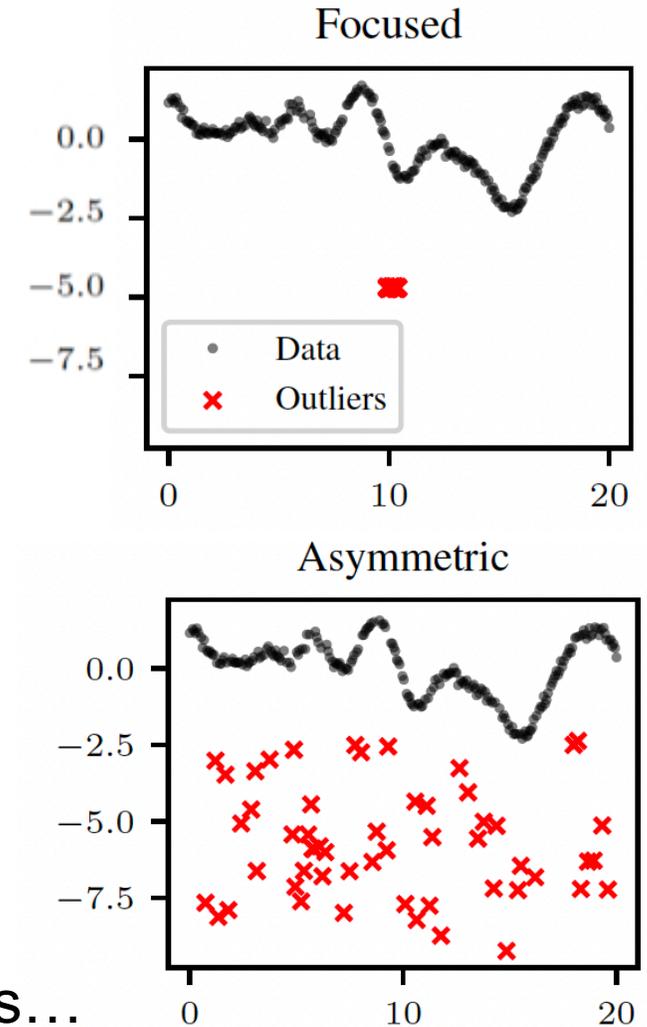
This is not true for other robust methods based on heavy-tailed likelihoods...

Performance when misspecified

	GP	RCGP	t-GP	m-GP
Focused Outliers				
Synthetic	0.19 (0.00)	0.16 (0.00)	0.20 (0.00)	0.23 (0.0)
Boston	0.27 (0.12)	0.22 (0.03)	0.25 (0.01)	0.27 (0.0)
Energy	0.06 (0.06)	0.02 (0.00)	0.03 (0.00)	0.24 (0.0)
Yacht	0.28 (0.19)	0.10 (0.06)	0.24 (0.08)	0.24 (0.0)
Asymmetric Outliers				
Synthetic	1.14 (0.00)	0.82 (0.00)	1.06 (0.00)	0.61 (0.0)
Boston	0.64 (0.04)	0.49 (0.01)	0.52 (0.00)	0.52 (0.0)
Energy	0.55 (0.05)	0.50 (0.16)	0.44 (0.04)	0.41 (0.0)
Yacht	0.54 (0.06)	0.36 (0.05)	0.41 (0.00)	0.40 (0.0)

RCGPs are robust!

Heavy-tailed likelihoods are not suitable for this type of outliers...



RCGPs are fast!

	GP	RCGP	t-GP	m-GP
Synthetic	1.5 (0.1)	1.2 (0.0)	2.2 (0.0)	3.0 (0.0)
Boston	1.9 (0.5)	5.1 (0.9)	30.7 (6.1)	16.7 (1.7)
Energy	3.8 (0.9)	4.6 (2.0)	34.0 (11)	33.8 (0.3)
Yacht	1.6 (0.3)	2.1 (0.2)	5.6 (0.7)	4.5 (0.4)



RCGPs are much faster than other robust alternatives!

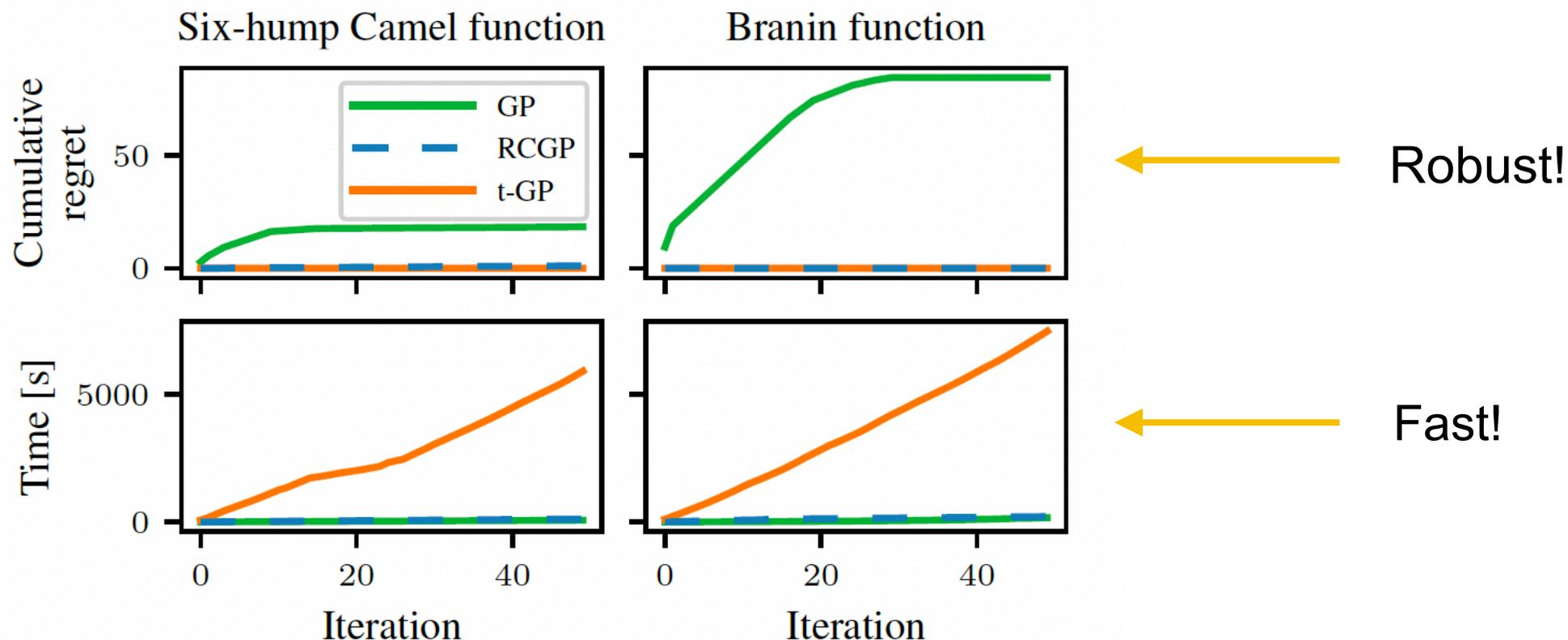
RCGPs are roughly as fast as GPs

	GP	RCGP	t-GP	m-GP
Synthetic	1.5 (0.1)	1.2 (0.0)	2.2 (0.0)	3.0 (0.0)
Boston	1.9 (0.5)	5.1 (0.9)	30.7 (6.1)	16.7 (1.7)
Energy	3.8 (0.9)	4.6 (2.0)	34.0 (11)	33.8 (0.3)
Yacht	1.6 (0.3)	2.1 (0.2)	5.6 (0.7)	4.5 (0.4)

Most of the difference between GP and RCGP comes down to adaptive optimisers for hyper parameter optimisation

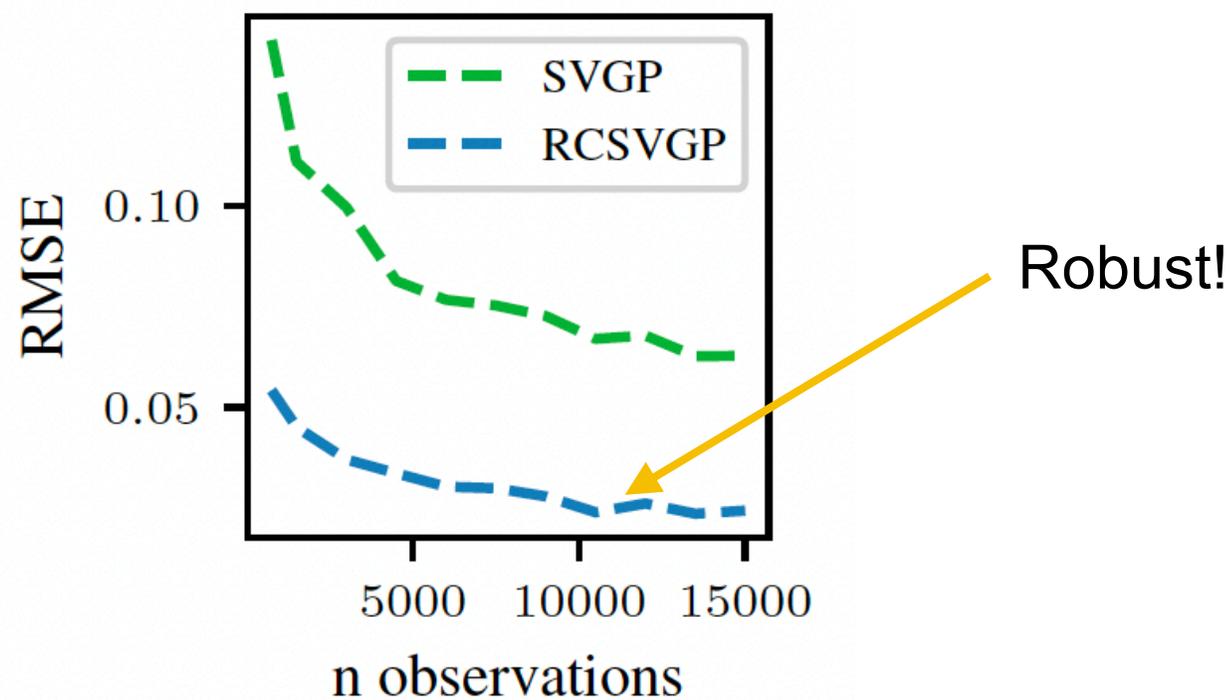
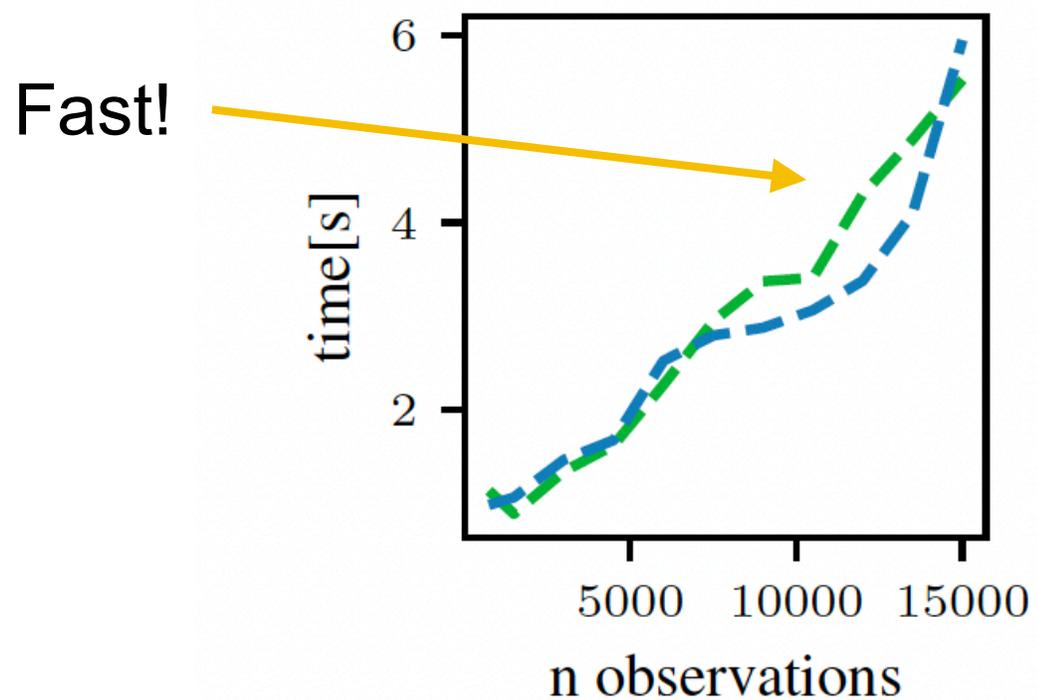
Robust Bayesian Optimisation

- In Bayesian optimisation, the GP posterior is used to create an acquisition function. Our RCGPs naturally lead to robust acquisition functions!



Robust SVGPs

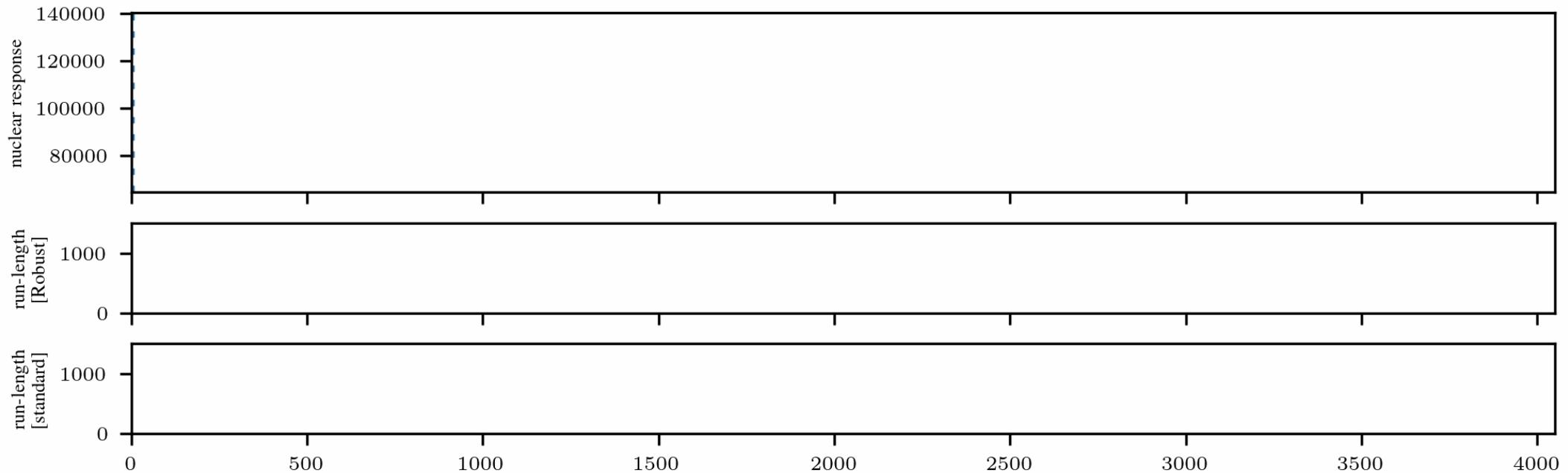
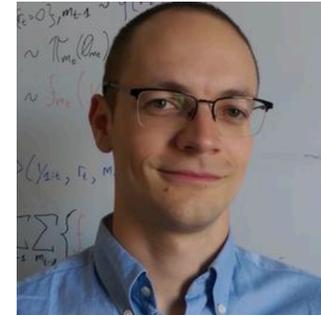
- Sparse Variational GPs (SVGPs) is an approximate GP method which reduces significantly the cost of GPs from $O(n^3)$ to $O(nm^2)$ where m is small. Our approach naturally leads to a robust version!



Conclusion

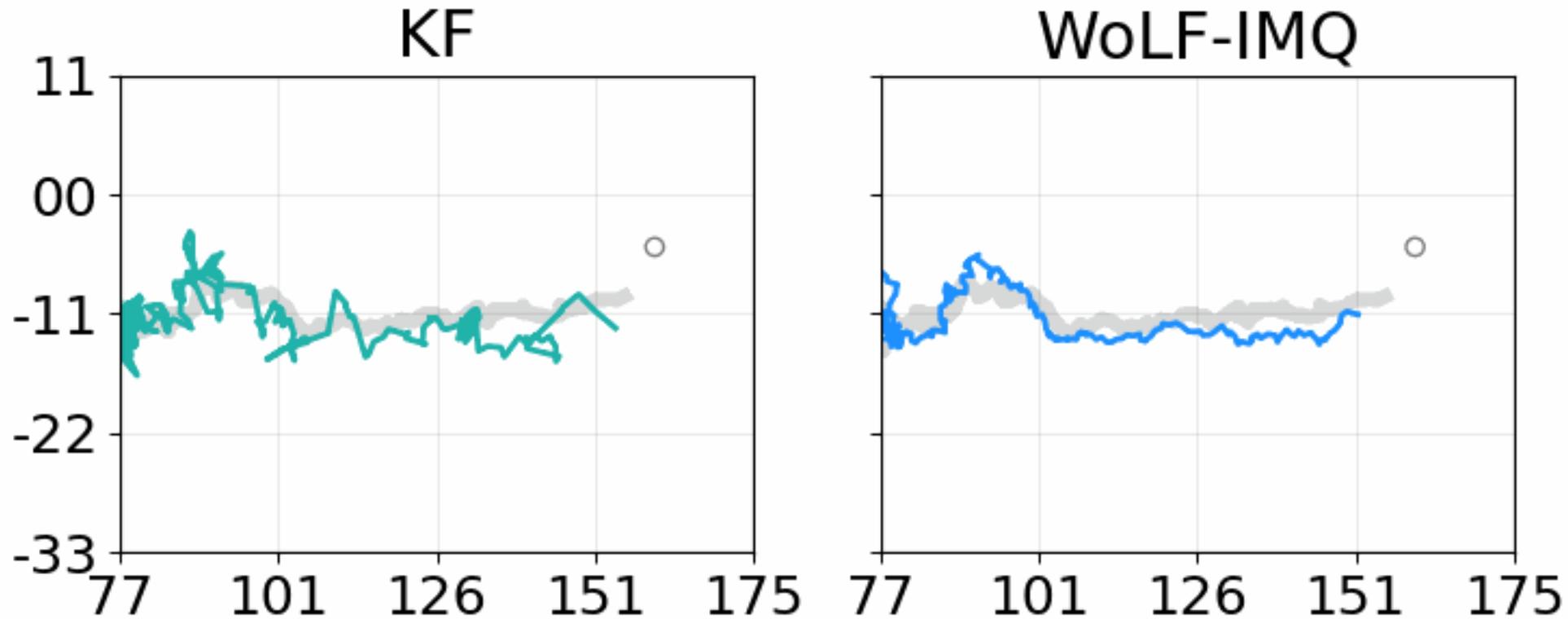
- With careful choices of loss functions, Generalised Bayes can bring both **robustness** and **computational efficiency**!
- RCGPs are an example in the case of GP regression where we get **both robustness and conjugacy**, something no other competitor has managed!
- RCGPs can be developed for any case where standard GPs, and could hence be used for multi-output GPs, multi-fidelity GPs, GPs with derivative or integral information, etc...
- This type of approach is also useful way beyond the GP world....!

Related work (online change point detection)

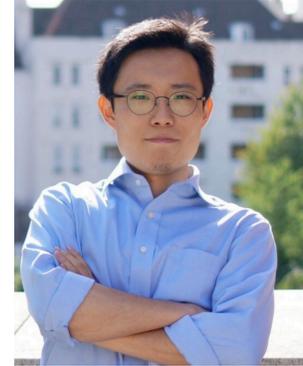


Altamirano, M., Briol, F.-X., & Knoblauch, J. (2023). Robust and scalable Bayesian online changepoint detection. ICML, 642–663.

Related work (Kalman filtering)



Related work (intractable likelihoods)



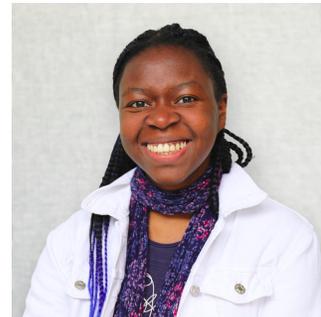
- Robust and conjugate generalised Bayes for **continuous doubly intractable models!**

Matsubara, T., Knoblauch, J., Briol, F.-X., & Oates, C. J. (2022). Robust generalised Bayesian inference for intractable likelihoods. *JRSBB*, 84(3), 997–1022.

- Robust (non-conjugate but fast!) generalised Bayes for **discrete doubly intractable models.**

Matsubara, T., Knoblauch, J., Briol, F.-X., & Oates, C. J. (2023). Generalised Bayesian inference for discrete intractable likelihood. *JASA*, to appear.

More soon.....





Any Questions?

arXiv > stat > arXiv:2311.00463

Statistics > Machine Learning

[Submitted on 1 Nov 2023]

Robust and Conjugate Gaussian Process Regression

Matias Altamirano, François-Xavier Briol, Jeremias Knoblauch