

Models and inference for temporal Gaussian processes

(*i.e.*, GPs for signal processing)

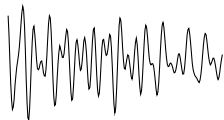
William Wilkinson

Postdoctoral researcher in machine learning
Department of computer science
Aalto University, Finland

Winter 2019

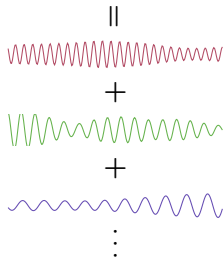
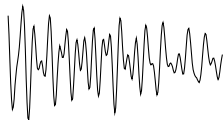
Problem domain

We are interested in applying **Gaussian process** models to time-domain **signals**.



Problem domain

We are interested in applying **Gaussian process** models to time-domain **signals**.



Detailed
latent
structure

Problem domain

We are interested in applying **Gaussian process** models to time-domain **signals**.



Long / unbounded
time duration

A contradiction

GPs for signal processing

Often acquire data at a rate
of many '000s per second

$\mathcal{O}(n^3)$ inference

A contradiction

GPs for signal processing

Often acquire data at a rate
of many '000s per second

$\mathcal{O}(n^3)$ inference

Sparse GPs are not a
natural fit for signals:

- still \sim cubic scaling
- smoothing causes a loss
of perceptually important
information

Solution - State space models

To get around this contradiction, we exploit the **sequential** nature of our data to reformulate our model.

First, some setting up:

First, some setting up:

Think of **signal processing as learning a function of time**, allowing us extrapolate, interpolate, reduce noise, . . .

First, some setting up:

Think of **signal processing as learning a function of time**, allowing us extrapolate, interpolate, reduce noise, . . .

$$f(t) \sim \mathcal{GP}(\mathbf{0}, K_{\theta}(t, t'))$$

First, some setting up:

Think of **signal processing as learning a function of time**, allowing us extrapolate, interpolate, reduce noise, . . .

$$f(t) \sim \mathcal{GP}(\mathbf{0}, K_{\theta}(t, t'))$$

A GP prior states:

- evaluations of $f(\cdot)$ are jointly Gaussian
- covariance between time steps is determined by $K_{\theta}(\cdot, \cdot)$

*“The **state** of a dynamic system is the smallest collection of numbers which must be specified at time t_k in order to be able to predict the behaviour of the system for any time $t_{k+1} \geq t_k$.”*

Rudolf E. Kalman

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

*“The **state** of a dynamic system is the smallest collection of numbers which must be specified at time t_k in order to be able to predict the behaviour of the system for any time $t_{k+1} \geq t_k$.”*

Rudolf E. Kalman

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

For a GP, the “numbers” required might be:

*“The **state** of a dynamic system is the smallest collection of numbers which must be specified at time t_k in order to be able to predict the behaviour of the system for any time $t_{k+1} \geq t_k$.”*

Rudolf E. Kalman

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

For a GP, the “numbers” required might be:

- some time derivatives, $\mathbf{f}_k = (f(t_k), \dot{f}(t_k), \ddot{f}(t_k), \dots)^\top$

$$\mathbf{f}_{k+1} = \mathbf{f}_k$$

*“The **state** of a dynamic system is the smallest collection of numbers which must be specified at time t_k in order to be able to predict the behaviour of the system for any time $t_{k+1} \geq t_k$.”*

Rudolf E. Kalman

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

For a GP, the “numbers” required might be:

- some time derivatives, $\mathbf{f}_k = (f(t_k), \dot{f}(t_k), \ddot{f}(t_k), \dots)^\top$
- a transition model between time steps, $\mathbf{A}_{\theta,k}$

$$\mathbf{f}_{k+1} = \mathbf{A}_{\theta,k} \mathbf{f}_k$$

“The *state* of a dynamic system is the smallest collection of numbers which must be specified at time t_k in order to be able to predict the behaviour of the system for any time $t_{k+1} \geq t_k$.”

Rudolf E. Kalman

$$f(t) \sim \mathcal{GP}(\mathbf{0}, K_\theta(t, t'))$$

For a GP, the “numbers” required might be:

- some time derivatives, $\mathbf{f}_k = (f(t_k), \dot{f}(t_k), \ddot{f}(t_k), \dots)^\top$
- a transition model between time steps, $\mathbf{A}_{\theta,k}$
- a process noise covariance, $\mathbf{Q}_{\theta,k}$

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

*“The **state** of a dynamic system is the smallest collection of numbers which must be specified at time t_k in order to be able to predict the behaviour of the system for any time $t_{k+1} \geq t_k$.”*

Rudolf E. Kalman

$$f(t) \sim \mathcal{GP}(\mathbf{0}, K_\theta(t, t'))$$

For a GP, the “numbers” required might be:

- some time derivatives, $\mathbf{f}_k = (f(t_k), \dot{f}(t_k), \ddot{f}(t_k), \dots)^\top$
- a transition model between time steps, $\mathbf{A}_{\theta,k}$
- a process noise covariance, $\mathbf{Q}_{\theta,k}$
- an observation model, $\mathbf{h}^\top = (1, 0, 0, \dots)$

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

$$f(t_{k+1}) = \mathbf{h}^\top \mathbf{f}_{k+1}$$

Stochastic differential equations

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

This state space model is the **discrete-time solution** to the linear time-invariant (LTI) **stochastic differential equation** (SDE):

$$\frac{d\mathbf{f}(t)}{dt} = \mathbf{F}_{\theta}\mathbf{f}(t) + \mathbf{L}\mathbf{w}(t)$$

Stochastic differential equations

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

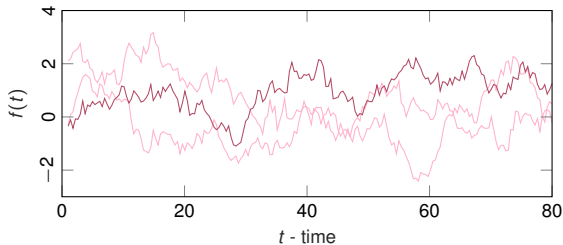
This state space model is the **discrete-time solution** to the linear time-invariant (LTI) **stochastic differential equation** (SDE):

$$\frac{d\mathbf{f}(t)}{dt} = \mathbf{F}_{\theta}\mathbf{f}(t) + \mathbf{L}\mathbf{w}(t)$$

with initial state $\mathbf{f}(t_0) = \mathbf{N}(\mathbf{0}, \mathbf{P}_{\infty})$, for some stationary covariance \mathbf{P}_{∞} .

- $\mathbf{w}(t)$ is white noise with spectral density \mathbf{Q}_c
- \mathbf{F}_{θ} is a *feedback* matrix, $\mathbf{A}_{\theta,k} = \exp(\mathbf{F}_{\theta}(t_{k'} - t_k))$
- \mathbf{L} is a *noise-effect* matrix

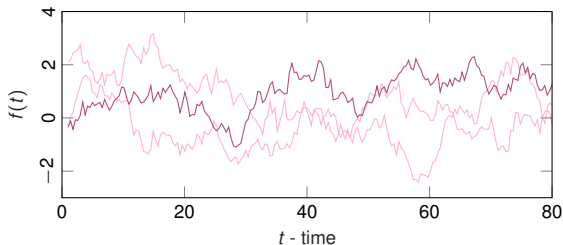
1D example



Samples from the Ornstein-Uhlenbeck process prior with $q = 0.2$, $\lambda = 0.1$.

$$\frac{df(t)}{dt} = -\lambda f(t) + w(t)$$

1D example

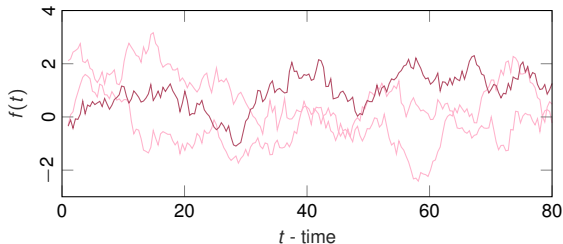


Samples from the Ornstein-Uhlenbeck process prior with $q = 0.2$, $\lambda = 0.1$.

$$\frac{df(t)}{dt} = -\lambda f(t) + w(t)$$

we have set $\mathbf{F}_\theta = -\lambda$, $\mathbf{L} = 1$, $\mathbf{Q}_c = q$, $\mathbf{h}^\top = 1 \implies \mathbf{P}_\infty = q/2\lambda$.

1D example



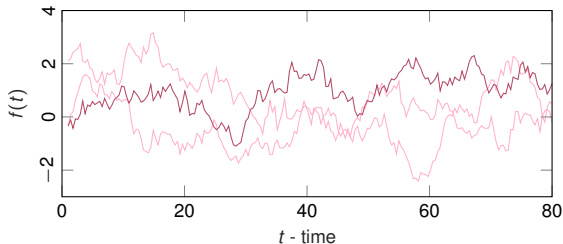
Samples from the Ornstein-Uhlenbeck process prior with $q = 0.2$, $\lambda = 0.1$.

$$\frac{df(t)}{dt} = -\lambda f(t) + w(t)$$

we have set $\mathbf{F}_\theta = -\lambda$, $\mathbf{L} = 1$, $\mathbf{Q}_c = q$, $\mathbf{h}^\top = 1 \implies \mathbf{P}_\infty = q/2\lambda$.

$$K(t, t') = \begin{cases} \mathbf{h}^\top \mathbf{P}_\infty \exp((t' - t)\mathbf{F}_\theta)^\top \mathbf{h}, & \text{if } t' - t \geq 0 \\ \mathbf{h}^\top \exp(-(t' - t)\mathbf{F}_\theta) \mathbf{P}_\infty \mathbf{h}, & \text{if } t' - t < 0 \end{cases}$$

1D example



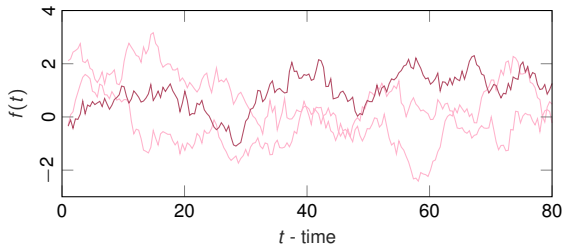
Samples from the Ornstein-Uhlenbeck process prior with $q = 0.2$, $\lambda = 0.1$.

$$\frac{df(t)}{dt} = -\lambda f(t) + w(t)$$

we have set $\mathbf{F}_\theta = -\lambda$, $\mathbf{L} = 1$, $\mathbf{Q}_c = q$, $\mathbf{h}^\top = 1 \implies \mathbf{P}_\infty = q/2\lambda$.

$$K(t, t') = \begin{cases} q/2\lambda \exp(-(t' - t)\lambda), & \text{if } t' - t \geq 0 \\ \exp((t' - t)\lambda) q/2\lambda, & \text{if } t' - t < 0 \end{cases}$$

1D example



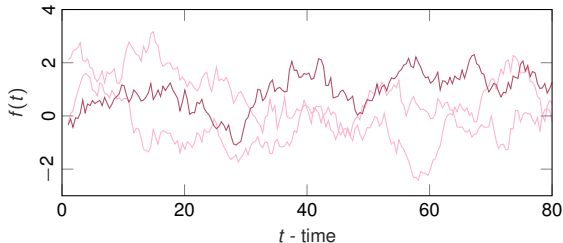
Samples from the Ornstein-Uhlenbeck process prior with $q = 0.2$, $\lambda = 0.1$.

$$\frac{df(t)}{dt} = -\lambda f(t) + w(t)$$

we have set $\mathbf{F}_\theta = -\lambda$, $\mathbf{L} = 1$, $\mathbf{Q}_c = q$, $\mathbf{h}^\top = 1 \implies \mathbf{P}_\infty = q/2\lambda$.

$$K(t, t') = q/2\lambda \exp(-|t' - t|\lambda)$$

1D example



Samples from the Ornstein-Uhlenbeck process prior with $q = 0.2$, $\lambda = 0.1$.

$$\frac{df(t)}{dt} = -\lambda f(t) + w(t)$$

we have set $\mathbf{F}_\theta = -\lambda$, $\mathbf{L} = 1$, $\mathbf{Q}_c = q$, $\mathbf{h}^\top = 1 \implies \mathbf{P}_\infty = q/2\lambda$.

$$K(t, t') = \sigma^2 \exp(-|t' - t|/\ell), \quad \sigma^2 = q/2\lambda, \quad \ell = 1/\lambda$$

There exists a **dual kernel / SDE form** for most popular GP models

$$f(t) \sim \mathcal{GP}(\mathbf{0}, K_{\theta}(t, t')),$$

$$y_k \sim \mathcal{N}(f(t_k), \sigma_y^2)$$

$$\mathbf{f}_k = \mathbf{A}_{\theta,k} \mathbf{f}_{k-1} + \mathbf{q}_{k-1},$$

$$y_k = \mathbf{h}^{\top} \mathbf{f}_k + \sigma_y \epsilon_k$$

There exists a **dual kernel / SDE form** for most popular GP models

$$f(t) \sim \mathcal{GP}(\mathbf{0}, K_{\theta}(t, t')),$$

$$y_k \sim \mathcal{N}(f(t_k), \sigma_y^2)$$

$$\mathbf{f}_k = \mathbf{A}_{\theta, k} \mathbf{f}_{k-1} + \mathbf{q}_{k-1},$$

$$y_k = \mathbf{h}^{\top} \mathbf{f}_k + \sigma_y \epsilon_k$$

and inference is performed in $\mathcal{O}(nm^3)$ via **Kalman filtering and smoothing** ($n = \#$ data, $m =$ dimensionality of \mathbf{f}_k)

There exists a **dual kernel / SDE form** for most popular GP models

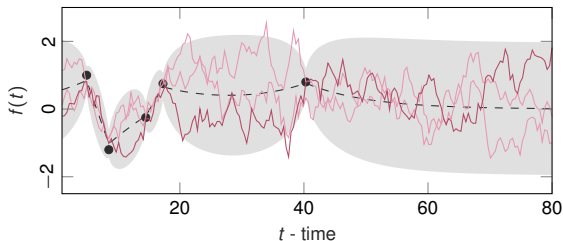
$$f(t) \sim \mathcal{GP}(0, K_{\theta}(t, t')),$$

$$y_k \sim \mathcal{N}(f(t_k), \sigma_y^2)$$

$$\mathbf{f}_k = \mathbf{A}_{\theta, k} \mathbf{f}_{k-1} + \mathbf{q}_{k-1},$$

$$y_k = \mathbf{h}^{\top} \mathbf{f}_k + \sigma_y \epsilon_k$$

and inference is performed in $\mathcal{O}(nm^3)$ via **Kalman filtering and smoothing** ($n = \#$ data, $m =$ dimensionality of \mathbf{f}_k)



Samples from the posterior distribution with $\sigma_y^2 = 0.05$.

Back to signal processing:

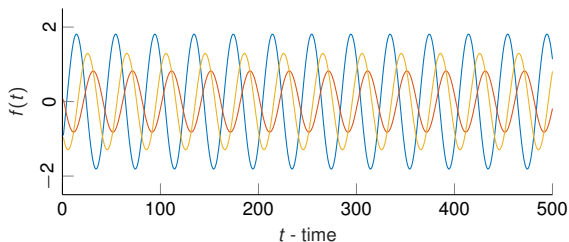
we are interested in models that capture [periodic latent structure](#).

The cosine kernel

$$K(t, t') = \cos(\omega(t - t'))$$

The cosine kernel

$$K(t, t') = \cos(\omega(t - t'))$$



Sample paths are pure sinusoids, *i.e.*, there is no noise process

The cosine kernel

$$K(t, t') = \cos(\omega(t - t'))$$

state space representation is actually a **complex ODE**:

$$\begin{pmatrix} \operatorname{Re}[\dot{f}(t)] \\ \operatorname{Im}[\dot{f}(t)] \end{pmatrix} = \begin{pmatrix} 0 & -\omega \\ \omega & 0 \end{pmatrix} \begin{pmatrix} \operatorname{Re}[f(t)] \\ \operatorname{Im}[f(t)] \end{pmatrix}$$

The cosine kernel

$$K(t, t') = \cos(\omega(t - t'))$$

state space representation is actually a **complex ODE**:

$$\begin{pmatrix} \operatorname{Re}[\dot{f}(t)] \\ \operatorname{Im}[\dot{f}(t)] \end{pmatrix} = \begin{pmatrix} 0 & -\omega \\ \omega & 0 \end{pmatrix} \begin{pmatrix} \operatorname{Re}[f(t)] \\ \operatorname{Im}[f(t)] \end{pmatrix}$$

with **discrete-time solution**:

$$\begin{pmatrix} \operatorname{Re}[f_{k+1}] \\ \operatorname{Im}[f_{k+1}] \end{pmatrix} = \begin{pmatrix} \cos \omega \Delta_k & -\sin \omega \Delta_k \\ \sin \omega \Delta_k & \cos \omega \Delta_k \end{pmatrix} \begin{pmatrix} \operatorname{Re}[f_k] \\ \operatorname{Im}[f_k] \end{pmatrix}$$

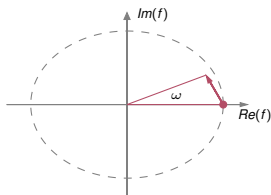
where $\Delta_k = t_{k+1} - t_k$

The cosine kernel

$$K(t, t') = \cos(\omega(t - t'))$$

state space representation is actually a **complex ODE**:

$$\begin{pmatrix} \operatorname{Re}[\dot{f}(t)] \\ \operatorname{Im}[\dot{f}(t)] \end{pmatrix} = \begin{pmatrix} 0 & -\omega \\ \omega & 0 \end{pmatrix} \begin{pmatrix} \operatorname{Re}[f(t)] \\ \operatorname{Im}[f(t)] \end{pmatrix}$$



The quasi-periodic kernel

So far we have seen the **exponential** kernel, and the **cosine** kernel.

We can construct a flexible **quasi-periodic** kernel, and its SDE form, via their product.

$$K(t, t') = \sigma^2 \exp(-|t - t'|/\ell) \cos(\omega(t - t'))$$

The quasi-periodic kernel

So far we have seen the **exponential** kernel, and the **cosine** kernel.

We can construct a flexible **quasi-periodic** kernel, and its SDE form, via their product.

$$K(t, t') = \sigma^2 \exp(-|t - t'|/\ell) \cos(\omega(t - t'))$$

$$\mathbf{F}_{\text{cos}} = \begin{pmatrix} 0 & -\omega \\ \omega & 0 \end{pmatrix},$$

$$\mathbf{F}_{\text{exp}} = -\frac{1}{\ell},$$

$$\mathbf{Q}_{\text{C,cos}} = \text{N/A},$$

$$\mathbf{Q}_{\text{C,exp}} = \frac{2\sigma^2}{\ell},$$

$$\mathbf{L}_{\text{cos}} = \text{N/A},$$

$$\mathbf{L}_{\text{exp}} = 1,$$

$$\mathbf{P}_{\infty,\text{cos}} = \mathbf{I}_2,$$

$$\mathbf{P}_{\infty,\text{exp}} = \sigma^2,$$

$$\mathbf{h}_{\text{cos}}^{\top} = (1 \quad 0),$$

$$\mathbf{h}_{\text{exp}}^{\top} = 1.$$

The quasi-periodic kernel

So far we have seen the **exponential** kernel, and the **cosine** kernel.

We can construct a flexible **quasi-periodic** kernel, and its SDE form, via their product.

$$K(t, t') = \sigma^2 \exp(-|t - t'|/\ell) \cos(\omega(t - t'))$$

$$\frac{df(t)}{dt} = \mathbf{F}f(t) + \mathbf{L}w(t)$$

$$\mathbf{F} = \mathbf{F}_{\cos} \oplus \mathbf{F}_{\exp} = \mathbf{F}_{\cos} \otimes \mathbf{I}_2 + \mathbf{I}_1 \otimes \mathbf{F}_{\exp} = \begin{pmatrix} -\frac{1}{\ell} & -\omega \\ \omega & -\frac{1}{\ell} \end{pmatrix},$$

$$\mathbf{Q}_c = \mathbf{I}_2 \otimes \mathbf{Q}_{c,\exp} = \frac{2\sigma^2}{\ell} \mathbf{I}_2,$$

$$\mathbf{L} = \mathbf{I}_2 \otimes \mathbf{L}_{\exp} = \mathbf{I}_2,$$

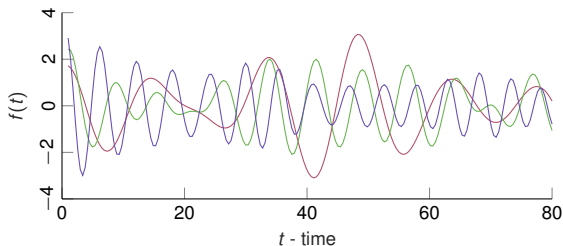
$$\mathbf{P}_{\infty} = \mathbf{I}_2 \otimes \mathbf{P}_{\infty,\exp} = \sigma^2 \mathbf{I}_2.$$

The quasi-periodic kernel

So far we have seen the **exponential** kernel, and the **cosine** kernel.

We can construct a flexible **quasi-periodic** kernel, and its SDE form, via their product.

$$K(t, t') = \sigma^2 \exp(-|t - t'|/\ell) \cos(\omega(t - t'))$$



Sample paths from the quasi-periodic GP

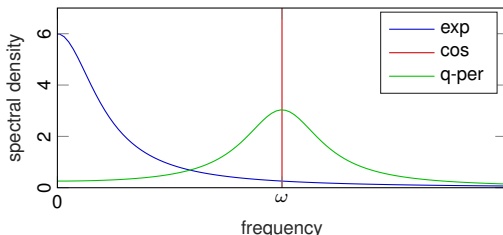
The quasi-periodic kernel

So far we have seen the **exponential** kernel, and the **cosine** kernel.

We can construct a flexible **quasi-periodic** kernel, and its SDE form, via their product.

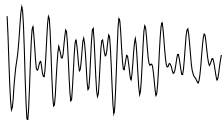
$$K(t, t') = \sigma^2 \exp(-|t - t'|/\ell) \cos(\omega(t - t'))$$

Constructing the quasi-periodic kernel in the spectral domain



The **cosine kernel** acts as a frequency shift operator on the **exponential kernel** to produce the **quasi-periodic kernel** ($\sigma^2 = 1, \ell = 3, \omega = \pi/2$).

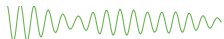
Spectral mixtures



||



+



+



⋮

Spectral mixtures

$$f(t) \sim \mathcal{GP} \left(\mathbf{0}, \sum_{d=1}^D K_{\text{q-per}}^{(d)}(t, t') \right),$$

$$y_k \sim \mathcal{N}(f(t_k), \sigma_y^2)$$

$$\mathbf{f}_k = \mathbf{A}_k \mathbf{f}_{k-1} + \mathbf{q}_{k-1},$$

$$y_k = \mathbf{h}^\top \mathbf{f}_k + \sigma_y \epsilon_{y,k}$$

Spectral mixtures

$$f(t) \sim \mathcal{GP} \left(\mathbf{0}, \sum_{d=1}^D K_{q\text{-per}}^{(d)}(t, t') \right),$$

$$y_k \sim \mathcal{N}(f(t_k), \sigma_y^2)$$

$$\mathbf{f}_k = \mathbf{A}_k \mathbf{f}_{k-1} + \mathbf{q}_{k-1},$$

$$y_k = \mathbf{h}^\top \mathbf{f}_k + \sigma_y \epsilon_{y,k}$$

$$\mathbf{A}_k = \begin{pmatrix} \exp\left(-\frac{\Delta_k}{\ell_1}\right) \begin{pmatrix} \cos \omega_1 \Delta_k & -\sin \omega_1 \Delta_k \\ \sin \omega_1 \Delta_k & \cos \omega_1 \Delta_k \end{pmatrix} & & & \\ & \ddots & & \\ & & \exp\left(-\frac{\Delta_k}{\ell_D}\right) \begin{pmatrix} \cos \omega_D \Delta_k & -\sin \omega_D \Delta_k \\ \sin \omega_D \Delta_k & \cos \omega_D \Delta_k \end{pmatrix} & \\ & & & \end{pmatrix},$$

$$\mathbf{Q}_k = \begin{pmatrix} \sigma_1^2 (1 - \exp(-\frac{2\Delta_k}{\ell_1})) \mathbf{I}_2 & & & \\ & \ddots & & \\ & & & \sigma_D^2 (1 - \exp(-\frac{2\Delta_k}{\ell_D})) \mathbf{I}_2 \end{pmatrix}.$$

$$\mathbf{h}^\top = (1 \quad 0 \quad 1 \quad 0 \quad \dots \quad 1 \quad 0).$$

Spectral mixtures

$$f(t) \sim \mathcal{GP} \left(0, \sum_{d=1}^D K_{\text{q-per}}^{(d)}(t, t') \right),$$

$$y_k \sim \mathcal{N}(f(t_k), \sigma_y^2)$$

$$\mathbf{f}_k = \mathbf{A}_k \mathbf{f}_{k-1} + \mathbf{q}_{k-1},$$

$$y_k = \mathbf{h}^\top \mathbf{f}_k + \sigma_y \epsilon_{y,k}$$

Can equivalently be written as:

$$f_{d,k} = \psi_d e^{i\omega_d \Delta t} f_{d,k-1} + \rho_d \epsilon_{d,k},$$

$$y_k = \sum_{d=1}^D \text{Re}[f_{d,k}] + \sigma_y \epsilon_{y,k},$$

with $\psi_d = \exp(-\Delta t/\ell_d)$ and $\rho_d = \sigma_d^2(1 - \exp(-2\Delta t/\ell_d))$.

Spectral mixtures

$$f(t) \sim \mathcal{GP} \left(0, \sum_{d=1}^D K_{\text{q-per}}^{(d)}(t, t') \right),$$

$$y_k \sim \mathcal{N}(f(t_k), \sigma_y^2)$$

$$\mathbf{f}_k = \mathbf{A}_k \mathbf{f}_{k-1} + \mathbf{q}_{k-1},$$

$$y_k = \mathbf{h}^\top \mathbf{f}_k + \sigma_y \epsilon_{y,k}$$

Can equivalently be written as:

$$f_{d,k} = \psi_d e^{i\omega_d \Delta t} f_{d,k-1} + \rho_d \epsilon_{d,k},$$

$$y_k = \sum_{d=1}^D \text{Re}[f_{d,k}] + \sigma_y \epsilon_{y,k},$$

with $\psi_d = \exp(-\Delta t / \ell_d)$ and $\rho_d = \sigma_d^2 (1 - \exp(-2\Delta t / \ell_d))$.

This model is known as the **probabilistic phase vocoder**.

Comparing the two model interpretations

Probabilistic phase vocoder

Fast inference via
Kalman smoothing

Fast frequency-domain
parameter learning

Interpreting the model
can be challenging

Changing the model is
hard

Spectral mixture GP

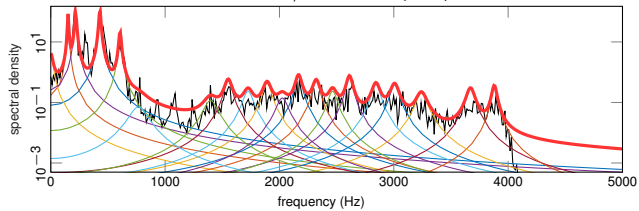
Inference slow for long
time-series

Freq.-domain parameter
learning possible

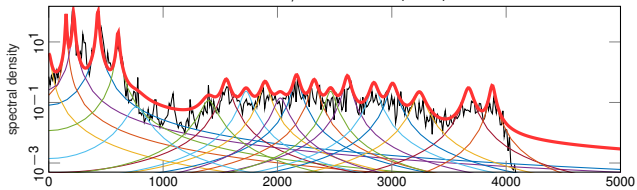
All model assumptions
encoded in the kernel

Changing the model is
easy

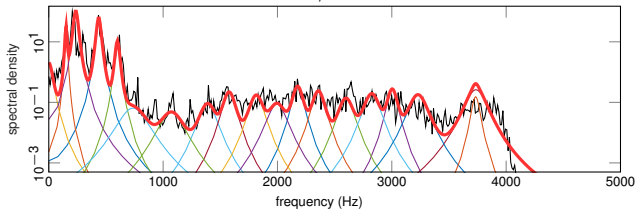
Matérn-1/2 filter bank (PPV)



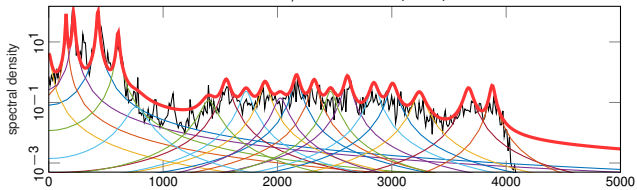
Matérn-1/2 filter bank (PPV)



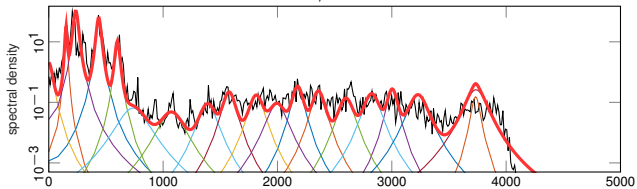
Matérn-3/2 filter bank



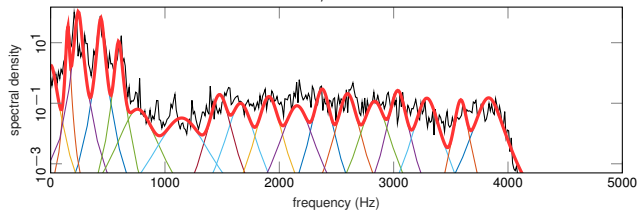
Matérn-1/2 filter bank (PPV)



Matérn-3/2 filter bank

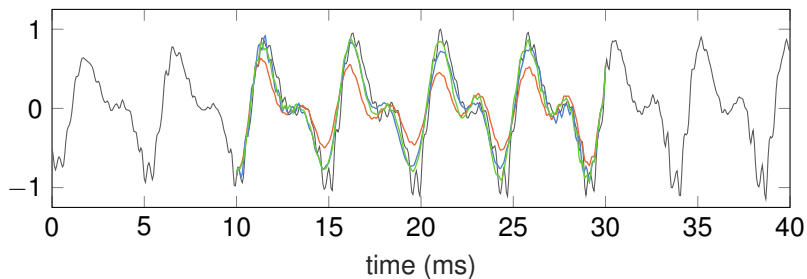


Matérn-5/2 filter bank



Missing Data Synthesis

Audio signal reconstruction



Data imputation using a filter bank composed of the following kernels:

Matérn^{1/2} (exponential) - 1st order state space form

Matérn^{3/2} - 2nd order state space form

Matérn^{5/2} - 3rd order state space form

PART II - APPROXIMATE INFERENCE IN TEMPORAL GPs

Non-conjugate state space models

$$f(t) \sim \mathcal{GP}(\mathbf{0}, K_{\theta}(t, t')),$$

$$y_k \sim p(y_k | f(t_k))$$

$$\mathbf{f}_k = \mathbf{A}_{\theta, k} \mathbf{f}_{k-1} + \mathbf{q}_{k-1},$$

$$y_k \sim p(y_k | \mathbf{h}^{\top} \mathbf{f}_k)$$

Non-conjugate state space models

$$f(t) \sim \mathcal{GP}(\mathbf{0}, K_\theta(t, t')),$$

$$y_k \sim p(y_k | f(t_k))$$

$$\mathbf{f}_k = \mathbf{A}_{\theta,k} \mathbf{f}_{k-1} + \mathbf{q}_{k-1},$$

$$y_k \sim p(y_k | \mathbf{h}^\top \mathbf{f}_k)$$

Many [approximate inference](#) methods (VB, EP) can be applied in the state space regime.

Expectation propagation (EP)

Intuitively a good fit for time series — we can process the data sequentially.

Expectation propagation (EP)

A closer look at the [Kalman filter](#):

Expectation propagation (EP)

A closer look at the [Kalman filter](#):

predict step:

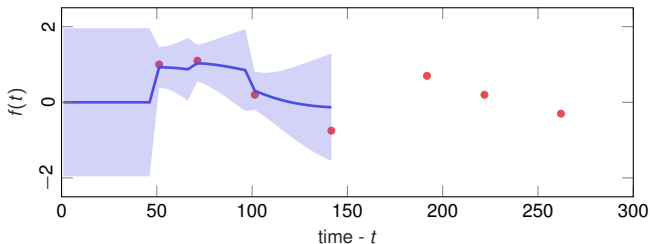
$$\begin{aligned} p(\mathbf{f}_k | y_{1:k-1}) &= N(\mathbf{m}_k^-, \mathbf{P}_k^-) = \int p(\mathbf{f}_k, \mathbf{f}_{k-1} | y_{1:k-1}) d\mathbf{f}_{k-1} \\ &= \int p(\mathbf{f}_k | \mathbf{f}_{k-1}) p(\mathbf{f}_{k-1} | y_{1:k-1}) d\mathbf{f}_{k-1} \\ &= \int N(\mathbf{A}_k, \mathbf{Q}_k) N(\mathbf{m}_{k-1}, \mathbf{P}_{k-1}) d\mathbf{f}_{k-1} \\ &= N(\mathbf{A}_k \mathbf{m}_{k-1}, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^\top + \mathbf{Q}_k) \end{aligned}$$

Expectation propagation (EP)

A closer look at the [Kalman filter](#):

predict step:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:k-1}) &= N(\mathbf{m}_k^-, \mathbf{P}_k^-) = \int p(\mathbf{f}_k, \mathbf{f}_{k-1} | y_{1:k-1}) d\mathbf{f}_{k-1} \\ &= \int p(\mathbf{f}_k | \mathbf{f}_{k-1}) p(\mathbf{f}_{k-1} | y_{1:k-1}) d\mathbf{f}_{k-1} \\ &= \int N(\mathbf{A}_k, \mathbf{Q}_k) N(\mathbf{m}_{k-1}, \mathbf{P}_{k-1}) d\mathbf{f}_{k-1} \\ &= N(\mathbf{A}_k \mathbf{m}_{k-1}, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^\top + \mathbf{Q}_k) \end{aligned}$$

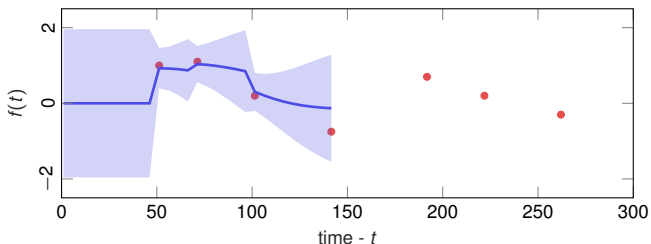


Expectation propagation (EP)

A closer look at the [Kalman filter](#):

update step:

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



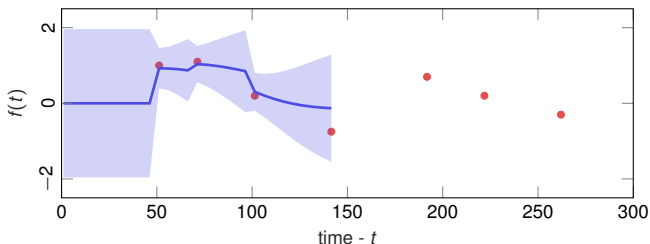
Expectation propagation (EP)

A closer look at the [Kalman filter](#):

update step:

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

if $p(y_k | f(t_k)) \sim N(\cdot, \cdot)$, the Kalman update equations are just a stable way to calculate this product of Gaussian densities.



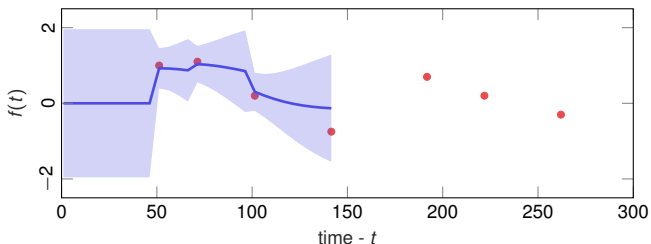
Expectation propagation (EP)

A closer look at the [Kalman filter](#):

update step:

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

How about if $p(y_k | f(t_k))$ is not Gaussian?



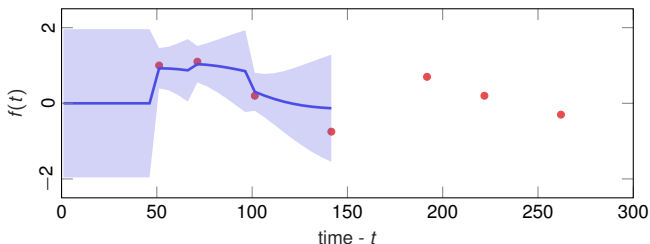
Expectation propagation (EP)

A closer look at the [Kalman filter](#):

update step:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:k}) &= N(\mathbf{m}_k, \mathbf{P}_k) \propto p(\mathbf{f}_k | y_{1:k-1}) p(y_k | f(t_k)) \\ &= \underbrace{N(\mathbf{m}_k^-, \mathbf{P}_k^-)}_{\text{"prior" over } \mathbf{f}_k \text{ conditioned on past data}} p(y_k | f(t_k)) \end{aligned}$$

"prior" over \mathbf{f}_k
conditioned on
past data

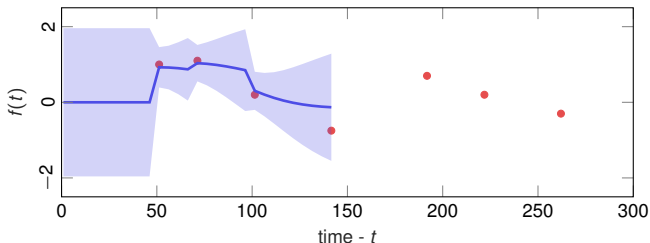


Expectation propagation (EP)

A closer look at the [Kalman filter](#):

update step:

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

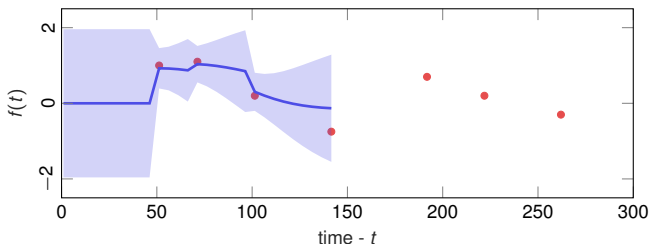


Expectation propagation (EP)

A closer look at the [Kalman filter](#):

update step:

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

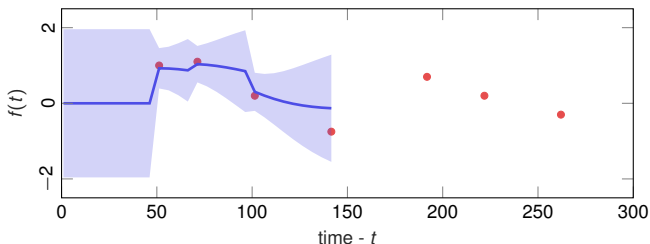


Expectation propagation (EP)

A closer look at the [Kalman filter](#):

update step:

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



Expectation propagation (EP)

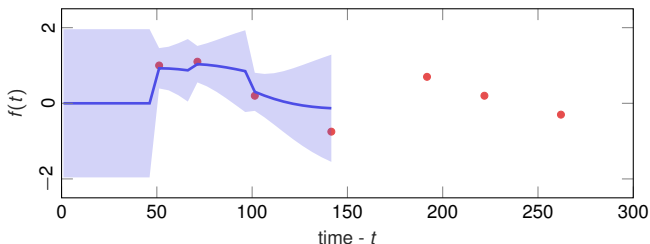
A closer look at the **Kalman filter**:

update step:

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

EP update:

↖ ↗ match moments



Expectation propagation (EP)

A closer look at the [Kalman filter](#):

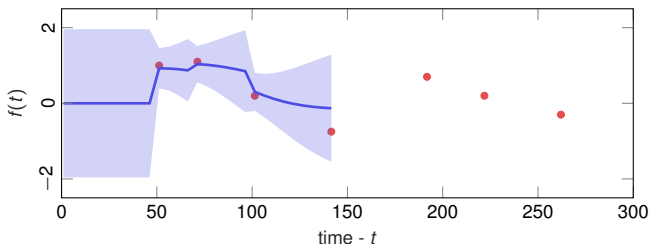
update step:

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

EP update:

↔ match moments

i.e., choose $s(\mathbf{f}_k) \sim N(\mathbf{m}_k^{\text{site}}, \mathbf{P}_k^{\text{site}})$ such that the moments are matched. [Store to be refined later.](#)



Expectation propagation (EP)

A closer look at the [Kalman filter](#):

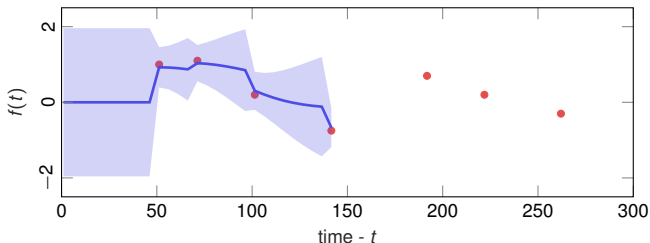
update step:

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

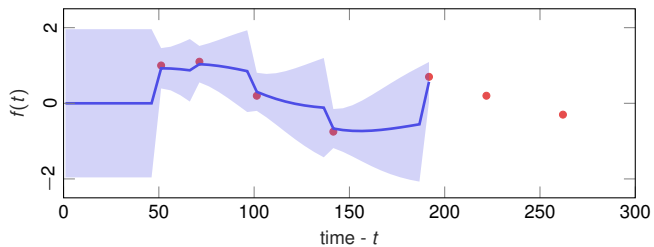
EP update:

↔ match moments

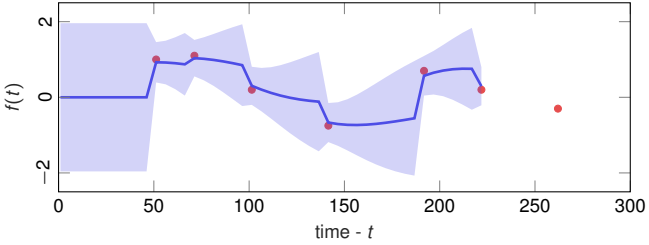
i.e., choose $s(\mathbf{f}_k) \sim N(\mathbf{m}_k^{\text{site}}, \mathbf{P}_k^{\text{site}})$ such that the moments are matched. [Store to be refined later.](#)



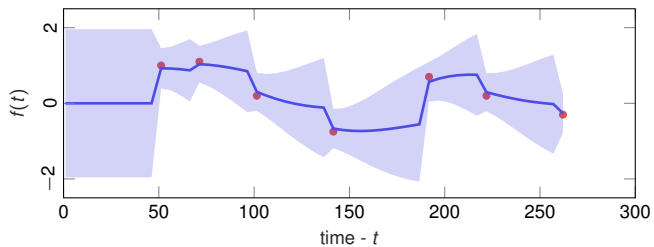
Expectation propagation (EP)



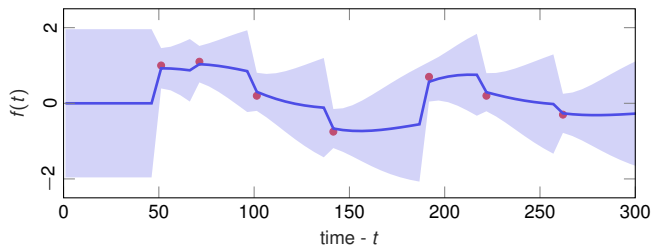
Expectation propagation (EP)



Expectation propagation (EP)



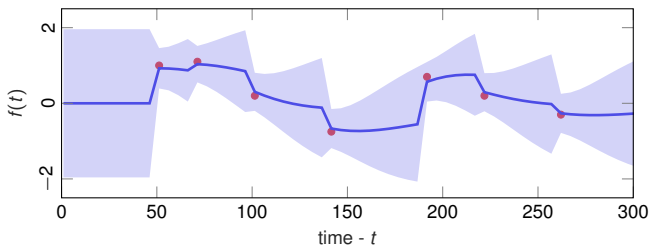
Expectation propagation (EP)



Expectation propagation (EP)

Now consider the **RTS Smoother**:

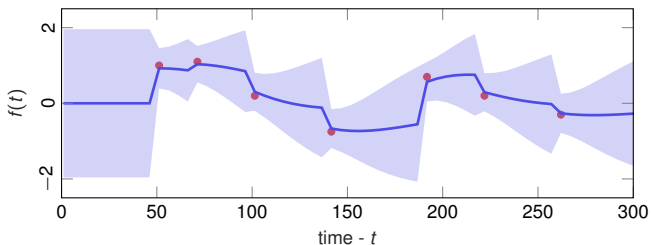
- update predictions with future observations
- **refine the EP sites** along the way



Expectation propagation (EP)

Now consider the **RTS Smoother**:

$$p(\mathbf{f}_k | \mathbf{y}_{1:T}) \propto p(\mathbf{f}_k | \mathbf{y}_{k+1:N}) p(\mathbf{f}_k | \mathbf{y}_{1:k})$$

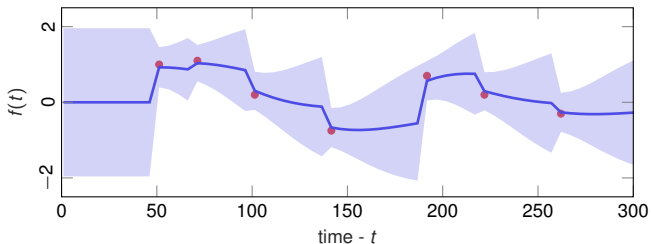


Expectation propagation (EP)

Now consider the **RTS Smoother**:

We have the full (marginal) posterior, so we must explicitly **remove the sites**:

$$p_{\text{cavity}}(\mathbf{f}_k) = p(\mathbf{f}_k | y_{1:T}) / s_{\text{old}}(\mathbf{f}_k)$$



Expectation propagation (EP)

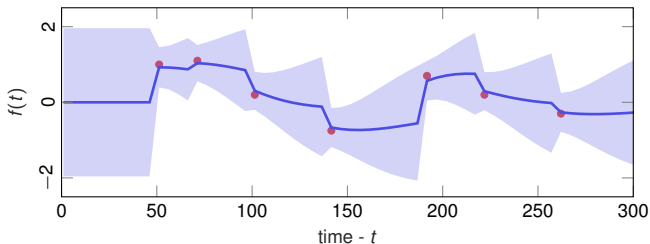
Now consider the **RTS Smoother**:

We have the full (marginal) posterior, so we must explicitly **remove the sites**:

$$p_{\text{cavity}}(\mathbf{f}_k) = p(\mathbf{f}_k | y_{1:T}) / s_{\text{old}}(\mathbf{f}_k)$$

Tilted distribution:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:T}) &= p_{\text{cavity}}(\mathbf{f}_k) p(y_k | f(t_k)) \\ &\approx p_{\text{cavity}}(\mathbf{f}_k) s_{\text{new}}(\mathbf{f}_k) \end{aligned}$$



Expectation propagation (EP)

Now consider the **RTS Smoother**:

We have the full (marginal) posterior, so we must explicitly **remove the sites**:

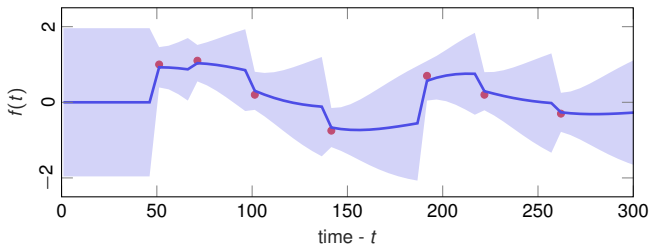
$$p_{\text{cavity}}(\mathbf{f}_k) = p(\mathbf{f}_k | y_{1:T}) / s_{\text{old}}(\mathbf{f}_k)$$

Tilted distribution:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:T}) &= p_{\text{cavity}}(\mathbf{f}_k) p(y_k | f(t_k)) \\ &\approx p_{\text{cavity}}(\mathbf{f}_k) s_{\text{new}}(\mathbf{f}_k) \end{aligned}$$

EP update:

↔ match moments



Expectation propagation (EP)

Now consider the **RTS Smoother**:

We have the full (marginal) posterior, so we must explicitly **remove the sites**:

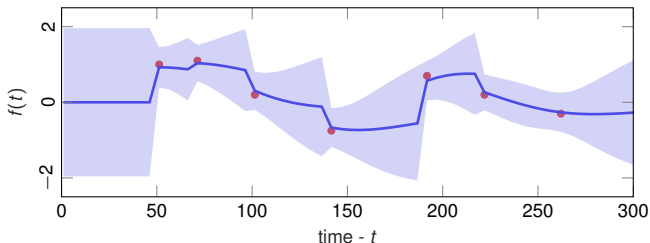
$$p_{\text{cavity}}(\mathbf{f}_k) = p(\mathbf{f}_k | y_{1:T}) / s_{\text{old}}(\mathbf{f}_k)$$

Tilted distribution:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:T}) &= p_{\text{cavity}}(\mathbf{f}_k) p(y_k | f(t_k)) \\ &\approx p_{\text{cavity}}(\mathbf{f}_k) s_{\text{new}}(\mathbf{f}_k) \end{aligned}$$

EP update:

↔ match moments



Expectation propagation (EP)

Now consider the **RTS Smoother**:

We have the full (marginal) posterior, so we must explicitly **remove the sites**:

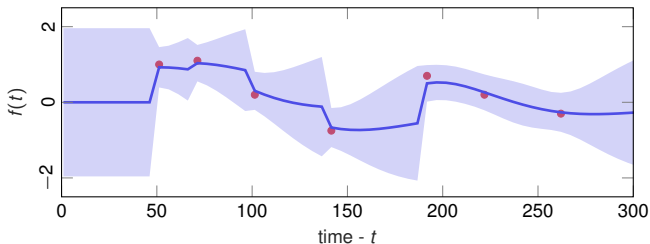
$$p_{\text{cavity}}(\mathbf{f}_k) = p(\mathbf{f}_k | y_{1:T}) / s_{\text{old}}(\mathbf{f}_k)$$

Tilted distribution:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:T}) &= p_{\text{cavity}}(\mathbf{f}_k) p(y_k | f(t_k)) \\ &\approx p_{\text{cavity}}(\mathbf{f}_k) s_{\text{new}}(\mathbf{f}_k) \end{aligned}$$

EP update:

↔ match moments



Expectation propagation (EP)

Now consider the **RTS Smoother**:

We have the full (marginal) posterior, so we must explicitly **remove the sites**:

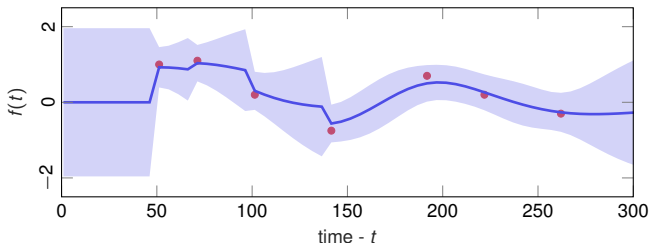
$$p_{\text{cavity}}(\mathbf{f}_k) = p(\mathbf{f}_k | y_{1:T}) / s_{\text{old}}(\mathbf{f}_k)$$

Tilted distribution:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:T}) &= p_{\text{cavity}}(\mathbf{f}_k) p(y_k | f(t_k)) \\ &\approx p_{\text{cavity}}(\mathbf{f}_k) s_{\text{new}}(\mathbf{f}_k) \end{aligned}$$

EP update:

↔ match moments



Expectation propagation (EP)

Now consider the **RTS Smoother**:

We have the full (marginal) posterior, so we must explicitly **remove the sites**:

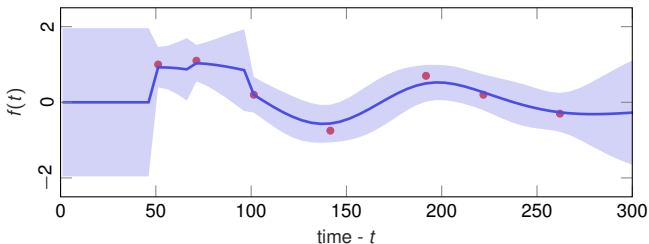
$$p_{\text{cavity}}(\mathbf{f}_k) = p(\mathbf{f}_k | y_{1:T}) / s_{\text{old}}(\mathbf{f}_k)$$

Tilted distribution:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:T}) &= p_{\text{cavity}}(\mathbf{f}_k) p(y_k | f(t_k)) \\ &\approx p_{\text{cavity}}(\mathbf{f}_k) s_{\text{new}}(\mathbf{f}_k) \end{aligned}$$

EP update:

↔ match moments



Expectation propagation (EP)

Now consider the **RTS Smoother**:

We have the full (marginal) posterior, so we must explicitly **remove the sites**:

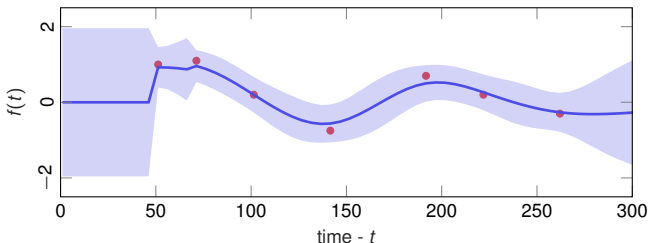
$$p_{\text{cavity}}(\mathbf{f}_k) = p(\mathbf{f}_k | y_{1:T}) / s_{\text{old}}(\mathbf{f}_k)$$

Tilted distribution:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:T}) &= p_{\text{cavity}}(\mathbf{f}_k) p(y_k | f(t_k)) \\ &\approx p_{\text{cavity}}(\mathbf{f}_k) s_{\text{new}}(\mathbf{f}_k) \end{aligned}$$

EP update:

↔ match moments



Expectation propagation (EP)

Now consider the **RTS Smoother**:

We have the full (marginal) posterior, so we must explicitly **remove the sites**:

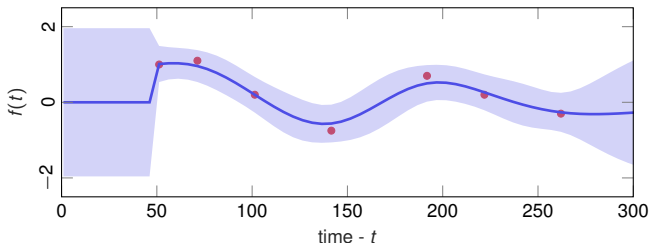
$$p_{\text{cavity}}(\mathbf{f}_k) = p(\mathbf{f}_k | y_{1:T}) / s_{\text{old}}(\mathbf{f}_k)$$

Tilted distribution:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:T}) &= p_{\text{cavity}}(\mathbf{f}_k) p(y_k | f(t_k)) \\ &\approx p_{\text{cavity}}(\mathbf{f}_k) s_{\text{new}}(\mathbf{f}_k) \end{aligned}$$

EP update:

↔ match moments



Expectation propagation (EP)

Now consider the **RTS Smoother**:

We have the full (marginal) posterior, so we must explicitly **remove the sites**:

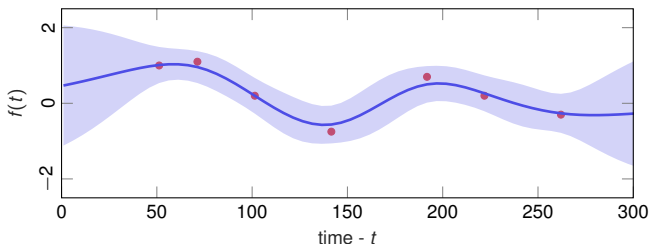
$$p_{\text{cavity}}(\mathbf{f}_k) = p(\mathbf{f}_k | y_{1:T}) / s_{\text{old}}(\mathbf{f}_k)$$

Tilted distribution:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:T}) &= p_{\text{cavity}}(\mathbf{f}_k) p(y_k | f(t_k)) \\ &\approx p_{\text{cavity}}(\mathbf{f}_k) s_{\text{new}}(\mathbf{f}_k) \end{aligned}$$

EP update:

↔ match moments



Expectation propagation (EP)

Now consider the **RTS Smoother**:

We have the full (marginal) posterior, so we must explicitly **remove the sites**:

$$p_{\text{cavity}}(\mathbf{f}_k) = p(\mathbf{f}_k | y_{1:T}) / s_{\text{old}}(\mathbf{f}_k)$$

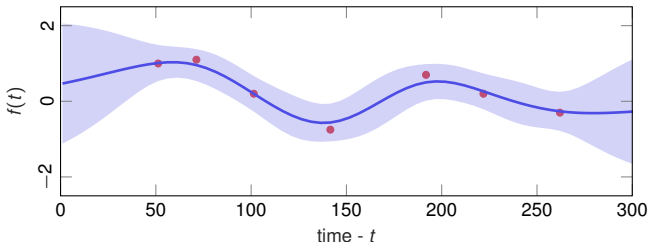
Tilted distribution:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:T}) &= p_{\text{cavity}}(\mathbf{f}_k) p(y_k | f(t_k)) \\ &\approx p_{\text{cavity}}(\mathbf{f}_k) s_{\text{new}}(\mathbf{f}_k) \end{aligned}$$

EP update:

↔ match moments

The new sites $s_{\text{new}}(\mathbf{f}_k)$ will be used on the next **forward** pass.



Expectation propagation (EP)

Now consider the **RTS Smoother**:

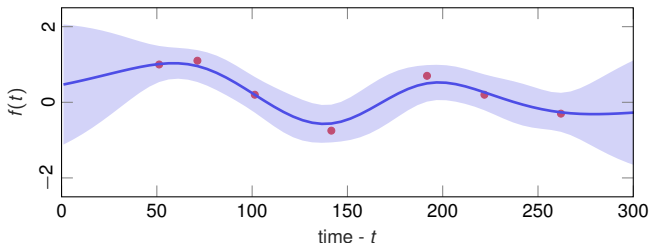
We have the full (marginal) posterior, so we must explicitly **remove the sites**:

$$p_{\text{cavity}}(\mathbf{f}_k) = p(\mathbf{f}_k | y_{1:T}) / s_{\text{old}}^\alpha(\mathbf{f}_k)$$

Tilted distribution:

$$\begin{aligned} p(\mathbf{f}_k | y_{1:T}) &= p_{\text{cavity}}(\mathbf{f}_k) p^\alpha(y_k | f(t_k)) \\ &\approx p_{\text{cavity}}(\mathbf{f}_k) s_{\text{new}}^\alpha(\mathbf{f}_k) \end{aligned}$$

Can add in the usual EP extras: power (α) and damping



An example: nonstationary TF-analysis

We apply this power EP method to a nonstationary extension of [time-frequency analysis model](#):

An example: nonstationary TF-analysis

We apply this power EP method to a nonstationary extension of [time-frequency analysis model](#):

prior:

$$f_d(t) \sim \text{GP}(0, K_{\text{q-periodic}}^{(d)}(t, t')), \quad d = 1, 2, \dots, D$$
$$\log a_d(t) \sim \text{GP}(0, K_{\text{Matérn}}^{(n)}(t, t')),$$

An example: nonstationary TF-analysis

We apply this power EP method to a nonstationary extension of [time-frequency analysis model](#):

prior:

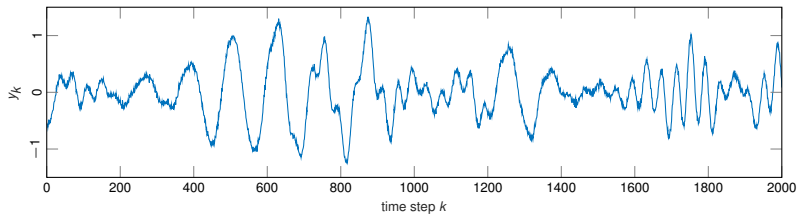
$$f_d(t) \sim \text{GP}(0, K_{\text{q-periodic}}^{(d)}(t, t')), \quad d = 1, 2, \dots, D$$
$$\log a_d(t) \sim \text{GP}(0, K_{\text{Matérn}}^{(n)}(t, t')),$$

likelihood:

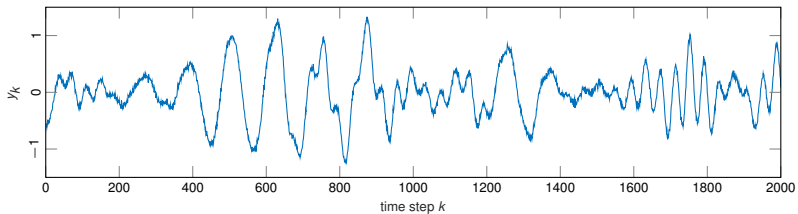
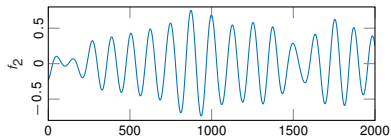
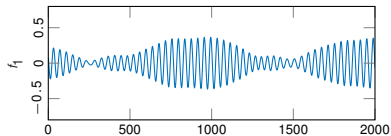
$$y_k = \sum_d a_d(t_k) f_d(t_k) + \sigma_y \epsilon_k$$

- $f_d(t)$: frequency components
- $a_d(t)$: positive amplitudes

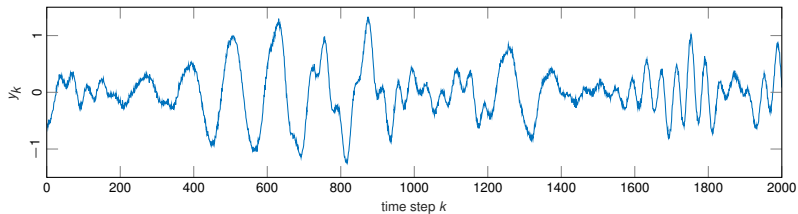
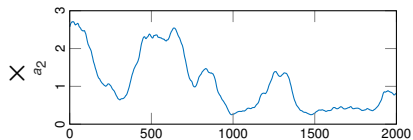
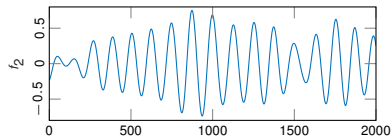
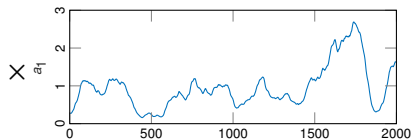
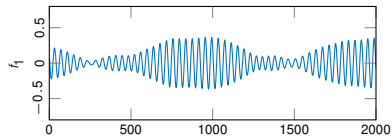
Toy data



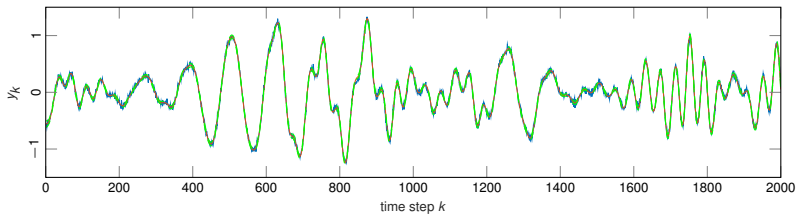
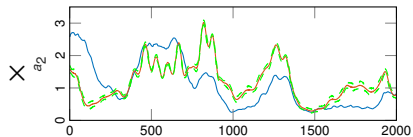
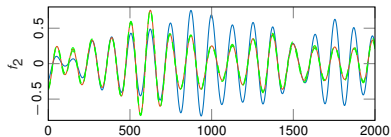
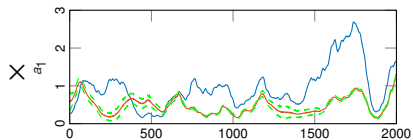
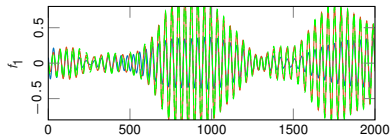
Toy data



Toy data



Toy data



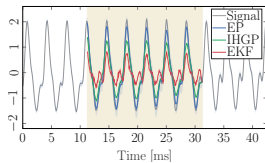
Applications

The model can, **without modification**, be applied to:

Applications

The model can, **without modification**, be applied to:

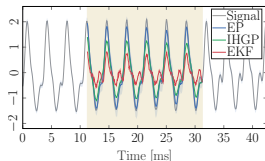
Missing Data Synthesis



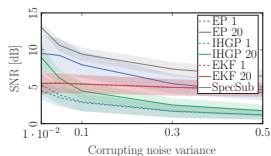
Applications

The model can, **without modification**, be applied to:

Missing Data Synthesis



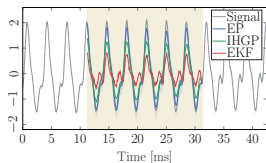
Denoising



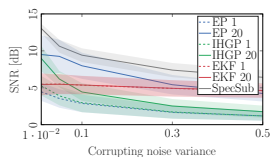
Applications

The model can, **without modification**, be applied to:

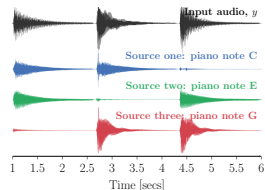
Missing Data Synthesis



Denoising



Source Separation



One issue with the power EP method

The crucial **moment matching** step involves calculating the following (intractable) expectation:

$$\mathcal{Z}_{\text{tilted}} = \mathbb{E}_{p_{\text{cavity}}(\mathbf{f}_k)} [\rho(y_k | \mathbf{f}_k)^\alpha]$$

One issue with the power EP method

The crucial **moment matching** step involves calculating the following (intractable) expectation:

$$\mathcal{Z}_{\text{tilted}} = \mathbb{E}_{p_{\text{cavity}}(\mathbf{f}_k)} [\rho(y_k | \mathbf{f}_k)^\alpha]$$

This can be a very high-dimensional integral for some models!

PART III - SCALABLE GLOBAL INFERENCE VIA LOCAL LINEARISATION

Another contradiction

Approximate inference in practice

In practical (industrial) applications, the extended Kalman filter (EKF) is still the tool of choice

Modern day approximate inference (VB, EP, ...) is more general and better approximates the true posterior

Pseudo-code for Kalman filter EP

for $k = 1 : T$

Kalman predict:

$$p(\mathbf{f}_k | \mathbf{y}_{1:k-1}) = N(\mathbf{m}_k^-, \mathbf{P}_k^-) = N(\mathbf{A}_k \mathbf{m}_{k-1}, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^T + \mathbf{Q}_k)$$

Moment match:

$$\mathcal{L}_k = \log \mathbb{E}_{N(\mathbf{m}_k^-, \mathbf{P}_k^-)} [\rho(\mathbf{y}_k | \mathbf{f}_k)^\alpha],$$

$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^- - \left(\frac{d^2 \mathcal{L}_k}{d\mathbf{m}_k^2} \right)^{-1} \frac{d\mathcal{L}_k}{d\mathbf{m}_k}, \quad \mathbf{P}_k^{\text{site}} = \alpha \left(-\mathbf{P}_k^- - \left(\frac{d^2 \mathcal{L}_k}{d\mathbf{m}_k^2} \right)^{-1} \right).$$

Kalman update:

$$\mathbf{S}_k = \mathbf{P}_k^- + \mathbf{P}_k^{\text{site}},$$

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{S}_k^{-1},$$

$$\mathbf{m}_k = \mathbf{m}_k^- + \mathbf{K}_k (\mathbf{m}_k^{\text{site}} - \mathbf{m}_k^-),$$

$$\mathbf{P}_k = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T.$$

$$p(\mathbf{f}_k | \mathbf{y}_{1:k}) = N(\mathbf{m}_k, \mathbf{P}_k)$$

end for

Pseudo-code for Kalman filter EP

for $k = 1 : T$

Kalman predict:

$$p(\mathbf{f}_k | \mathbf{y}_{1:k-1}) = N(\mathbf{m}_k^-, \mathbf{P}_k^-) = N(\mathbf{A}_k \mathbf{m}_{k-1}, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^T + \mathbf{Q}_k)$$

Moment match:

$$\mathcal{L}_k = \log \mathbb{E}_{N(\mathbf{m}_k^-, \mathbf{P}_k^-)} [\rho(\mathbf{y}_k | \mathbf{f}_k)^\alpha],$$

$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^- - \left(\frac{d^2 \mathcal{L}_k}{d\mathbf{m}_k^2} \right)^{-1} \frac{d\mathcal{L}_k}{d\mathbf{m}_k}, \quad \mathbf{P}_k^{\text{site}} = \alpha \left(-\mathbf{P}_k^- - \left(\frac{d^2 \mathcal{L}_k}{d\mathbf{m}_k^2} \right)^{-1} \right).$$

Kalman update:

$$\mathbf{S}_k = \mathbf{P}_k^- + \mathbf{P}_k^{\text{site}},$$

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{S}_k^{-1},$$

$$\mathbf{m}_k = \mathbf{m}_k^- + \mathbf{K}_k (\mathbf{m}_k^{\text{site}} - \mathbf{m}_k^-),$$

$$\mathbf{P}_k = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T.$$

$$p(\mathbf{f}_k | \mathbf{y}_{1:k}) = N(\mathbf{m}_k, \mathbf{P}_k)$$

end for

Pseudo-code for Kalman filter EP

$$\mathcal{L}_k = \log \mathbb{E}_{\mathbf{N}(\mathbf{m}_k^-, \mathbf{P}_k^-)} [\rho(\mathbf{y}_k | \mathbf{f}_k)^\alpha]$$

Pseudo-code for Kalman filter EP

Likelihood $p(\mathbf{y}_k | \mathbf{f}_k) = h(\mathbf{f}_k, \mathbf{r}_k)$ is a nonlinear function of Gaussian process \mathbf{f}_k and Gaussian observation noise $\mathbf{r}_k \sim \mathbf{N}(\mathbf{0}, \mathbf{R}_k)$

$$\mathcal{L}_k = \log \mathbb{E}_{\mathbf{N}(\mathbf{m}_k^-, \mathbf{P}_k^-)} [\rho(\mathbf{y}_k | \mathbf{f}_k)^\alpha]$$

Pseudo-code for Kalman filter EP

Likelihood $p(\mathbf{y}_k | \mathbf{f}_k) = h(\mathbf{f}_k, \mathbf{r}_k)$ is a nonlinear function of Gaussian process \mathbf{f}_k and Gaussian observation noise $\mathbf{r}_k \sim \mathbf{N}(\mathbf{0}, \mathbf{R}_k)$

$$\mathcal{L}_k = \log \mathbb{E}_{\mathbf{N}(\mathbf{m}_k^-, \mathbf{P}_k^-)} [\rho(\mathbf{y}_k | \mathbf{f}_k)^\alpha]$$

Linearisation w.r.t. \mathbf{f}_k and \mathbf{r}_k via first-order Taylor expansion leads to a Gaussian approximation:

$$\begin{aligned} h(\mathbf{f}_k, \mathbf{r}_k) &\approx h(\mathbf{m}_k^-, \mathbf{0}) + \mathbf{J}_{\mathbf{f}_k}(\mathbf{f}_k - \mathbf{m}_k^-) + \mathbf{J}_{\mathbf{r}_k} \mathbf{r}_k \\ &= \mathbf{N}(\mathbf{y}_k | h(\mathbf{m}_k^-, \mathbf{0}) + \mathbf{J}_{\mathbf{f}_k}(\mathbf{f}_k - \mathbf{m}_k^-), \mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top) \end{aligned}$$

Pseudo-code for Kalman filter EP

$$\mathcal{L}_k = \log \mathbb{E}_{\mathbf{N}(\mathbf{m}_k^-, \mathbf{P}_k^-)} \left[\mathbf{N} \left(\mathbf{y}_k \mid h(\mathbf{m}_k^-, \mathbf{0}) + \mathbf{J}_{\mathbf{f}_k} (\mathbf{f}_k - \mathbf{m}_k^-), \mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top \right)^\alpha \right]$$

Pseudo-code for Kalman filter EP

$$\mathcal{L}_k = c + \log N(\mathbf{y}_k | h(\mathbf{m}_k^-, \mathbf{0}), \mathbf{J}_{f_k} \mathbf{P}_k^- \mathbf{J}_{f_k}^\top + \frac{1}{\alpha} \mathbf{J}_{r_k} \mathbf{R}_k \mathbf{J}_{r_k}^\top)$$

Pseudo-code for Kalman filter EP

for $k = 1 : T$

Kalman predict:

$$p(\mathbf{f}_k | \mathbf{y}_{1:k-1}) = N(\mathbf{m}_k^-, \mathbf{P}_k^-) = N(\mathbf{A}_k \mathbf{m}_{k-1}, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^\top + \mathbf{Q}_k)$$

Moment match:

$$\mathcal{L}_k = c + \log N(\mathbf{y}_k | h(\mathbf{m}_k^-, \mathbf{0}), \mathbf{J}_{\mathbf{f}_k} \mathbf{P}_k^- \mathbf{J}_{\mathbf{f}_k}^\top + \frac{1}{\alpha} \mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top),$$

$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^- - \left(\frac{d^2 \mathcal{L}_k}{d\mathbf{m}_k^2} \right)^{-1} \frac{d\mathcal{L}_k}{d\mathbf{m}_k}, \quad \mathbf{P}_k^{\text{site}} = \alpha \left(-\mathbf{P}_k^- - \left(\frac{d^2 \mathcal{L}_k}{d\mathbf{m}_k^2} \right)^{-1} \right).$$

Kalman update:

$$\mathbf{S}_k = \mathbf{P}_k^- + \mathbf{P}_k^{\text{site}},$$

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{S}_k^{-1},$$

$$\mathbf{m}_k = \mathbf{m}_k^- + \mathbf{K}_k (\mathbf{m}_k^{\text{site}} - \mathbf{m}_k^-),$$

$$\mathbf{P}_k = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^\top.$$

$$p(\mathbf{f}_k | \mathbf{y}_{1:k}) = N(\mathbf{m}_k, \mathbf{P}_k)$$

end for

Pseudo-code for Kalman filter EP

for $k = 1 : T$

Kalman predict:

$$p(\mathbf{f}_k | \mathbf{y}_{1:k-1}) = N(\mathbf{m}_k^-, \mathbf{P}_k^-) = N(\mathbf{A}_k \mathbf{m}_{k-1}, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^\top + \mathbf{Q}_k)$$

Moment match:

$$\mathbf{P}_k^{\text{site}} = \left(\mathbf{J}_{\mathbf{f}_k}^\top (\mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top)^{-1} \mathbf{J}_{\mathbf{f}_k} \right)^{-1},$$

$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^- + (\mathbf{P}_k^{\text{site}} + \mathbf{P}_k^-) \mathbf{J}_{\mathbf{f}_k}^\top (\mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top + \mathbf{J}_{\mathbf{f}_k} \mathbf{P}_k^- \mathbf{J}_{\mathbf{f}_k}^\top)^{-1} (\mathbf{y}_k - h(\mathbf{m}_k^-, \mathbf{0}))$$

Kalman update:

$$\mathbf{S}_k = \mathbf{P}_k^- + \mathbf{P}_k^{\text{site}},$$

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{S}_k^{-1},$$

$$\mathbf{m}_k = \mathbf{m}_k^- + \mathbf{K}_k (\mathbf{m}_k^{\text{site}} - \mathbf{m}_k^-),$$

$$\mathbf{P}_k = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^\top.$$

$$p(\mathbf{f}_k | \mathbf{y}_{1:k}) = N(\mathbf{m}_k, \mathbf{P}_k)$$

end for

Pseudo-code for Kalman filter EP

for $k = 1 : T$

Kalman predict:

$$p(\mathbf{f}_k | \mathbf{y}_{1:k-1}) = N(\mathbf{m}_k^-, \mathbf{P}_k^-) = N(\mathbf{A}_k \mathbf{m}_{k-1}, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^\top + \mathbf{Q}_k)$$

Moment match:

Kalman update:

$$\mathbf{S}_k = \mathbf{J}_{\mathbf{f}_k} \mathbf{P}_k^- \mathbf{J}_{\mathbf{f}_k}^\top + \mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top,$$

$$\mathbf{m}_k = \mathbf{m}_k^- + \mathbf{K}_k (\mathbf{y}_k - h(\mathbf{m}_k^-, \mathbf{0})),$$

$$p(\mathbf{f}_k | \mathbf{y}_{1:k}) = N(\mathbf{m}_k, \mathbf{P}_k)$$

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{J}_{\mathbf{f}_k}^\top \mathbf{S}_k^{-1},$$

$$\mathbf{P}_k = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^\top.$$

end for

Pseudo-code for Kalman filter EP

for $k = 1 : T$

Kalman predict:

$$p(\mathbf{f}_k | \mathbf{y}_{1:k-1}) = N(\mathbf{m}_k^-, \mathbf{P}_k^-) = N(\mathbf{A}_k \mathbf{m}_{k-1}, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^\top + \mathbf{Q}_k)$$

Moment match:

Kalman update:

$$\mathbf{S}_k = \mathbf{J}_{\mathbf{f}_k} \mathbf{P}_k^- \mathbf{J}_{\mathbf{f}_k}^\top + \mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top,$$

$$\mathbf{m}_k = \mathbf{m}_k^- + \mathbf{K}_k (\mathbf{y}_k - h(\mathbf{m}_k^-, \mathbf{0})),$$

$$p(\mathbf{f}_k | \mathbf{y}_{1:k}) = N(\mathbf{m}_k, \mathbf{P}_k)$$

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{J}_{\mathbf{f}_k}^\top \mathbf{S}_k^{-1},$$

$$\mathbf{P}_k = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^\top.$$

**This is exactly
the EKF**

end for

Pseudo-code for Kalman filter EP

for $k = 1 : T$

Kalman predict:

$$p(\mathbf{f}_k | \mathbf{y}_{1:k-1}) = N(\mathbf{m}_k^-, \mathbf{P}_k^-) = N(\mathbf{A}_k \mathbf{m}_{k-1}, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^\top + \mathbf{Q}_k)$$

Moment match:

But storing the sites
allowed us to iterate

Kalman update:

$$\mathbf{S}_k = \mathbf{J}_{\mathbf{f}_k} \mathbf{P}_k^- \mathbf{J}_{\mathbf{f}_k}^\top + \mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top,$$

$$\mathbf{m}_k = \mathbf{m}_k^- + \mathbf{K}_k (\mathbf{y}_k - h(\mathbf{m}_k^-, \mathbf{0})),$$

$$p(\mathbf{f}_k | \mathbf{y}_{1:k}) = N(\mathbf{m}_k, \mathbf{P}_k)$$

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{J}_{\mathbf{f}_k}^\top \mathbf{S}_k^{-1},$$

$$\mathbf{P}_k = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^\top.$$

This is exactly
the EKF

end for

Pseudo-code for Kalman filter EP

for $k = 1 : T$

Kalman predict:

$$p(\mathbf{f}_k | \mathbf{y}_{1:k-1}) = N(\mathbf{m}_k^-, \mathbf{P}_k^-) = N(\mathbf{A}_k \mathbf{m}_{k-1}, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^\top + \mathbf{Q}_k)$$

Moment match:

$$\mathbf{P}_k^{\text{site}} = \left(\mathbf{J}_{\mathbf{f}_k}^\top (\mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top)^{-1} \mathbf{J}_{\mathbf{f}_k} \right)^{-1},$$

$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^- + (\mathbf{P}_k^{\text{site}} + \mathbf{P}_k^-) \mathbf{J}_{\mathbf{f}_k}^\top (\mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top + \mathbf{J}_{\mathbf{f}_k} \mathbf{P}_k^- \mathbf{J}_{\mathbf{f}_k}^\top)^{-1} (\mathbf{y}_k - h(\mathbf{m}_k^-, \mathbf{0}))$$

Kalman update:

$$\mathbf{S}_k = \mathbf{P}_k^- + \mathbf{P}_k^{\text{site}},$$

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{S}_k^{-1},$$

$$\mathbf{m}_k = \mathbf{m}_k^- + \mathbf{K}_k (\mathbf{m}_k^{\text{site}} - \mathbf{m}_k^-),$$

$$\mathbf{P}_k = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^\top.$$

$$p(\mathbf{f}_k | \mathbf{y}_{1:k}) = N(\mathbf{m}_k, \mathbf{P}_k)$$

end for

Pseudo-code for Kalman filter EP

for $k = 1 : T$

Kalman predict:

$$p(\mathbf{f}_k | \mathbf{y}_{1:k-1}) = N(\mathbf{m}_k^-, \mathbf{P}_k^-) = N(\mathbf{A}_k \mathbf{m}_{k-1}, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^\top + \mathbf{Q}_k)$$

Moment match:

$$\mathbf{P}_k^{\text{site}} = \left(\mathbf{J}_{\mathbf{f}_k}^\top (\mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top)^{-1} \mathbf{J}_{\mathbf{f}_k} \right)^{-1},$$

$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^- + (\mathbf{P}_k^{\text{site}} + \mathbf{P}_k^-) \mathbf{J}_{\mathbf{f}_k}^\top (\mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top + \mathbf{J}_{\mathbf{f}_k} \mathbf{P}_k^- \mathbf{J}_{\mathbf{f}_k}^\top)^{-1} (\mathbf{y}_k - h(\mathbf{m}_k^-, \mathbf{0}))$$

Kalman update:

$$\mathbf{S}_k = \mathbf{P}_k^- + \mathbf{P}_k^{\text{site}},$$

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{S}_k^{-1},$$

$$\mathbf{m}_k = \mathbf{m}_k^- + \mathbf{K}_k (\mathbf{m}_k^{\text{site}} - \mathbf{m}_k^-),$$

$$\mathbf{P}_k = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^\top.$$

$$p(\mathbf{f}_k | \mathbf{y}_{1:k}) = N(\mathbf{m}_k, \mathbf{P}_k)$$

This is now a
globally iterated EKF

end for

Pseudo-code for Kalman filter EP

for $k = 1 : T$

Kalman predict:

$$p(\mathbf{f}_k | \mathbf{y}_{1:k-1}) = N(\mathbf{m}_k^-, \mathbf{P}_k^-) = N(\mathbf{A}_k \mathbf{m}_{k-1}, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^\top + \mathbf{Q}_k)$$

Moment match:

$$\mathbf{P}_k^{\text{site}} = \left(\mathbf{J}_{\mathbf{f}_k}^\top (\mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top)^{-1} \mathbf{J}_{\mathbf{f}_k} \right)^{-1},$$

Site updates don't
depend on α

$$\mathbf{m}_k^{\text{site}} = \mathbf{m}_k^- + (\mathbf{P}_k^{\text{site}} + \mathbf{P}_k^-) \mathbf{J}_{\mathbf{f}_k}^\top (\mathbf{J}_{\mathbf{r}_k} \mathbf{R}_k \mathbf{J}_{\mathbf{r}_k}^\top + \mathbf{J}_{\mathbf{f}_k} \mathbf{P}_k^- \mathbf{J}_{\mathbf{f}_k}^\top)^{-1} (\mathbf{y}_k - h(\mathbf{m}_k^-, \mathbf{0}))$$

Kalman update:

$$\mathbf{S}_k = \mathbf{P}_k^- + \mathbf{P}_k^{\text{site}},$$

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{S}_k^{-1},$$

$$\mathbf{m}_k = \mathbf{m}_k^- + \mathbf{K}_k (\mathbf{m}_k^{\text{site}} - \mathbf{m}_k^-),$$

$$\mathbf{P}_k = \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^\top.$$

$$p(\mathbf{f}_k | \mathbf{y}_{1:k}) = N(\mathbf{m}_k, \mathbf{P}_k)$$

This is now a
globally iterated EKF

end for

Unifying EP and the EKF

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

Unifying EP and the EKF

- For sequential data, the EKF is equivalent to single-sweep EP where the moment matching integral is solved via linearisation.
- Our algorithm iteratively refines the EKF by linearising about the cavity, rather than the filter predictions (prior).

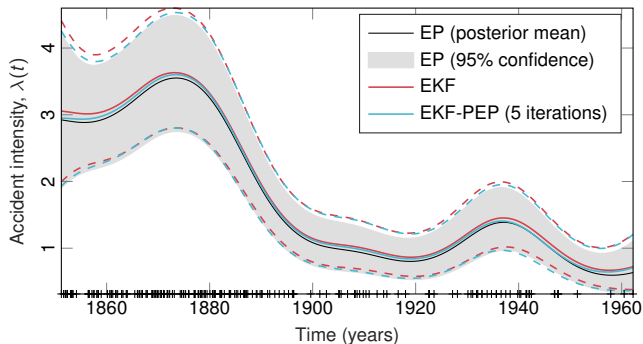
Unifying EP and the EKF

- For sequential data, the EKF is equivalent to single-sweep EP where the moment matching integral is solved via linearisation.
- Our algorithm iteratively refines the EKF by linearising about the cavity, rather than the filter predictions (prior).
- When $\alpha = 0$ we recover the Posterior Linearisation Filter. Lack of cavity calculation makes it very stable.

Unifying EP and the EKF

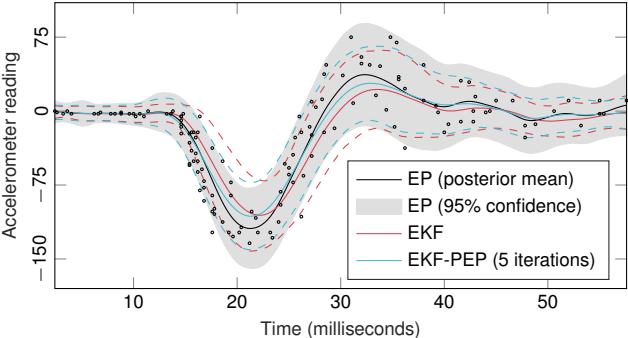
- For sequential data, the EKF is equivalent to single-sweep EP where the moment matching integral is solved via linearisation.
- Our algorithm iteratively refines the EKF by linearising about the cavity, rather than the filter predictions (prior).
- When $\alpha = 0$ we recover the Posterior Linearisation Filter. Lack of cavity calculation makes it very stable.
- Linearisation allows for straightforward calculation of cross-covariances vs. quadrature methods.

Example: log-Gaussian Cox process



The coal mining accident task (log-Gaussian Cox process) is well approximated by local linearisations, and iteration improves the match to the EP posterior.

Example: heteroscedastic noise



Linearisation in the motorcycle crash task (heteroscedastic noise) is a crude approximation, but iterating still improves the posterior.

Conclusions

- Spectral mixture GPs can be written as SDEs and have a close connection to probabilistic time-frequency analysis
- We can perform full power EP in the state space GP setting ($\mathcal{O}(n)$)
- Linearisation in state space power EP recovers the EKF and PLF (more connections to be found ...)

Thanks for listening!

contact: **william.wilkinson@aalto.fi**

Bibliography



Hartikainen, J. and Särkkä, S. (2010).

Kalman filtering and smoothing solutions to temporal Gaussian process regression models. In *International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 379–384. IEEE.



Solin, A. and Särkkä, S. (2014).

Explicit link between periodic covariance functions and state space models. In *Artificial Intelligence and Statistics*, pages 904–912.



Turner, R. E. and Sahani, M. (2014).

Time-frequency analysis as probabilistic inference. *IEEE Transactions on Signal Processing*, 62(23):6171.



Wilkinson, W. J., Andersen, M. R., Reiss, J. D., Stowell, D., and Solin, A. (2019a).

End-to-end probabilistic inference for nonstationary audio analysis. In *International Conference on Machine Learning*, pages 6776–6785.



Wilkinson, W. J., Chang, P. E., Andersen, M. R., and Solin, A. (2019b).

Global approximate inference via local linearisation for temporal gaussian processes. In *Second Symposium on Advances in Approximate Bayesian Inference*.



Wilkinson, W. J., Riis Andersen, M., Reiss, J. D., Stowell, D., and Solin, A. (2019c).

Unifying probabilistic models for time-frequency analysis. In *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pages 3352–3356.

contact: **william.wilkinson@aalto.fi**