

State Space Expectation Propagation

Efficient Inference Schemes for Temporal Gaussian Processes

William Wilkinson*, Paul Chang*, Michael Riis Andersen[†],
Arno Solin*

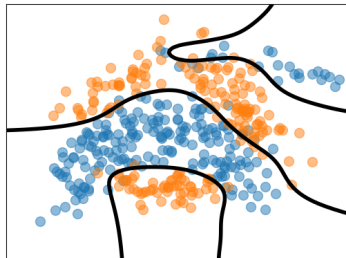
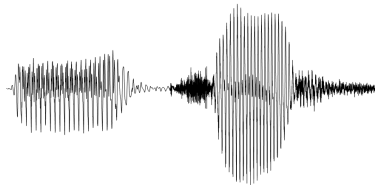
Aalto University*, Technical University of Denmark[†]

ICML 2020



Motivation

- We're interested in long **temporal and spatio-temporal** data with interesting **non-conjugate** GP models (e.g. classification, log-Gaussian Cox processes).
- **Idea: We should treat the temporal dimension in a fundamentally different manner to other dimensions.**



Approximate Inference in Temporal GPs

There exists a **dual kernel / SDE form** for most popular Gaussian process (GP) models

$$f(t) \sim \mathcal{GP}(0, K_\theta(t, t')),$$
$$y_k \sim p(y_k | f(t_k))$$

$$\mathbf{f}_k = \mathbf{A}_{\theta, k} \mathbf{f}_{k-1} + \mathbf{q}_k, \quad \mathbf{q}_k \sim \mathbf{N}(0, \mathbf{Q}_k)$$
$$y_k = h(\mathbf{f}_k, \sigma_k), \quad \sigma_k \sim \mathbf{N}(0, \Sigma_k)$$

Approximate Inference in Temporal GPs

There exists a **dual kernel / SDE form** for most popular Gaussian process (GP) models

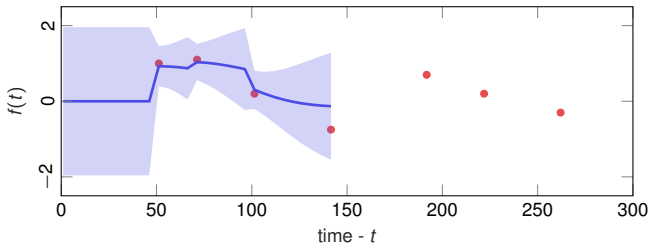
$$\begin{aligned} f(t) &\sim \mathcal{GP}(0, K_\theta(t, t')), & \mathbf{f}_k &= \mathbf{A}_{\theta, k} \mathbf{f}_{k-1} + \mathbf{q}_k, & \mathbf{q}_k &\sim \mathbf{N}(0, \mathbf{Q}_k) \\ y_k &\sim p(y_k | f(t_k)) & y_k &= h(\mathbf{f}_k, \sigma_k), & \sigma_k &\sim \mathbf{N}(0, \Sigma_k) \end{aligned}$$

inference in $\mathcal{O}(n)$ via **Kalman filtering and smoothing**

Approximate Inference

Kalman filter update step:

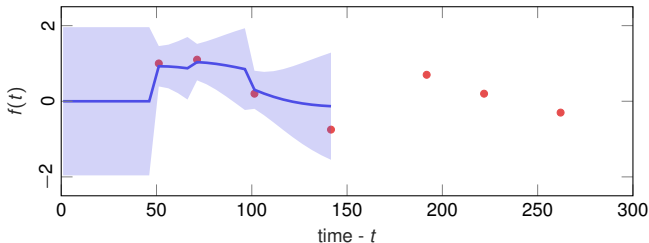
$$p(\mathbf{f}_k | y_{1:k}) \propto N(\mathbf{m}_k^{\text{predict}}, \mathbf{P}_k^{\text{predict}}) p(y_k | f(t_k))$$



Approximate Inference

Kalman filter update step:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:k}) &\propto N(\mathbf{m}_k^{\text{predict}}, \mathbf{P}_k^{\text{predict}}) p(y_k | f(t_k)) \\ &\approx N(\mathbf{m}_k^{\text{predict}}, \mathbf{P}_k^{\text{predict}}) \underbrace{N(\mathbf{m}_k^{\text{site}}, \mathbf{P}_k^{\text{site}})}_{\text{"site"}} \end{aligned}$$

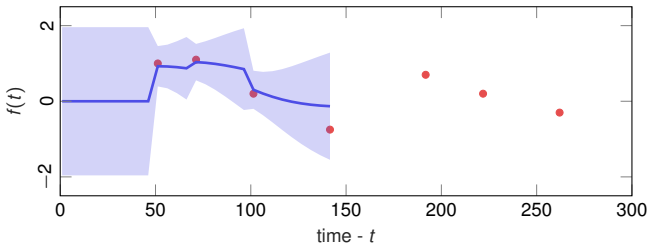


Approximate Inference

Kalman filter update step:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:k}) &\propto N(\mathbf{m}_k^{\text{predict}}, \mathbf{P}_k^{\text{predict}}) p(y_k | f(t_k)) \\ &\approx N(\mathbf{m}_k^{\text{predict}}, \mathbf{P}_k^{\text{predict}}) N(\mathbf{m}_k^{\text{site}}, \mathbf{P}_k^{\text{site}}) \end{aligned}$$

Approx. Inference:
← select parameters

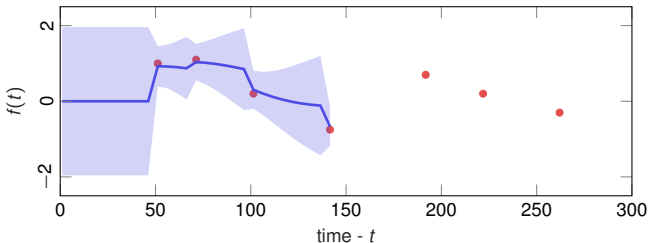


Approximate Inference

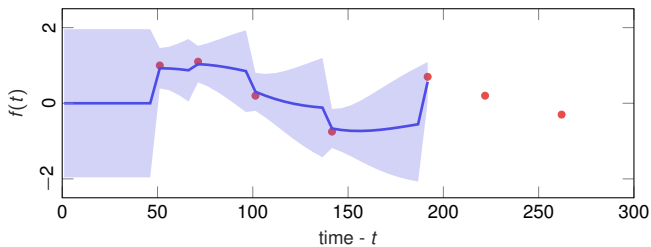
Kalman filter update step:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:k}) &\propto N(\mathbf{m}_k^{\text{predict}}, \mathbf{P}_k^{\text{predict}}) p(y_k | f(t_k)) \\ &\approx N(\mathbf{m}_k^{\text{predict}}, \mathbf{P}_k^{\text{predict}}) N(\mathbf{m}_k^{\text{site}}, \mathbf{P}_k^{\text{site}}) \end{aligned}$$

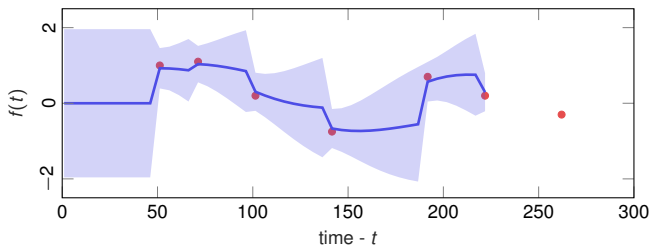
Approx. Inference:
← select parameters



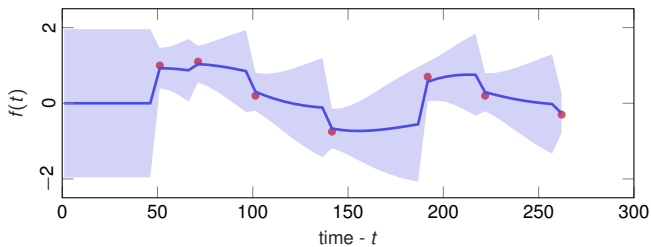
Approximate Inference



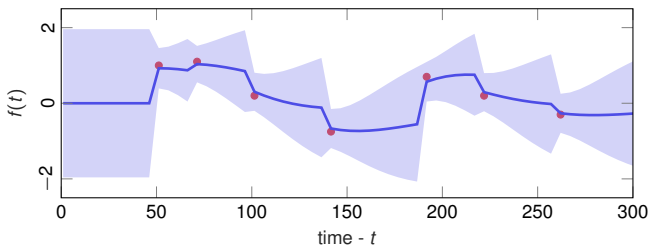
Approximate Inference



Approximate Inference



Approximate Inference

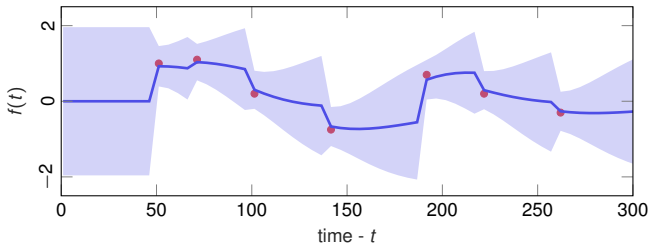


Approximate Inference

Smoothing:

- update posterior with future observations,

$$p(\mathbf{f}_k | y_{1:N}) = N(\mathbf{m}_k^{\text{post.}}, \mathbf{P}_k^{\text{post.}})$$

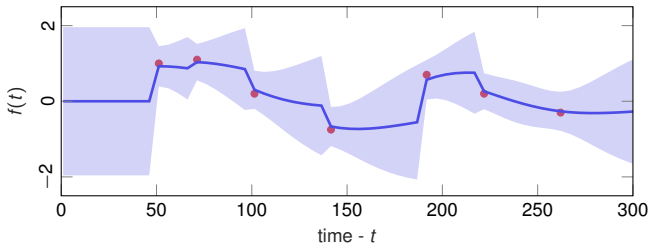


Approximate Inference

Smoothing:

- update posterior with future observations,

$$p(\mathbf{f}_k | y_{1:N}) = N(\mathbf{m}_k^{\text{post.}}, \mathbf{P}_k^{\text{post.}})$$

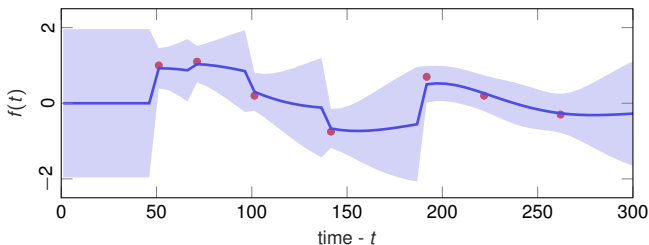


Approximate Inference

Smoothing:

- update posterior with future observations,

$$p(\mathbf{f}_k | y_{1:N}) = N(\mathbf{m}_k^{\text{post.}}, \mathbf{P}_k^{\text{post.}})$$

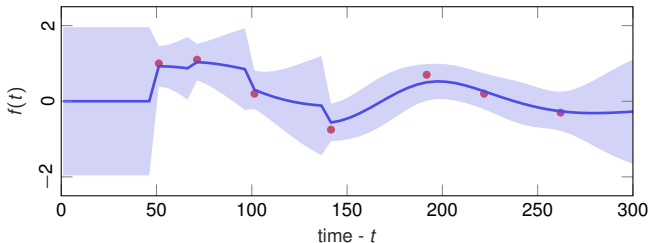


Approximate Inference

Smoothing:

- update posterior with future observations,

$$p(\mathbf{f}_k | y_{1:N}) = N(\mathbf{m}_k^{\text{post.}}, \mathbf{P}_k^{\text{post.}})$$

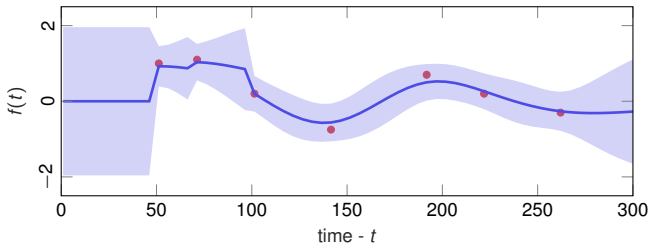


Approximate Inference

Smoothing:

- update posterior with future observations,

$$p(\mathbf{f}_k | y_{1:N}) = N(\mathbf{m}_k^{\text{post.}}, \mathbf{P}_k^{\text{post.}})$$

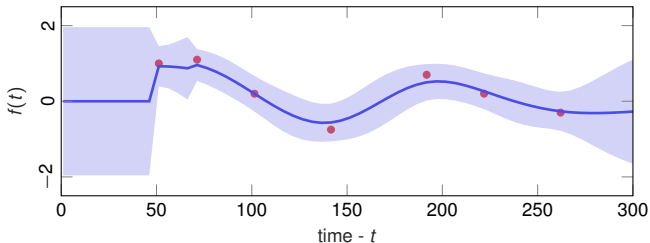


Approximate Inference

Smoothing:

- update posterior with future observations,

$$p(\mathbf{f}_k | y_{1:N}) = N(\mathbf{m}_k^{\text{post.}}, \mathbf{P}_k^{\text{post.}})$$

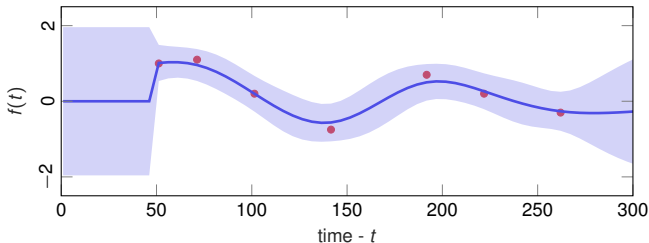


Approximate Inference

Smoothing:

- update posterior with future observations,

$$p(\mathbf{f}_k | y_{1:N}) = N(\mathbf{m}_k^{\text{post.}}, \mathbf{P}_k^{\text{post.}})$$

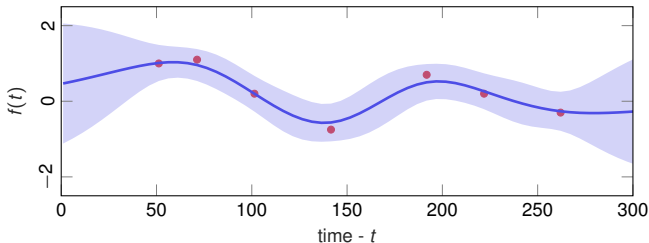


Approximate Inference

Smoothing:

- update posterior with future observations,

$$p(\mathbf{f}_k | y_{1:N}) = N(\mathbf{m}_k^{\text{post.}}, \mathbf{P}_k^{\text{post.}})$$



Approximate Inference

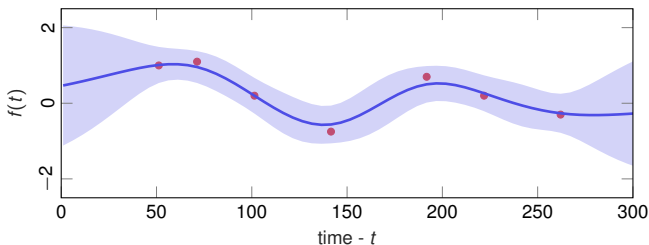
Smoothing:

- update posterior with future observations,

$$p(\mathbf{f}_k | y_{1:N}) = N(\mathbf{m}_k^{\text{post.}}, \mathbf{P}_k^{\text{post.}})$$

Our Contribution:

Given marginal posterior $N(\mathbf{m}_k^{\text{post.}}, \mathbf{P}_k^{\text{post.}})$, we show how approximate inference amounts to a **simple site parameter update** rule during smoothing.



Approximate Inference

Smoothing:

- update posterior with future observations,
 $p(\mathbf{f}_k | y_{1:N}) = N(\mathbf{m}_k^{\text{post.}}, \mathbf{P}_k^{\text{post.}})$

Our Contribution:

Given marginal posterior $N(\mathbf{m}_k^{\text{post.}}, \mathbf{P}_k^{\text{post.}})$, we show how approximate inference amounts to a **simple site parameter update** rule during smoothing.

This encompasses:

- Power Expectation Propagation
- Variational Inference (with natural gradients)
- Extended Kalman Smoothing
- Unscented / Gauss-Hermite Kalman Smoothing
- Posterior Linearisation

Parameter Update Rules

$$\text{for } \nabla \mathcal{L}_k = \frac{d\mathcal{L}_k}{d\mathbf{m}_k}$$

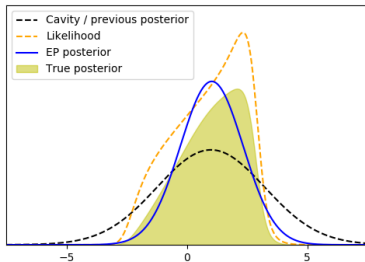
Power Expectation Propagation:

$$q_{\text{cavity}}(\mathbf{f}_k) = q_{\text{post.}}(\mathbf{f}_k) / q_{\text{site}}^\alpha(\mathbf{f}_k)$$

$$\mathcal{L}_k = \log \mathbb{E}_{q_{\text{cavity}}} [\rho^\alpha(\mathbf{y}_k | \mathbf{f}_k)]$$

$$\mathbf{P}_k^{\text{site}} = -\alpha \left(\mathbf{P}_k^{\text{cavity}} + (\nabla^2 \mathcal{L}_k)^{-1} \right)$$

$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^{\text{cavity}} - (\nabla^2 \mathcal{L}_k)^{-1} \nabla \mathcal{L}_k$$



Parameter Update Rules

$$\text{for } \nabla \mathcal{L}_k = \frac{d\mathcal{L}_k}{d\mathbf{m}_k}$$

Power Expectation Propagation:

$$q_{\text{cavity}}(\mathbf{f}_k) = q_{\text{post.}}(\mathbf{f}_k) / q_{\text{site}}^\alpha(\mathbf{f}_k)$$

$$\mathcal{L}_k = \log \mathbb{E}_{q_{\text{cavity}}} [\rho^\alpha(\mathbf{y}_k | \mathbf{f}_k)]$$

$$\mathbf{P}_k^{\text{site}} = -\alpha \left(\mathbf{P}_k^{\text{cavity}} + (\nabla^2 \mathcal{L}_k)^{-1} \right)$$

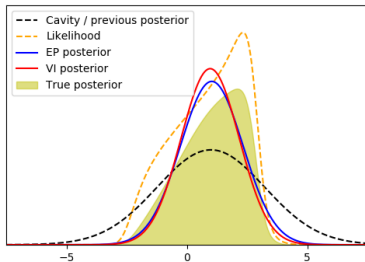
$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^{\text{cavity}} - (\nabla^2 \mathcal{L}_k)^{-1} \nabla \mathcal{L}_k$$

Variational Inference:

$$\mathcal{L}_k = \mathbb{E}_{q_{\text{post.}}} [\log p(\mathbf{y}_k | \mathbf{f}_k)]$$

$$\mathbf{P}_k^{\text{site}} = -(\nabla^2 \mathcal{L}_k)^{-1}$$

$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^{\text{post.}} - (\nabla^2 \mathcal{L}_k)^{-1} \nabla \mathcal{L}_k$$



Parameter Update Rules

$$\text{for } \nabla \mathcal{L}_k = \frac{d\mathcal{L}_k}{d\mathbf{m}_k}$$

Power Expectation Propagation:

$$q_{\text{cavity}}(\mathbf{f}_k) = q_{\text{post.}}(\mathbf{f}_k) / q_{\text{site}}^\alpha(\mathbf{f}_k)$$

$$\mathcal{L}_k = \log \mathbb{E}_{q_{\text{cavity}}} [\rho^\alpha(\mathbf{y}_k | \mathbf{f}_k)]$$

$$\mathbf{P}_k^{\text{site}} = -\alpha \left(\mathbf{P}_k^{\text{cavity}} + (\nabla^2 \mathcal{L}_k)^{-1} \right)$$

$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^{\text{cavity}} - (\nabla^2 \mathcal{L}_k)^{-1} \nabla \mathcal{L}_k$$

Variational Inference:

$$\mathcal{L}_k = \mathbb{E}_{q_{\text{post.}}} [\log p(\mathbf{y}_k | \mathbf{f}_k)]$$

$$\mathbf{P}_k^{\text{site}} = -(\nabla^2 \mathcal{L}_k)^{-1}$$

$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^{\text{post.}} - (\nabla^2 \mathcal{L}_k)^{-1} \nabla \mathcal{L}_k$$

Extended Kalman Smoother:

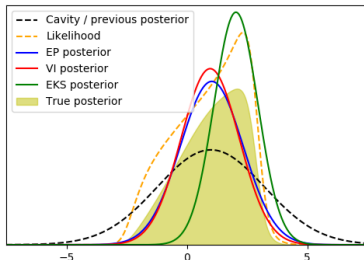
$$\mathbf{v}_k = \mathbf{y}_k - \mathbf{h}(\mathbf{m}_k^{\text{post.}}, \mathbf{0})$$

$$\mathbf{S}_k = \mathbf{H}_f^\top \mathbf{P}_k^{\text{post.}} \mathbf{H}_f + \mathbf{H}_\sigma \Sigma_k \mathbf{H}_\sigma^\top$$

$$\mathbf{P}_k^{\text{site}} = \left(\mathbf{H}_f^\top \left(\mathbf{H}_\sigma \Sigma_k \mathbf{H}_\sigma^\top \right)^{-1} \mathbf{H}_f \right)^{-1}$$

$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^{\text{post.}} + (\mathbf{P}_k^{\text{site}} + \mathbf{P}_k^{\text{post.}}) \mathbf{H}_f^\top \mathbf{S}_k^{-1} \mathbf{v}_k$$

$$\text{for } \mathbf{H}_f = \frac{d\mathbf{h}}{d\mathbf{f}} \text{ and } \mathbf{H}_\sigma = \frac{d\mathbf{h}}{d\sigma}, \sigma_k \sim N(0, \Sigma_k)$$



A Unifying Perspective

- For sequential data, the EKF / UKF / GHKF are equivalent to single-sweep EP where the moment matching is solved via linearisation.

A Unifying Perspective

- For sequential data, the **EKF / UKF / GHKF** are equivalent to **single-sweep EP** where the moment matching is solved via **linearisation**.
- The **iterated Kalman smoothers** (**EKS / UKS / GHKS**) can also be recovered under certain parameter choices. But note that they optimise a different objective to EP (see paper for details).

A Unifying Perspective

- For sequential data, the **EKF / UKF / GHKF** are equivalent to **single-sweep EP** where the moment matching is solved via **linearisation**.
- The **iterated Kalman smoothers** (EKS / UKS / GHKS) can also be recovered under certain parameter choices. But note that they optimise a different objective to EP (see paper for details).
- We show how **natural gradient VI** updates are surprisingly similar to the EP updates (when using a similar parametrisation).

New Algorithms

- We propose to mix the beneficial properties of EP with the efficiency of classical smoothers.

New Algorithms

- We propose to mix the beneficial properties of EP with the efficiency of classical smoothers.
- For example, using [linearisation to speed up the updates](#), whilst also introducing the [EP cavity and fractional updates](#).

New Algorithms

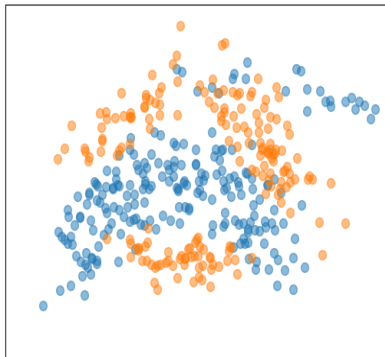
- We propose to mix the beneficial properties of EP with the efficiency of classical smoothers.
- For example, using [linearisation to speed up the updates](#), whilst also introducing the [EP cavity and fractional updates](#).
- We call this [Extended Kalman Expectation Propagation \(EK-EP\)](#).

New Algorithms

- We propose to mix the beneficial properties of EP with the efficiency of classical smoothers.
- For example, using [linearisation to speed up the updates](#), whilst also introducing the [EP cavity and fractional updates](#).
- We call this [Extended Kalman Expectation Propagation \(EK-EP\)](#).
- It has clear computational benefits when the parameter updates are high-dimensional, *e.g.*, in [spatio-temporal problems](#).

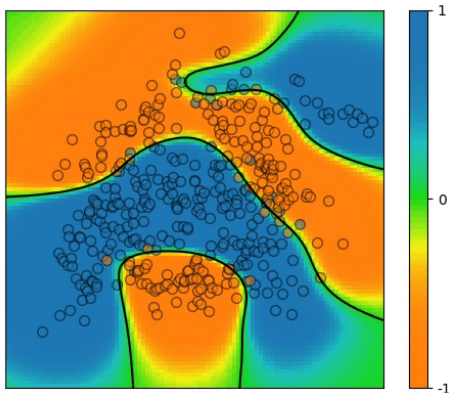
Spatio-Temporal Classification

- We show that our smoothing methods can be applied to tasks with **more than one input dimension**.



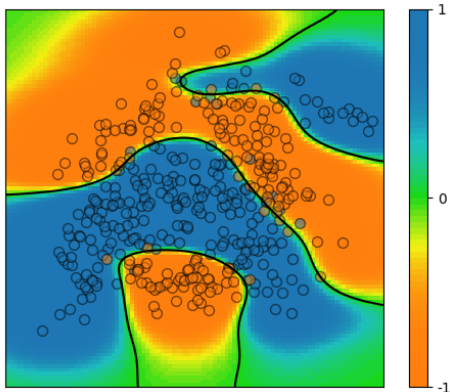
Spatio-Temporal Classification

- We show that our smoothing methods can be applied to tasks with more than one input dimension.



Spatio-Temporal Classification

- We show that our smoothing methods can be applied to tasks with **more than one input dimension**.



We treat the first dimension (x-axis) as time, and run iterated spatio-temporal smoothing (this demo uses EP).

Fast Learning Using JAX

- Temporal GP methods have been limited by a lack of appropriate software for hyperparameter learning:

Automatic differentiation + massive for loops

`https://github.com/AaltoML/kalman-jax`

Fast Learning Using JAX

- Temporal GP methods have been limited by a lack of appropriate software for hyperparameter learning:

Automatic differentiation + massive for loops

- We provide a temporal GP framework in **JAX** with all inference methods implemented.

`https://github.com/AaltoML/kalman-jax`

Fast Learning Using JAX

- Temporal GP methods have been limited by a lack of appropriate software for hyperparameter learning:

Automatic differentiation + massive for loops

- We provide a temporal GP framework in **JAX** with all inference methods implemented.
 - i)* avoid loop “unrolling” to reduce compilation overheads

`https://github.com/AaltoML/kalman-jax`

Fast Learning Using JAX

- Temporal GP methods have been limited by a lack of appropriate software for hyperparameter learning:

Automatic differentiation + massive for loops

- We provide a temporal GP framework in **JAX** with all inference methods implemented.
 - i)* avoid loop “unrolling” to reduce compilation overheads
 - ii)* JIT compilation to avoid graph retracing

`https://github.com/AaltoML/kalman-jax`

Fast Learning Using JAX

- Temporal GP methods have been limited by a lack of appropriate software for hyperparameter learning:

Automatic differentiation + massive for loops

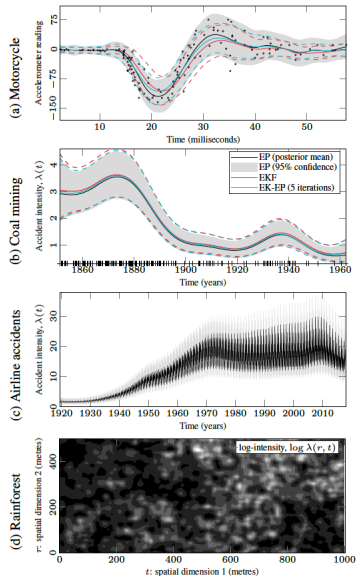
- We provide a temporal GP framework in **JAX** with all inference methods implemented.
 - i)* avoid loop “unrolling” to reduce compilation overheads
 - ii)* JIT compilation to avoid graph retracing
 - iii)* Exploits accelerated linear algebra (XLA) ops

<https://github.com/AaltoML/kalman-jax>

Results

We run extensive analysis on synthetic and real world data:

- Heteroscedastic Noise
- 1D & 2D Log Gaussian Cox Process
- 1D & 2D Classification
- Audio Amplitude Demodulation



Results

- No consistently best inference method or EP power value:
 - EK-EP the only practical method when updates are high dimensional (rainforest)
 - EP or VI needed when likelihood is highly nonlinear

Results

- No consistently best inference method or EP power value:
 - EK-EP the only practical method when updates are high dimensional (rainforest)
 - EP or VI needed when likelihood is highly nonlinear
- We compare against non-sequential baselines (SVGP and EP).
 - Sequential learning methods match the performance of batch methods, whilst scaling to larger data.

Results

- No consistently best inference method or EP power value:
 - EK-EP the only practical method when updates are high dimensional (rainforest)
 - EP or VI needed when likelihood is highly nonlinear
- We compare against non-sequential baselines (SVGP and EP).
 - Sequential learning methods match the performance of batch methods, whilst scaling to larger data.
- See the paper for full results table.

Thanks for Listening

Take home messages:

- Any approximate inference method can be framed as a simple parameter update rule during Kalman smoothing.
- Sequential methods match the performance of batch methods, and can be extended to multiple dimensions.
- We provide fast JAX code for all methods.

Contact: `william.wilkinson@aalto.fi`

JAX code: `https://github.com/AaltoML/kalman-jax`