# Gaussian processes approximations for time series

Dr. Richard E. Turner (ret26@cam.ac.uk)

Computational and Biological Learning Lab, Department of Engineering, University of Cambridge





### Using Gaussian processes for time-series modelling



### Using Gaussian processes for time-series modelling

Generative model (like non-linear regression)  $y(t) = f(t) + \epsilon \sigma_y$  (independent Gaussian noise)  $p(\epsilon) = \mathcal{N}(0, 1)$ 

place GP prior over the non-linear function

 $p(f(\mathbf{t})|\theta) = \mathcal{GP}(0, \mathbf{K}(\mathbf{t}, \mathbf{t}'))$ 

$$\mathsf{K}(\mathsf{t},\mathsf{t}') = \sigma^2 \cos(\omega(\mathsf{t}-\mathsf{t}')) \exp\left(-\frac{1}{2l^2}(\mathsf{t}-\mathsf{t}')^2\right)$$



typical time-series covariance sinusoids with SE envelopes power in Gaussian subband

Generative model (like non-linear regression)  $y(t) = f(t) + \epsilon \sigma_y$  (independent Gaussian noise)  $p(\epsilon) = \mathcal{N}(0, 1)$ 

place GP prior over the non-linear function

 $p(f(\mathbf{t})|\theta) = \mathcal{GP}(0, \mathbf{K}(\mathbf{t}, \mathbf{t}'))$ 

$$\mathsf{K}(\mathsf{t},\mathsf{t}') = \sigma^2 \cos(\omega(\mathsf{t}-\mathsf{t}')) \exp\left(-\frac{1}{2l^2}(\mathsf{t}-\mathsf{t}')^2\right)$$



typical time-series covariance sinusoids with SE envelopes power in Gaussian subband

another popular class of GP time-series models  $\begin{aligned} \mathbf{x}_t &= g(\mathbf{x}_{t-1}) + \sigma_{\mathbf{x}} \eta_t \\ \mathbf{y}_t &= f(\mathbf{x}_t) + \sigma_{\mathbf{y}} \epsilon_t \\ f(\mathbf{x}), g(\mathbf{x}) \sim \mathcal{GP}(0, \mathsf{K}) \end{aligned}$ 

Generative model (like non-linear regression)  $y(t) = f(t) + \epsilon \sigma_y$  (independent Gaussian noise)  $p(\epsilon) = \mathcal{N}(0, 1)$ 

place GP prior over the non-linear function

$$p(f(t)|\theta) = \mathcal{GP}(0, \mathsf{K}(t, t'))$$
$$\mathsf{K}(t, t') = \sigma^2 \cos(\omega(t - t')) \exp\left(-\frac{1}{2l^2}(t - t')^2\right)$$

$$(1) \quad (1) \quad (1) \quad (1) \quad (1) \quad (2) \quad (2)$$

sum of two Gaussians is a Gaussian,  $\implies$  induces GP over y(t)

$$p(\mathbf{y}(\mathbf{t})|\theta) = \mathcal{GP}(0, \mathbf{K}(\mathbf{t}, \mathbf{t}') + \mathbf{I}\sigma_{\mathbf{y}}^{2})$$
$$\Sigma(\mathbf{t}, \mathbf{t}')$$



typical time-series covariance sinusoids with SE envelopes power in Gaussian subband

another popular class of GP time-series models  $x_{t} = g(x_{t-1}) + \sigma_{x}\eta_{t}$  $y_{t} = f(x_{t}) + \sigma_{y}\epsilon_{t}$  $f(x), g(x) \sim \mathcal{GP}(0, \mathsf{K})$ non-Gaussian distribution on y(t)

Generative model (like non-linear regression)  $y(t) = f(t) + \epsilon \sigma_y$  (independent Gaussian noise)  $p(\epsilon) = \mathcal{N}(0, 1)$ 

place GP prior over the non-linear function

$$p(f(\mathbf{t})|\theta) = \mathcal{GP}(0, \mathbf{K}(\mathbf{t}, \mathbf{t}'))$$
$$\mathbf{K}(\mathbf{t}, \mathbf{t}') = \sigma^2 \cos(\omega(\mathbf{t} - \mathbf{t}')) \exp\left(-\frac{1}{2l^2}(\mathbf{t} - \mathbf{t}')^2\right)$$

sum of two Gaussians is a Gaussian,  $\implies$  induces GP over y(t)

$$p(\mathbf{y}(\mathbf{t})|\theta) = \mathcal{GP}(0, \mathbf{K}(\mathbf{t}, \mathbf{t}') + \mathbf{I}\sigma_{\mathbf{y}}^{2})$$
$$\Sigma(\mathbf{t}, \mathbf{t}')$$

How do we make predictions? How do we learn hyper-parameters?



typical time-series covariance sinusoids with SE envelopes power in Gaussian subband

another popular class of GP time-series models  $x_{t} = g(x_{t-1}) + \sigma_{x}\eta_{t}$  $y_{t} = f(x_{t}) + \sigma_{y}\epsilon_{t}$  $f(x), g(x) \sim \mathcal{GP}(0, \mathsf{K})$ non-Gaussian distribution on y(t)

How do we make predictions?

$$p(\mathbf{y}_{1}, \mathbf{y}_{2}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y}_{1} \\ \mathbf{y}_{2} \end{bmatrix}; \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$
$$p(\mathbf{y}_{1}|\mathbf{y}_{2}) = \frac{p(\mathbf{y}_{1}, \mathbf{y}_{2})}{p(\mathbf{y}_{2})} \qquad p(\mathbf{y}_{2}) = \mathcal{N}\left(\mathbf{y}_{2}; \mathbf{0}, \Sigma_{22}\right)$$



 $\implies p(\mathbf{y}_1|\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1; \Sigma_{12}\Sigma_{22}^{-1}\mathbf{y}_2, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^{\mathsf{T}})$ 

How do we make predictions?

$$p(\mathbf{y}_{1}, \mathbf{y}_{2}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y}_{1} \\ \mathbf{y}_{2} \end{bmatrix}; \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$
$$p(\mathbf{y}_{1}|\mathbf{y}_{2}) = \frac{p(\mathbf{y}_{1}, \mathbf{y}_{2})}{p(\mathbf{y}_{2})} \longrightarrow p(\mathbf{y}_{2}) = \mathcal{N}\left(\mathbf{y}_{2}; \mathbf{0}, \Sigma_{22}\right)$$

У

$$\implies p(\mathbf{y}_1|\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1; \Sigma_{12}\Sigma_{22}^{-1}\mathbf{y}_2, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^{\mathsf{T}})$$

How do we learn hyper-parameters?

 $p(\theta|\mathbf{y}_{1:N}) = \frac{p(\mathbf{y}_{1:N}|\theta)p(\theta)}{p(\mathbf{y}_{1:N})} \quad \text{(Bayes' Rule)}$ 

 $p(\mathbf{y}_{1:N}|\boldsymbol{\theta}) = \mbox{ likelihood of the parameters }$ 

= how well did  $\theta$  predict the data we observed

$$p(\mathbf{y}_{1:N}|\theta) = \frac{1}{\det(2\pi\Sigma(\theta))^{-1/2}} \exp\left(-\frac{1}{2}\mathbf{y}_{1:N}^{\mathsf{T}}\Sigma^{-1}(\theta)\mathbf{y}_{1:N}\right)$$

У

X

require matrix inversion

(Cholesky)

How do we make predictions?

$$p(\mathbf{y}_1, \mathbf{y}_2) = \mathcal{N}\left(\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix}; \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$
$$p(\mathbf{y}_1 | \mathbf{y}_2) = \frac{p(\mathbf{y}_1, \mathbf{y}_2)}{p(\mathbf{y}_2)} \qquad p(\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_2; \mathbf{0}, \Sigma_{22})$$

$$\implies p(\mathbf{y}_1|\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1; \Sigma_{12}\Sigma_{22}^{-1}\mathbf{y}_2, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^{\mathsf{T}})$$

How do we learn hyper-parameters?

$$\begin{split} p(\theta|\mathbf{y}_{1:N}) &= \frac{p(\mathbf{y}_{1:N}|\theta)p(\theta)}{p(\mathbf{y}_{1:N})} \quad \text{(Bayes' Rule)} \quad \Rightarrow \text{O(1000) datapoints} \\ p(\mathbf{y}_{1:N}|\theta) &= \text{likelihood of the parameters} \\ &= \text{how well did } \theta \text{ predict the data we observed} \\ p(\mathbf{y}_{1:N}|\theta) &= \frac{1}{\det(2\pi\Sigma(\theta))^{-1/2}} \exp\left(-\frac{1}{2}\mathbf{y}_{1:N}^{\mathsf{T}}\Sigma^{-1}(\theta)\mathbf{y}_{1:N}\right) \end{split}$$

DFT: Turner."Statistical models for natural sounds"

DTC: Seeger et al. "Fast Forward Selection to Speed Up Sparse Gaussian Process Regression"

EP: Qi et al. "Sparse-posterior Gaussian Processes for general likelihoods."

ET: Sudderth et al. "Embedded Trees: Estimation of Gaussian Processes on Graphs with Cycles"

FITC: Snelson et al. "Sparse Gaussian Processes using Pseudo-inputs"

ISL: Gibbs et al. "Efficient implementation of Gaussian processes"

PITC: Snelson et al. "Local and global sparse Gaussian process approximations"

SDE: Sarkka et al. "Spatiotemporal Learning via Infinite Dimensional Bayesian Filtering and Smoothing"

SS: Lazaro-Gredilla et al. "Sparse spectrum Gaussian process regression".



DFT: Turner."Statistical models for natural sounds"

DTC: Seeger et al. "Fast Forward Selection to Speed Up Sparse Gaussian Process Regression"

EP: Qi et al. "Sparse-posterior Gaussian Processes for general likelihoods."

ET: Sudderth et al. "Embedded Trees: Estimation of Gaussian Processes on Graphs with Cycles"

FITC: Snelson et al. "Sparse Gaussian Processes using Pseudo-inputs"

ISL: Gibbs et al. "Efficient implementation of Gaussian processes"

PITC: Snelson et al. "Local and global sparse Gaussian process approximations"

SDE: Sarkka et al. "Spatiotemporal Learning via Infinite Dimensional Bayesian Filtering and Smoothing"

SS: Lazaro-Gredilla et al. "Sparse spectrum Gaussian process regression".



DFT: Turner."Statistical models for natural sounds"

DTC: Seeger et al. "Fast Forward Selection to Speed Up Sparse Gaussian Process Regression"

EP: Qi et al. "Sparse-posterior Gaussian Processes for general likelihoods."

ET: Sudderth et al. "Embedded Trees: Estimation of Gaussian Processes on Graphs with Cycles"

FITC: Snelson et al. "Sparse Gaussian Processes using Pseudo-inputs"

ISL: Gibbs et al. "Efficient implementation of Gaussian processes"

PITC: Snelson et al. "Local and global sparse Gaussian process approximations"

SDE: Sarkka et al. "Spatiotemporal Learning via Infinite Dimensional Bayesian Filtering and Smoothing"

SS: Lazaro-Gredilla et al. "Sparse spectrum Gaussian process regression".



DFT: Turner."Statistical models for natural sounds"

DTC: Seeger et al. "Fast Forward Selection to Speed Up Sparse Gaussian Process Regression"

EP: Qi et al. "Sparse-posterior Gaussian Processes for general likelihoods."

ET: Sudderth et al. "Embedded Trees: Estimation of Gaussian Processes on Graphs with Cycles"

FITC: Snelson et al. "Sparse Gaussian Processes using Pseudo-inputs"

ISL: Gibbs et al. "Efficient implementation of Gaussian processes"

PITC: Snelson et al. "Local and global sparse Gaussian process approximations"

SDE: Sarkka et al. "Spatiotemporal Learning via Infinite Dimensional Bayesian Filtering and Smoothing"

SS: Lazaro-Gredilla et al. "Sparse spectrum Gaussian process regression".







1. augment model with M<T pseudo data

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathsf{K}_{\mathbf{ff}} & \mathsf{K}_{\mathbf{fu}} \\ \mathsf{K}_{\mathbf{uf}} & \mathsf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$



1. augment model with M<T pseudo data

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

2. remove some of the dependencies

(results in simpler model)



all factors

1. augment model with M<T pseudo data

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathsf{K}_{\mathbf{ff}} & \mathsf{K}_{\mathbf{fu}} \\ \mathsf{K}_{\mathbf{uf}} & \mathsf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

2. remove some of the dependencies

(results in simpler model)



 $(f_i) \bullet (f_j) \longrightarrow (f_i) (f_j)$  all factors

1. augment model with M<T pseudo data

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

2. remove some of the dependencies (results in simpler model)



(e.g. using KL divergence, many choices)

$$\underset{q(\mathbf{u}),\{q(\mathbf{f}_t|\mathbf{u})\}_{t=1}^T}{\arg\min} \mathsf{KL}(p(\mathbf{f},\mathbf{u})||q(\mathbf{u})\prod_{t=1}^T q(\mathbf{f}_t|\mathbf{u})) \implies \frac{q(\mathbf{u}) = p(\mathbf{u})}{q(\mathbf{f}_t|\mathbf{u}) = p(\mathbf{f}_t|\mathbf{u})}$$

equal to exact conditionals

1. augment model with M<T pseudo data

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

2. remove some of the dependencies (results in simpler model)



(e.g. using KL divergence, many choices)

 $\underset{q(\mathbf{u}),\{q(\mathbf{f}_t|\mathbf{u})\}_{t=1}^T}{\operatorname{arg\,min}} \operatorname{KL}(p(\mathbf{f},\mathbf{u})||q(\mathbf{u})\prod_{t=1}^T q(\mathbf{f}_t|\mathbf{u})) \implies \begin{array}{l} q(\mathbf{u}) = p(\mathbf{u}) \\ q(\mathbf{f}_t|\mathbf{u}) = p(\mathbf{f}_t|\mathbf{u}) \end{array}$ 

equal to exact conditionals



 $q(\mathbf{u}) = p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; 0, \mathsf{K}_{\mathsf{uu}})$ 



 $\begin{aligned} q(\mathbf{u}) &= p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; 0, \mathsf{K}_{\mathsf{uu}}) \\ q(\mathsf{f}_t | \mathbf{u}) &= p(\mathsf{f}_t | \mathbf{u}) \end{aligned}$ 



3

$$q(\mathbf{u}) = p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; 0, \mathsf{K}_{\mathsf{uu}})$$
$$q(\mathsf{f}_t | \mathbf{u}) = p(\mathsf{f}_t | \mathbf{u})$$





$$q(\mathbf{u}) = p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; 0, \mathsf{K}_{\mathsf{uu}})$$
$$q(\mathsf{f}_t | \mathbf{u}) = p(\mathsf{f}_t | \mathbf{u})$$
$$= \mathcal{N}(\mathsf{f}_t; \mathsf{K}_{\mathsf{f}_t \mathsf{u}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathbf{u}, \mathsf{K}_{\mathsf{f}_t \mathsf{f}_t} - \mathsf{K}_{\mathsf{f}_t \mathsf{u}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathsf{K}_{\mathsf{uf}_t})$$



$$q(\mathbf{u}) = p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; 0, \mathsf{K}_{uu})$$

$$q(\mathbf{f}_t | \mathbf{u}) = p(\mathbf{f}_t | \mathbf{u})$$

$$= \mathcal{N}(\mathbf{f}_t; \mathsf{K}_{f_t u} \mathsf{K}_{uu}^{-1} \mathbf{u}, \mathsf{K}_{f_t f_t} - \mathsf{K}_{f_t u} \mathsf{K}_{uu}^{-1} \mathsf{K}_{uf_t})$$

$$D_{tt}$$

$$(\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3)$$

$$(\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3)$$

$$q(\mathbf{u}) = p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; 0, \mathsf{K}_{uu})$$

$$q(\mathbf{f}_t | \mathbf{u}) = p(\mathbf{f}_t | \mathbf{u})$$

$$= \mathcal{N}(\mathbf{f}_t; \mathsf{K}_{\mathbf{f}_t u} \mathsf{K}_{uu}^{-1} \mathbf{u}, \mathsf{K}_{\mathbf{f}_t \mathbf{f}_t} - \mathsf{K}_{\mathbf{f}_t u} \mathsf{K}_{uu}^{-1} \mathsf{K}_{uf_t})$$

$$p(\mathbf{y}_t | \mathbf{f}_t) = p(\mathbf{y}_t | \mathbf{f}_t) = \mathcal{N}(\mathbf{y}_t; \mathbf{f}_t, \sigma_{\mathbf{y}}^2)$$

$$q(\mathbf{y}_t | \mathbf{f}_t) = p(\mathbf{y}_t | \mathbf{f}_t) = \mathcal{N}(\mathbf{y}_t; \mathbf{f}_t, \sigma_{\mathbf{y}}^2)$$

$$(\mathbf{y}_1 | \mathbf{y}_2 | \mathbf{y}_3)$$

$$(\mathbf{y}_1 | \mathbf{y}_2 | \mathbf{y}_3)$$

$$q(\mathbf{u}) = p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; 0, \mathsf{K}_{\mathsf{uu}})$$

$$q(\mathbf{f}_t | \mathbf{u}) = p(\mathbf{f}_t | \mathbf{u})$$

$$= \mathcal{N}(\mathbf{f}_t; \mathsf{K}_{\mathsf{f}_t \mathsf{u}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathbf{u}, \mathsf{K}_{\mathsf{f}_t \mathsf{f}_t} - \mathsf{K}_{\mathsf{f}_t \mathsf{u}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathsf{K}_{\mathsf{uf}_t})$$

$$\mathbf{D}_{tt}$$

$$q(\mathbf{y}_t | \mathbf{f}_t) = p(\mathbf{y}_t | \mathbf{f}_t) = \mathcal{N}(\mathbf{y}_t; \mathbf{f}_t, \sigma_{\mathbf{y}}^2)$$

cost of computing likelihood is  $O(TM^2)$ 





$$q(\mathbf{u}) = p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; 0, \mathsf{K}_{\mathsf{uu}})$$

$$q(\mathbf{f}_t | \mathbf{u}) = p(\mathbf{f}_t | \mathbf{u})$$

$$= \mathcal{N}(\mathbf{f}_t; \mathsf{K}_{\mathsf{f}_t \mathsf{u}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathbf{u}, \mathsf{K}_{\mathsf{f}_t \mathsf{f}_t} - \mathsf{K}_{\mathsf{f}_t \mathsf{u}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathsf{K}_{\mathsf{u}})$$

$$D_{tt}$$

$$q(\mathbf{y}_t | \mathbf{f}_t) = p(\mathbf{y}_t | \mathbf{f}_t) = \mathcal{N}(\mathbf{y}_t; \mathbf{f}_t, \sigma_{\mathbf{y}}^2)$$

cost of computing likelihood is  $O(TM^2)$ 

$$p(\mathbf{y}_t|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}; \mathbf{0}, \mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathsf{K}_{\mathsf{uu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathsf{K}_{\mathsf{uf}} + \mathsf{D} + \sigma_{\mathsf{y}}^2\mathsf{I})$$





construct new generative model (with pseudo-data) cheaper to perform exact learning and inference calibrated to original

indirect posterior approximation

$$\begin{split} q(\mathbf{u}) &= p(\mathbf{u}) = \mathcal{N}(\mathbf{u}; 0, \mathsf{K}_{\mathsf{uu}}) \\ q(\mathsf{f}_t | \mathbf{u}) &= p(\mathsf{f}_t | \mathbf{u}) \\ &= \mathcal{N}(\mathsf{f}_t; \mathsf{K}_{\mathsf{f}_t \mathsf{u}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathbf{u}, \mathsf{K}_{\mathsf{f}_t \mathsf{f}_t} - \mathsf{K}_{\mathsf{f}_t \mathsf{u}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathsf{K}_{\mathsf{uf}_t}) \\ &= \mathcal{N}(\mathsf{f}_t; \mathsf{K}_{\mathsf{f}_t \mathsf{u}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathbf{u}, \mathsf{K}_{\mathsf{f}_t \mathsf{f}_t} - \mathsf{K}_{\mathsf{f}_t \mathsf{u}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathsf{K}_{\mathsf{uf}_t}) \\ &= p(\mathsf{y}_t | \mathsf{f}_t) = p(\mathsf{y}_t | \mathsf{f}_t) = \mathcal{N}(\mathsf{y}_t; \mathsf{f}_t, \sigma_{\mathsf{y}}^2) \\ &\text{cost of computing likelihood is } \mathcal{O}(TM^2) \\ &= \mathcal{N}(\mathsf{y}; \mathbf{0}, \mathsf{K}_{\mathsf{fu}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathsf{K}_{\mathsf{uu}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathsf{K}_{\mathsf{uf}} + \mathsf{D} + \sigma_{\mathsf{y}}^2 \mathsf{I}) \\ &= \mathcal{N}(\mathsf{y}; \mathbf{0}, \mathsf{K}_{\mathsf{fu}} \mathsf{K}_{\mathsf{uu}}^{-1} \mathsf{K}_{\mathsf{uf}} + \mathsf{D} + \sigma_{\mathsf{y}}^2 \mathsf{I}) \end{split}$$





construct new generative model (with pseudo-data) cheaper to perform exact learning and inference calibrated to original

indirect posterior approximation



construct new generative model (with pseudo-data) cheaper to perform exact learning and inference calibrated to original

indirect posterior approximation



calibrated to original

posterior approximation

- parametric (although cleverly so)
- if I see more data, should I add extra pseudo-data?
## Fully independent training conditional (FITC) approximation

- parametric (although cleverly so)
- if I see more data, should I add extra pseudo-data?
  - unnatural from a generative modelling perspective
  - natural from a prediction perspective (need greater complexity/bandwidth)
- $\implies$  lost elegant separation of model, inference and approximation

## Fully independent training conditional (FITC) approximation

- parametric (although cleverly so)
- if I see more data, should I add extra pseudo-data?
  - unnatural from a generative modelling perspective
  - natural from a prediction perspective (need greater complexity/bandwidth)
- $\implies$  lost elegant separation of model, inference and approximation
- example of prior approximation

## Fully independent training conditional (FITC) approximation

- parametric (although cleverly so)
- if I see more data, should I add extra pseudo-data?
  - unnatural from a generative modelling perspective
  - natural from a prediction perspective (need greater complexity/bandwidth)
- $\implies$  lost elegant separation of model, inference and approximation
- example of prior approximation

#### **Extensions:**

- methods for optimising pseudo-inputs (indirect approximations tend to over-fit)
- partially independent training conditional...



construct new generative model (with pseudo-data) cheaper to perform exact learning and inference calibrated to original ap

indirect posterior approximation

1. augment model with M<T pseudo data

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathsf{K}_{\mathbf{f}\mathbf{f}} & \mathsf{K}_{\mathbf{f}\mathbf{u}} \\ \mathsf{K}_{\mathbf{u}\mathbf{f}} & \mathsf{K}_{\mathbf{u}\mathbf{u}} \end{bmatrix} \right)$$



construct new generative model (with pseudo-data)indirectcheaper to perform exact learning and inferenceposteriorcalibrated to originalapproximation

1. augment model with M<T pseudo data

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \left[ \begin{array}{c} \mathbf{f} \\ \mathbf{u} \end{array} \right]; \left[ \begin{array}{c} 0 \\ 0 \end{array} \right], \left[ \begin{array}{c} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{array} \right] \right)$$

2. remove some of the dependencies (results in simpler model)



 $(f_j)$  between blocks

construct new generative model (with pseudo-data)indirectcheaper to perform exact learning and inferenceposteriorcalibrated to originalapproximation

1. augment model with M<T pseudo data

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

2. remove some of the dependencies (results in simpler model)



 $\mathbf{f}_1 = \begin{bmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \end{bmatrix} \quad \mathbf{f}_2 = \mathbf{f}_3$ 

construct new generative model (with pseudo-data)indirectcheaper to perform exact learning and inferenceposteriorcalibrated to originalapproximation

1. augment model with M<T pseudo data  $p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathsf{K}_{\mathbf{ff}} & \mathsf{K}_{\mathbf{fu}} \\ \mathsf{K}_{\mathbf{uf}} & \mathsf{K}_{\mathbf{uu}} \end{bmatrix}\right)$ 

2. remove some of the dependencies (results in simpler model)

$$(f_i) \bullet (f_j) \longrightarrow (f_i) (f_j)$$
 between blocks

3. calibrate model

(e.g. using KL divergence, many choices)

$$\underset{q(\mathbf{u}),\{q(\mathbf{f}_{k}|\mathbf{u})\}_{k=1}^{K}}{\arg\min} \operatorname{KL}(p(\mathbf{f},\mathbf{u})||q(\mathbf{u})\prod_{k=1}^{K}q(\mathbf{f}_{k}|\mathbf{u})) \Longrightarrow \frac{q(\mathbf{u})=p(\mathbf{u})}{q(\mathbf{f}_{k}|\mathbf{u})=p(\mathbf{f}_{k}|\mathbf{u})}$$

equal to exact conditionals

construct new generative model (with pseudo-data)indirectcheaper to perform exact learning and inferenceposteriorcalibrated to originalapproximation

 $\begin{array}{c} u_1 & u_2 \\ \hline f_1 & f_2 \\ k = 1 \\ K = 2 \\ blocks \end{array}$ 

 $\mathbf{f}_1 = \begin{bmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \end{bmatrix} \quad \mathbf{f}_2 = \mathbf{f}_3$ 

- minimise variational KL between two posterior distributions
  - direct posterior approximation
  - likelihood approximation

$$\underset{q(\mathbf{y}|\mathbf{u})}{\arg\min} \mathsf{KL}\left(\frac{1}{Z}p(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(\mathbf{y}|\mathbf{u})||\frac{1}{Z}p(\mathbf{u})p(\mathbf{f}|\mathbf{u})p(\mathbf{y}|\mathbf{f})\right) \text{ such that } \int \mathsf{d}\mathbf{y} \ q(\mathbf{y}|\mathbf{u}) = 1$$

- minimise variational KL between two posterior distributions
  - direct posterior approximation
  - likelihood approximation

$$\underset{q(\mathbf{y}|\mathbf{u})}{\arg\min} \mathsf{KL}\left(\frac{1}{Z}p(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(\mathbf{y}|\mathbf{u})||\frac{1}{Z}p(\mathbf{u})p(\mathbf{f}|\mathbf{u})p(\mathbf{y}|\mathbf{f})\right) \text{ such that } \int d\mathbf{y} \ q(\mathbf{y}|\mathbf{u}) = 1$$

$$\implies q(\mathbf{y}|\mathbf{u}) = \mathcal{N}(\mathbf{y}; \mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathbf{u}, \sigma_{\mathsf{y}}^{2}\mathsf{I})$$

- minimise variational KL between two posterior distributions
  - direct posterior approximation
  - likelihood approximation

$$\underset{q(\mathbf{y}|\mathbf{u})}{\operatorname{arg\,min}} \operatorname{KL}\left(\frac{1}{Z}p(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(\mathbf{y}|\mathbf{u})||\frac{1}{Z}p(\mathbf{u})p(\mathbf{f}|\mathbf{u})p(\mathbf{y}|\mathbf{f})\right) \text{ such that } \int d\mathbf{y} \ q(\mathbf{y}|\mathbf{u}) = 1$$

$$\implies q(\mathbf{y}|\mathbf{u}) = \mathcal{N}(\mathbf{y}; \mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathbf{u}, \sigma_{\mathsf{y}}^{2}\mathsf{I})$$

can this approximation be categorised in another way?

- minimise variational KL between two posterior distributions
  - direct posterior approximation
  - likelihood approximation

$$\underset{q(\mathbf{y}|\mathbf{u})}{\arg\min} \mathsf{KL}\left(\frac{1}{Z}p(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(\mathbf{y}|\mathbf{u})||\frac{1}{Z}p(\mathbf{u})p(\mathbf{f}|\mathbf{u})p(\mathbf{y}|\mathbf{f})\right) \text{ such that } \int \mathsf{d}\mathbf{y} \ q(\mathbf{y}|\mathbf{u}) = 1$$

$$\implies q(\mathbf{y}|\mathbf{u}) = \mathcal{N}(\mathbf{y}; \mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathbf{u}, \sigma_{\mathsf{y}}^{2}\mathsf{I})$$

• probabilistic PCA to FITC's factor analysis

$$\begin{aligned} \mathsf{DTC}: \ p(\mathbf{y}|\theta) &= \mathcal{N}(\mathbf{y};\mathbf{0},\mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathsf{K}_{\mathsf{uf}} + \sigma_{\mathsf{y}}^{2}\mathsf{I}) \\ \mathsf{FITC}: \ p(\mathbf{y}|\theta) &= \mathcal{N}(\mathbf{y};\mathbf{0},\mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathsf{K}_{\mathsf{uf}} + \mathsf{D} + \sigma_{\mathsf{y}}^{2}\mathsf{I}) \end{aligned}$$

- minimise variational KL between two posterior distributions
  - direct posterior approximation
  - likelihood approximation

$$\underset{q(\mathbf{y}|\mathbf{u})}{\arg\min} \mathsf{KL}\left(\frac{1}{Z}p(\mathbf{u})p(\mathbf{f}|\mathbf{u})q(\mathbf{y}|\mathbf{u})||\frac{1}{Z}p(\mathbf{u})p(\mathbf{f}|\mathbf{u})p(\mathbf{y}|\mathbf{f})\right) \text{ such that } \int d\mathbf{y} \ q(\mathbf{y}|\mathbf{u}) = 1$$

$$\implies q(\mathbf{y}|\mathbf{u}) = \mathcal{N}(\mathbf{y}; \mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathbf{u}, \sigma_{\mathsf{y}}^{2}\mathsf{I})$$

• probabilistic PCA to FITC's factor analysis

$$\begin{aligned} \mathsf{DTC}: \ p(\mathbf{y}|\theta) &= \mathcal{N}(\mathbf{y};\mathbf{0},\mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathsf{K}_{\mathsf{uf}} + \sigma_{\mathsf{y}}^{2}\mathsf{I}) \\ \mathsf{FITC}: \ p(\mathbf{y}|\theta) &= \mathcal{N}(\mathbf{y};\mathbf{0},\mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathsf{K}_{\mathsf{uf}} + \mathsf{D} + \sigma_{\mathsf{y}}^{2}\mathsf{I}) \end{aligned}$$

• blurred division between direct/indirect and likelihood/prior approximation

augment the model with pseudo-data:  $p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, \mathbf{u})$ 

augment the model with pseudo-data:  $p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, \mathbf{u})$ 

lower bound the likelihood

 $\mathcal{L}(\theta) = \log p(\mathbf{y}|\theta)$ 





augment the model with pseudo-data:  $p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, \mathbf{u})$ 

lower bound the likelihood

$$\mathcal{L}(\theta) = \log p(\mathbf{y}|\theta) = \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u})$$



augment the model with pseudo-data:  $p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, \mathbf{u})$ lower bound the likelihood

 $\mathcal{L}(\theta) = \log p(\mathbf{y}|\theta) = \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u})$ 

$$= \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u}) \frac{q(\mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})}$$







augment the model with pseudo-data:  $p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, \mathbf{u})$ lower bound the likelihood  $\mathcal{L}(\theta) = \log p(\mathbf{y}|\theta) = \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u})$  $= \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u}) \frac{q(\mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \ge \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} = \mathcal{F}(q, \theta)$ 

assume approximate posterior factorisation with special form

 $q(\mathbf{f}, \mathbf{u}) = q(\mathbf{f}|\mathbf{u})q(\mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u}) \qquad \text{(exact } q(\mathbf{f}|\mathbf{u}) = p(\mathbf{f}|\mathbf{y})\text{)}$ 

augment the model with pseudo-data:  $p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, \mathbf{u})$ lower bound the likelihood  $\mathcal{L}(\theta) = \log p(\mathbf{y}|\theta) = \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u})$  $= \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u}) \frac{q(\mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \ge \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} = \mathcal{F}(q, \theta)$ 

assume approximate posterior factorisation with special form

$$q(\mathbf{f}, \mathbf{u}) = q(\mathbf{f}|\mathbf{u})q(\mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u}) \quad (\text{exact } q(\mathbf{f}|\mathbf{u}) = p(\mathbf{f}|\mathbf{y}))$$
$$\mathcal{F}(q, \theta) = \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u})}{p(\mathbf{f}|\mathbf{u})q(\mathbf{u})} = \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{f}|\mathbf{u})q(\mathbf{u})}$$

augment the model with pseudo-data:  $p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, \mathbf{u})$ lower bound the likelihood  $\mathcal{L}(\theta) = \log p(\mathbf{y}|\theta) = \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u})$  $= \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u}) \frac{q(\mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \ge \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} = \mathcal{F}(q, \theta)$ 

assume approximate posterior factorisation with special form

$$q(\mathbf{f}, \mathbf{u}) = q(\mathbf{f}|\mathbf{u})q(\mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u}) \quad (\text{exact } q(\mathbf{f}|\mathbf{u}) = p(\mathbf{f}|\mathbf{y}))$$
$$\mathcal{F}(q, \theta) = \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u})}{p(\mathbf{f}|\mathbf{u})q(\mathbf{u})} = \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}|\mathbf{f})\frac{p(\mathbf{f}|\mathbf{u})}{p(\mathbf{f}|\mathbf{u})}q(\mathbf{u})}{\frac{p(\mathbf{f}|\mathbf{u})}{q(\mathbf{u})}}$$

augment the model with pseudo-data:  $p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, \mathbf{u})$ lower bound the likelihood  $\mathcal{L}(\theta) = \log p(\mathbf{y}|\theta) = \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u})$  $= \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u}) \frac{q(\mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \ge \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} = \mathcal{F}(q, \theta)$ 

assume approximate posterior factorisation with special form

$$q(\mathbf{f}, \mathbf{u}) = q(\mathbf{f}|\mathbf{u})q(\mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u}) \quad (\text{exact } q(\mathbf{f}|\mathbf{u}) = p(\mathbf{f}|\mathbf{y}))$$
$$\mathcal{F}(q, \theta) = \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u})}{p(\mathbf{f}|\mathbf{u})q(\mathbf{u})} = \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}|\mathbf{f})\frac{p(\mathbf{f}|\mathbf{u})}{p(\mathbf{f}|\mathbf{u})}q(\mathbf{u})}{\frac{p(\mathbf{f}|\mathbf{u})}{q(\mathbf{u})}q(\mathbf{u})}$$
make bound as tight as possible:  $q^*(\mathbf{u}) = \arg \max \mathcal{F}(q, \theta)$ 

 $q(\mathbf{u})$ 

augment the model with pseudo-data:  $p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, \mathbf{u})$ lower bound the likelihood  $\mathcal{L}(\theta) = \log p(\mathbf{y}|\theta) = \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u})$  $= \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u}) \frac{q(\mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \ge \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} = \mathcal{F}(q, \theta)$ 

assume approximate posterior factorisation with special form

$$q(\mathbf{f}, \mathbf{u}) = q(\mathbf{f}|\mathbf{u})q(\mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u}) \quad (\text{exact } q(\mathbf{f}|\mathbf{u}) = p(\mathbf{f}|\mathbf{y}))$$
$$\mathcal{F}(q, \theta) = \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u})}{p(\mathbf{f}|\mathbf{u})q(\mathbf{u})} = \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}|\mathbf{f})\frac{p(\mathbf{f}|\mathbf{u})}{p(\mathbf{f}|\mathbf{u})}q(\mathbf{u})}{\frac{p(\mathbf{f}|\mathbf{u})}{q(\mathbf{u})}}$$
make bound as tight as possible:  $q^*(\mathbf{u}) = \underset{q(\mathbf{u})}{\arg \max} \mathcal{F}(q, \theta)$ 

 $q^*(\mathbf{u}) \propto p(\mathbf{u}) \mathcal{N}(\mathbf{y}; \mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathbf{u}, \sigma_{\mathsf{y}}^2\mathsf{I})$  (DTC)

augment the model with pseudo-data:  $p(\mathbf{y}, \mathbf{f}, \mathbf{u}|\theta) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}, \mathbf{u})$ lower bound the likelihood  $\mathcal{L}(\theta) = \log p(\mathbf{y}|\theta) = \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u})$  $\theta$  $= \log \int d\mathbf{f} \, d\mathbf{u} \, p(\mathbf{y}, \mathbf{f}, \mathbf{u}) \frac{q(\mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \geq \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} = \mathcal{F}(q, \theta)$ assume approximate posterior factorisation with special form

 $q(\mathbf{f}, \mathbf{u}) = q(\mathbf{f}|\mathbf{u})q(\mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u})$  (exact  $q(\mathbf{f}|\mathbf{u}) = p(\mathbf{f}|\mathbf{y})$ )  $\mathcal{F}(q,\theta) = \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f},\mathbf{u}) \log \frac{p(\mathbf{y},\mathbf{f},\mathbf{u})}{p(\mathbf{f}|\mathbf{u})q(\mathbf{u})} = \int d\mathbf{f} \, d\mathbf{u} \, q(\mathbf{f},\mathbf{u}) \log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})}{p(\mathbf{f}|\mathbf{u})}q(\mathbf{u})$ make bound as tight as possible:  $q^*(\mathbf{u}) = \arg \max \mathcal{F}(q, \theta)$  $q(\mathbf{u})$  $q^*(\mathbf{u}) \propto p(\mathbf{u}) \mathcal{N}(\mathbf{y}; \mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathbf{u}, \sigma_{\mathsf{v}}^2\mathsf{I})$  (DTC)  $\mathcal{F}(q^*, \theta) = \log \mathcal{N}(\mathbf{y}; \mathbf{0}, \mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathsf{K}_{\mathsf{uf}}, \sigma_{\mathsf{y}}^2\mathsf{I}) - \frac{1}{2\sigma_{\mathsf{v}}^2}\mathsf{trace}(\mathsf{K}_{\mathsf{ff}} - \mathsf{K}_{\mathsf{fu}}\mathsf{K}_{\mathsf{uu}}^{-1}\mathsf{K}_{\mathsf{uf}})$ **DTC** like

uncertainty based correction

## Summary of DTC/VFE methods

- optimisation pseudo point inputs much less prone to over-fitting in VFE methods (direct posterior approximation)
- variational methods known to **underfit** (and have othe **biases**)
- like FITC, VFE added u to generative model not merely parameters of the approximation
  - approximation will 'waste' some representational power capturing the distribution over the pseudo-data
  - coherrent way of adding pseudo-data in light of more data?
  - perhaps a more pure approach uses pseudo-data exclusively to parameterise the posterior:  $q(\mathbf{f}) = \frac{1}{Z}p(\mathbf{f})p(\tilde{\mathbf{y}}|\mathbf{f})$  (e.g. Qi et al.)
- how do these methods perform for time-series data-sets?













- effect of each pseudo-datapoint is local
- but computations involving them are **global**
- back of the envelope calculation:
  - M = number pseudo-points, T = data-range
  - -l = range of the (shortest) dependencies in the posterior

$$-M \ge \frac{T}{l}$$

- e.g. audio  $l = 10, T = 10^5 \implies M \ge 10^4$
- $\Rightarrow$  require large numbers of pseudo-points, but methods are  $\mathcal{O}(NM^2)$

#### $\implies$ many applications are out of reach

#### Using the discrete Fourier transform (DFT) to accelerate GPs

**stationary** covariance functions K(t, t') = k(t - t')



#### Using the discrete Fourier transform (DFT) to accelerate GPs

**stationary** covariance functions K(t, t') = k(t - t')

- eigen-functions are sines and cosines



#### Using the discrete Fourier transform (DFT) to accelerate GPs

**stationary** covariance functions K(t, t') = k(t - t')

- eigen-functions are sines and cosines

Bochner's theorem/Wiener Khintchine Theorem



$$\mathbf{k}(\tau) = \int \mathbf{d}\nu \ \gamma(\nu) \exp(2\pi i \tau \nu) - \mathbf{Fourier transform}$$
power spectral density
**stationary** covariance functions K(t, t') = k(t - t')

- eigen-functions are sines and cosines

Bochner's theorem/Wiener Khintchine Theorem

$$\mathbf{k}(\tau) = \int \mathbf{d}\nu \ \gamma(\nu) \exp(2\pi i \tau \nu) \checkmark \mathbf{Fourier transform}$$
  
power spectral density

for regularly sampled data, leads to approximations based on the FFT  $\implies \mathcal{O}(T \log T)$ 

**stationary** covariance functions K(t, t') = k(t - t')

- eigen-functions are sines and cosines

Bochner's theorem/Wiener Khintchine Theorem



$$\mathbf{k}(\tau) = \int \mathbf{d}\nu \ \gamma(\nu) \exp(2\pi i \tau \nu) \longleftarrow \mathbf{Fourier transform}$$
power spectral density

for regularly sampled data, leads to approximations based on the FFT  $\implies \mathcal{O}(T \log T)$   $\overset{T}{\mathsf{K}_{t,t'}} \approx \sum_{k=1}^{T} \mathsf{FT}_{t,k}^{-1} \gamma_{d,k} \mathsf{FT}_{k,t'}$ 

stationary covariance functions K(t, t') = k(t - t')

- eigen-functions are sines and cosines

Bochner's theorem/Wiener Khintchine Theorem



$$\mathbf{k}(\tau) = \int \mathbf{d}\nu \ \gamma(\nu) \exp(2\pi i \tau \nu) - \mathbf{Fourier transform}$$
power spectral density

for regularly sampled data, leads to approximations based on the FFT  $\implies \mathcal{O}(T \log T)$ T

$$\mathsf{K}_{t,t'} \approx \sum_{k=1}^{1} \mathsf{FT}_{t,k}^{-1} \gamma_{d,k} \mathsf{FT}_{k,t'}$$

 $\mathsf{FT}_{k,t} = e^{-2\pi i (k-1)(t-1)/T}$ 

**stationary** covariance functions K(t, t') = k(t - t')

- eigen-functions are sines and cosines

Bochner's theorem/Wiener Khintchine Theorem



$$\mathbf{k}(\tau) = \int \mathbf{d}\nu \ \gamma(\nu) \exp(2\pi i \tau \nu) - \mathbf{Fourier transform}$$
power spectral density

for regularly sampled data, leads to approximations based on the FFT  $\implies \mathcal{O}(T \log T)$   $\overset{T}{\mathsf{K}_{t,t'}} \approx \sum_{t=1}^{T} \mathsf{FT}_{t,k}^{-1} \gamma_{d,k} \mathsf{FT}_{k,t'}$ 

$$\mathsf{FT}_{k,t} = e^{-2\pi i (k-1)(t-1)/T}$$
$$p(\mathbf{y}_{1:T}|\theta) \propto \prod_{k=1}^{T} \gamma_k^{-1/2}(\theta) \exp\left(-\frac{1}{2} \left|\tilde{y}_k\right|^2 / \gamma_k\right)$$

**stationary** covariance functions K(t, t') = k(t - t')

- eigen-functions are sines and cosines

Bochner's theorem/Wiener Khintchine Theorem



$$\mathbf{k}(\tau) = \int \mathbf{d}\nu \ \gamma(\nu) \exp(2\pi i \tau \nu)$$
 Fourier transform power spectral density

for regularly sampled data, leads to approximations based on the FFT

$$\begin{split} & \underset{k=1}{\overset{T}{\leftarrow}} \mathsf{FT}_{t,k}^{-1} \gamma_{d,k} \mathsf{FT}_{k,t'} & \underset{k=1}{\overset{T}{\leftarrow}} \mathsf{FT}_{t,k}^{-1} \gamma_{d,k} \mathsf{FT}_{k,t'} & \underset{f_{11}}{\overset{f_{12}}{\leftarrow}} \mathfrak{f}_{1} \\ & \underset{f_{12}}{\overset{f_{11}}{\leftarrow}} \mathfrak{f}_{1} \\ & \underset{f_{12}}{\overset{f_{12}}{\leftarrow}} \mathfrak{f}_{1} \\ & \underset{f_{12}}{\overset{f_{13}}{\leftarrow}} \mathfrak{f}_{1} \\ & \underset{f_{10}}{\overset{f_{10}}{\leftarrow}} \mathfrak{f}_{4} \\ & \underset{f_{9}}{\overset{f_{10}}{\leftarrow}} \mathfrak{f}_{5} \\ & \underset{f_{9}}{\overset{f_{9}}{\leftarrow}} \mathfrak{f}_{5} \\ & \underset{f_{7}}{\overset{f_{9}}{\leftarrow}} \mathfrak{f}_{6} \\ & \underset{f_{8}}{\overset{f_{7}}{\leftarrow}} \mathfrak{f}_{6} \\ & \underset{f_{7}}{\overset{f_{9}}{\leftarrow}} \mathfrak{f}_{6} \\ & \underset{f_{7}}{\overset{f_{9}}{\leftarrow}} \mathfrak{f}_{6} \\ & \underset{f_{9}}{\overset{f_{9}}{\leftarrow}} \mathfrak{f}_{7} \\ & \underset{f_{8}}{\overset{f_{9}}{\leftarrow}} \mathfrak{f}_{7} \\ & \underset{f_{8}}{\overset{f_{8}}{\leftarrow}} \mathfrak{f}_$$

alternative view: new model with ring topography + exact inference

# Stochastic differential equations (SDE)

- core idea: convert into a classical smoothing problem
- linear, time-invariant stochastic differential equation:

$$\frac{\mathsf{d}^{M}}{\mathsf{d}t^{M}}f(t) + c_{M-1}\frac{\mathsf{d}^{M-1}}{\mathsf{d}t^{M-1}}f(t) + \ldots + c_{0}f(t) = w(t)$$

• defines a (stationary) GP covariance function (filters the white noise)

#### • procedure:

- convert GP covariance function to above form (spectrum matching procedure), truncate derivative expansion (if necessary)
- implement inference using Rauch–Tung–Striebel smoothing (truncation sets dimensionality of state space,  $\tilde{M}$ )
- Linear complexity in time (cubic in size of state-space, finding LDS parameters costly of learning):  $\mathcal{O}(T\tilde{M}^3)$





$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathsf{K}_{\mathbf{ff}} & \mathsf{K}_{\mathbf{fu}} \\ \mathsf{K}_{\mathbf{uf}} & \mathsf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$



•  $f_j$ 

1. augment model with M<T pseudo data

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

2. remove some of the dependencies (results in simpler model)



between blocks

chain structured pseudo-data

 $[f_i]$ 

**---(** U<sub>k</sub>)

1. augment model with M<T pseudo data

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

2. remove some of the dependencies (results in simpler model)



between blocks

chain structured pseudo-data

1. augment model with M<T pseudo data

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{ff}} & \mathbf{K}_{\mathbf{fu}} \\ \mathbf{K}_{\mathbf{uf}} & \mathbf{K}_{\mathbf{uu}} \end{bmatrix} \right)$$

2. remove some of the dependencies (results in simpler model)



K=3 blocks

between blocks

-**---(** u<sub>k</sub>)

chain structured pseudo-data

1. augment model with M<T pseudo data

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left( \begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix}; \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{\mathbf{f}\mathbf{f}} & \mathbf{K}_{\mathbf{f}\mathbf{u}} \\ \mathbf{K}_{\mathbf{u}\mathbf{f}} & \mathbf{K}_{\mathbf{u}\mathbf{u}} \end{bmatrix} \right)$$

2. remove some of the dependencies (results in simpler model)



 $K=3~{\rm blocks}$ 

between blocks

 $u_i = u_j = u_k$   $\rightarrow$   $u_i = u_j$  chain structured pseudo-data

 $(f_j)$ 

3. calibrate model

(e.g. using KL divergence, many choices)  $\underset{\{q(\mathbf{u}_{k}|\mathbf{u}_{k-1}),q(\mathbf{f}_{k}|\mathbf{u}_{k})\}_{k=1}^{K}}{\operatorname{arg\,min}} \operatorname{KL}(p(\mathbf{f},\mathbf{u})||\prod_{k=1}^{K}q(\mathbf{u}_{k}|\mathbf{u}_{k-1})q(\mathbf{f}_{k}|\mathbf{u}_{k}))$   $\implies \frac{q(\mathbf{u}_{k}|\mathbf{u}_{k-1}) = p(\mathbf{u}_{k}|\mathbf{u}_{k-1})}{q(\mathbf{f}_{k}|\mathbf{u}_{k}) = p(\mathbf{f}_{k}|\mathbf{u}_{k})} \quad \text{equal to exact conditionals}$ 

- **cost of inference is linear in T** (Rauch–Tung–Striebel smoothing)
- PITC, FITC and local versions are special cases (more edges deleted)
- sensible choices for the blocking and numbers of pseudo-points required

#### **Comparisons: training time (T=50,000)**



#### **Comparisons: training time (T=50,000)**



#### Comparisons: testing time (T=50,000)



# Summary

- taxonomy of GP approximations (direct/indirect, w/ or w/o pseudo-data, likelihood/prior)
- standard pseudo datapoint based approximations do not work well for long time-series
- chain structured variant performs better (outperforming other time-series GP approximations)
- **opportunity**: direct approximations in which pseudo-points only appear in the recognition model