

Scalable Exact Inference in Multi-Output Gaussian Processes

Wessel P. Bruinsma^{1,2}, Eric Perim², Will Tebbutt¹,
J. Scott Hosking^{3,4}, Arno Solin⁵, Richard E. Turner^{1,6}

¹University of Cambridge, ²Invenia Labs, ³British Antarctic Survey,
⁴Alan Turing Institute, ⁵Aalto University, ⁶Microsoft Research

International Conference on Machine Learning 2020

Collaborators



Wessel P.
Bruinsma



Eric Perim



Will Tebbutt



J. Scott
Hoskings



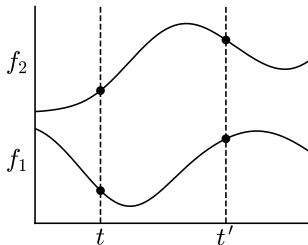
Arno Solin



Richard E.
Turner

Introduction and Motivation

- **Gaussian processes** are a powerful and popular probabilistic modelling framework for nonlinear functions.



Central modelling choice:

$$\mathbf{K}(t, t') = \begin{bmatrix} \text{cov}(f_1(t), f_1(t')) & \text{cov}(f_1(t), f_2(t')) \\ \text{cov}(f_2(t), f_1(t')) & \text{cov}(f_2(t), f_2(t')) \end{bmatrix}$$

- Inference and learning: $O(n^3 p^3)$ time and $O(n^2 p^2)$ memory.
- Often alleviated by **exploiting structure** in \mathbf{K} .

↑
number of
outputs

$$\mathbf{K}(t, t) = \mathbf{I}_m$$

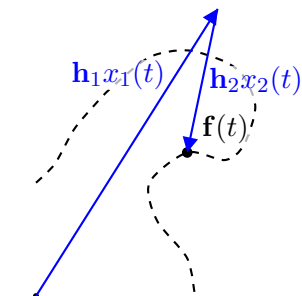
$$\mathbf{x} \sim \mathcal{GP}(\mathbf{0}, \mathbf{K}(t, t')),$$

$$\begin{aligned} \mathbf{f}(t) &= \mathbf{h}_1 x_1(t) + \mathbf{h}_2 x_2(t) \\ &= \mathbf{H}\mathbf{x}(t), \end{aligned}$$

$$\mathbf{y}(t) \sim \mathcal{N}(\mathbf{f}(t), \mathbf{\Sigma}),$$

\mathbf{x} : “latent processes”,

\mathbf{H} : “basis” or “mixing matrix”. $\mathbf{0}$



- Use $m \ll p$ basis vectors: data lives in “pancake” around $\text{col}(\mathbf{H})$.
- Generalisation of FA to time series setting.
- Captures many existing MOGPs from literature.
- Inference and learning: $O(m^3 n^3)$ instead of $O(p^3 n^3)$.

Inside the ILMM

high-dim. observation

$$p \left\{ \begin{bmatrix} \bullet \\ \bullet \\ \vdots \\ \bullet \\ \bullet \end{bmatrix} \right.$$

 \mathbf{y} ✗ inference in $p(\mathbf{y})$ noise: Σ “projected observation”
for $\mathbf{x} \sim \mathcal{GP}(\mathbf{0}, \mathbf{K}(t, t'))$ \mapsto

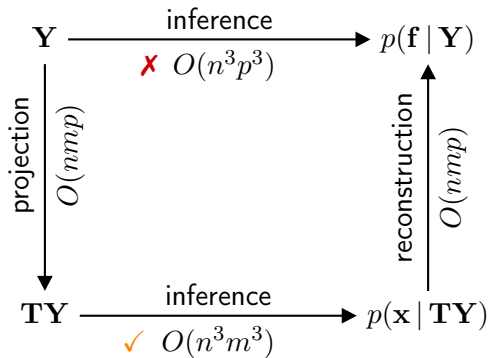
$$\left. \begin{bmatrix} \bullet \\ \vdots \\ \bullet \end{bmatrix} \right\} m (\ll p)$$

$$\mathbf{y}_{\text{proj}} = \mathbf{T}\mathbf{y}$$

✓ inference in $p(\mathbf{x})$ projected noise: $\Sigma_{\mathbf{T}}$ **Proposition:** This is exact!

Key Result (2)

4/17



likelihood of **projected observations** under **projected noise**

$$\log p(\mathbf{Y}) = \log \int p(\mathbf{x}) \prod_{i=1}^n \mathcal{N}(\mathbf{T}\mathbf{y}_i \mid \mathbf{x}_i, \Sigma_{\mathbf{T}}) d\mathbf{x}$$

$$- \underbrace{\frac{1}{2} \sum_{i=1}^n \|\mathbf{y}_i - \mathbf{H}\mathbf{T}\mathbf{y}_i\|_{\Sigma}^2}_{\text{data "lost" by projection (reconstruction error)}} - \underbrace{\frac{1}{2} n \log \frac{|\Sigma|}{|\Sigma_{\mathbf{T}}|}}_{\text{noise "lost" by projection}} + \text{const.}$$

- Learning $\mathbf{H} \Leftrightarrow$ learning $\mathbf{T} \Leftrightarrow$ learning a transform of the data!
- “Regularisation terms” prevent underfitting.

- Inference in ILMM: condition \mathbf{x} on \mathbf{Y}_{proj} under noise $\Sigma_{\mathbf{T}}$.
- Hence,
 - if \mathbf{x} are independent under the prior and the projected noise $\Sigma_{\mathbf{T}}$ is diagonal,
 - then \mathbf{x} remain independent upon observing data.



Treat latent processes independently:
condition x_i on $(\mathbf{Y}_{\text{proj}})_i$: under noise $(\Sigma_{\mathbf{T}})_{ii}$!

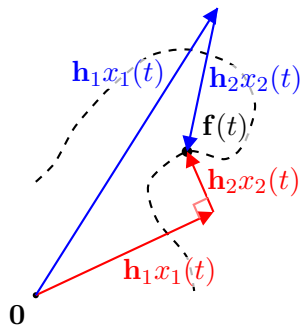
- Decouples inference into independent single-output problems.

“Decoupling” the ILMM

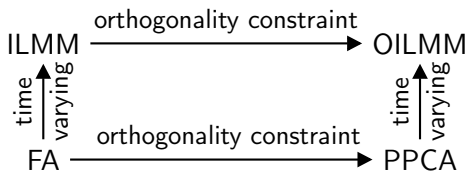
Orthogonal ILMM (OILMM)

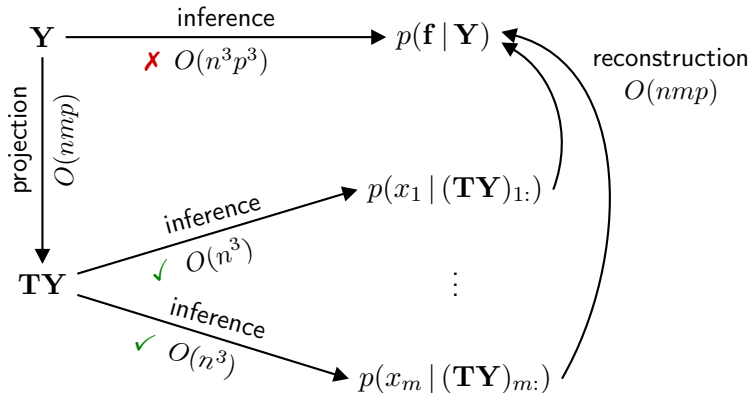
7/17

$$\begin{aligned} \mathbf{x} &\sim \mathcal{GP}(\mathbf{0}, \mathbf{K}(t, t')), \\ \mathbf{f}(t) &= \mathbf{H}\mathbf{x}(t) \\ &= \mathbf{U}\mathbf{S}^{\frac{1}{2}}\mathbf{x}(t), \\ &\begin{array}{cc} \nearrow & \nwarrow \\ \text{orthogonal} & \text{diagonal scaling} \end{array} \\ \mathbf{y}(t) &\sim \mathcal{N}(\mathbf{f}(t), \mathbf{\Sigma}). \end{aligned}$$



Key property: $\mathbf{\Sigma}_T$ is diagonal!

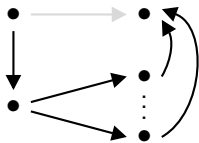




- **Linear** scaling in m !
- Trivially compatible with single-output scaling techniques!

- 1 Project data and compute proj. noise:

$$\mathbf{Y}_{\text{proj}} = \mathbf{S}^{-\frac{1}{2}} \mathbf{U}^T \mathbf{Y}, \quad \Sigma_{\mathbf{T}} = \sigma^{-2} \mathbf{S}^{-1} + \mathbf{D}.$$



- 2 For $i = 1, \dots, m$,


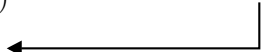
compute the log-probability LML_i of $(\mathbf{Y}_{\text{proj}})_{:i}$ under latent process x_i and observation noise $(\Sigma_{\mathbf{T}})_{ii}$.

- 3 Compute the “regularisation term”:

$$\text{reg.} = -\frac{n}{2} \log |\mathbf{S}| - \frac{n(p-m)}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \|(\mathbf{I}_p - \mathbf{U}\mathbf{U}^T)\mathbf{Y}\|_F^2$$

- 4 Construct the log-probability of the data \mathbf{Y} under the OILMM:

$$\log p(\mathbf{Y}) = \sum_{i=1}^m \text{LML}_i + \text{reg.}$$

	Class	Complexity	
more restrictive 	MOGP	$O(p^3 n^3)$	Use single-output scaling techniques to also bring down complexity in n .
	ILMM	$O(m^3 n^3)$	
	OILMM	$O(mn^3)$	
		$O(mnr^2)$	
$O(mnd^3)$		(d -dim. state-space approximation)	

Orthogonality gives excellent computational benefits.
But how restrictive is it?

Definition

An (O)ILMM is **separable** if $\mathbf{K}(t, t') = k(t, t')\mathbf{I}_m$. Example: ICM.

ILMM versus OILMM:

- Separable case: without loss of generality.
- Non-separable case: only affects **correlations through time**.
- ILMM can be approximated by an OILMM (in KL) if the right singular vectors of \mathbf{H} are close to unit vectors (in $\|\cdot\|_F$).

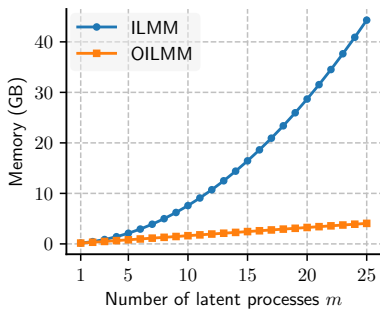
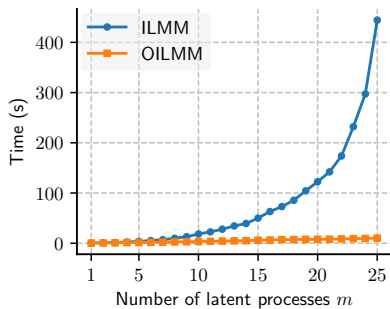
- Separable spatio-temporal GP is an OILMM.
- OILMM gives **non-separable relaxation** of separable models whilst retaining efficient inference.

- Missing data is troublesome: it breaks orthogonality of \mathbf{H} .
- In the paper, we derive a simple and effective approximation.

The OILMM in Practice

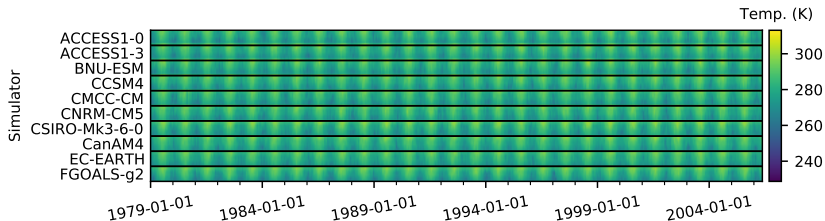
Demonstration of Scalability

13/17



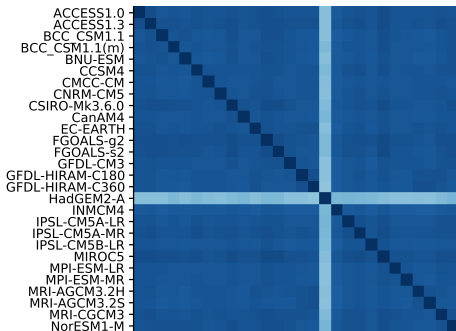
	EEG		FX	
	PPLP	SMSE	PPLP	SMSE
ILMM	-2.11	0.49	3.39	0.19
OILMM	-2.11	0.49	3.39	0.19

- Near identical performance on two real-world data sets.
- Demonstrates that missing data approximation works well.

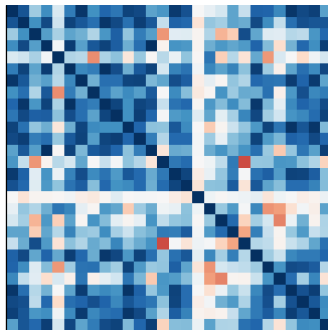


- Jointly model $p_s = 28$ climate simulators at $p_r = 247$ spatial locations and $n = 10\,000$ points in time.
- Equals $p = p_s p_r \approx 7\text{ k}$ outputs and $pn \approx 70\text{ M}$ observations.
- **Goal:** Learn covariance between simulators with $\mathbf{H} = \mathbf{H}_s \otimes \mathbf{H}_r$.
- Use $m = 50$ and inducing points to scale decoupled problems.

Empirical correlations



Learned by OILMM



Use **projection of the data** to accelerate inference in MOGPs with **orthogonal** bases:

- ✓ Linear scaling in m .
- ✓ Simple to implement.
- ✓ Trivially compatible with single-output scaling techniques.
- ✓ Does not sacrifice significant expressivity.