

The Art of Gaussian Processes: Classical and Contemporary

César Lincoln C. Mattos¹ Felipe Tobar²

¹Department of Computer Science
Federal University of Ceará
Fortaleza, Ceará, Brazil
cesarlincoln@dc.ufc.br

²Initiative for Data & Artificial Intelligence
University of Chile
Santiago, Chile
ftobar@dim.uchile.cl

2021

Welcome!

In this tutorial we will first guide you through the **classics** and fundamentals of Gaussian Processes (GP) to then present **contemporary** advances in the field in an illustrative way.



César Lincoln C. Mattos is an Associate Professor at the Department of Computer Science, at Federal University of Ceará (UFC), Brazil. He is also an associate researcher at the Logics and Artificial Intelligence Group (LOGIA). He has research interests in (deep) probabilistic learning, approximate inference, system identification and anomaly detection.



Felipe Tobar is an Associate Professor at the Initiative for Data & Artificial Intelligence, Universidad de Chile, and the Coordinator of the Master of Data Science at the same University. He is also an Associate Researcher at the Center for Mathematical Modeling and the Advanced Centre for Electrical and Electronic Engineering. He teaches courses on Probability, Statistics and Machine Learning, and his research interests include time series, Bayesian inference, computational optimal transport and the societal impacts of machine learning.

Goals

Main

To present **classic** and **contemporary** perspectives to Gaussian processes to those interested in **understand, use and research** GP models.

Specific

- To illustrate the construction of the GP as a probabilistic model.
- To present and exemplify the use of vanilla GPs for regression, including covariance design, learning and inference.
- To demo a practical procedure of implementing a GP in Python
- To expand GP modeling concepts to non-Gaussian data and uncertain inputs.
- To review different ways of scaling inference in GP models to larger datasets.
- To showcase recent approaches to work with complex datasets.
- To extend the reach of the GP concept to multiple outputs

Outline

Duration: ~3 hours of presentation + ~1 hour of Q&A + a 10-min break

I. Welcome

- 1 Introduction (~5 min)

II. Classics

- 2 Fundamentals of Bayesian inference (~15 min)
- 3 From one to many random variables (~15 min)
- **Q&A** (~15 min)
- 4 The Gaussian process (~20 min)
- 5 Learning the kernel and its parameters (~15 min)
- 6 Implementation of a GP (~10 min)
- **Q&A** (~15 min) + **Break** (~10 min)
- 7 Beyond the Gaussian likelihood (~15 min)

III. Contemporary

- 8 Sparse approximations (~20 min)
 - **Q&A** (~15 min)
- 9 Current trends on kernel design (~10 min)
- 10 From GPLVM to Deep GPs (~15 min)
- 11 Multioutput GPs (~15 min)

IV. Discussion

- 12 Concluding remarks (~5min)
 - **Q&A** (~15 min)

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

Learning from Data

- Data: partial observations of the world that can reveal its inner working.
- Models linking observations and latent objects are needed: what kind? how general? trainable?
- In such setting, uncertainty is commonplace: data are incomplete and noisy. Uncertainty is randomness
- Learning is an inverse problem: it can be argued that learning under uncertainty requires *inverse probability*



The real
world is
complex

Data is
fragmentary

The
Model

(taken with permission from Thoughtfulnz on Twitter)

Probability

General < ----- > Particular

Population < ----- > Sample

Model < ----- > Data

Statistics

The Statistical Model

Notation

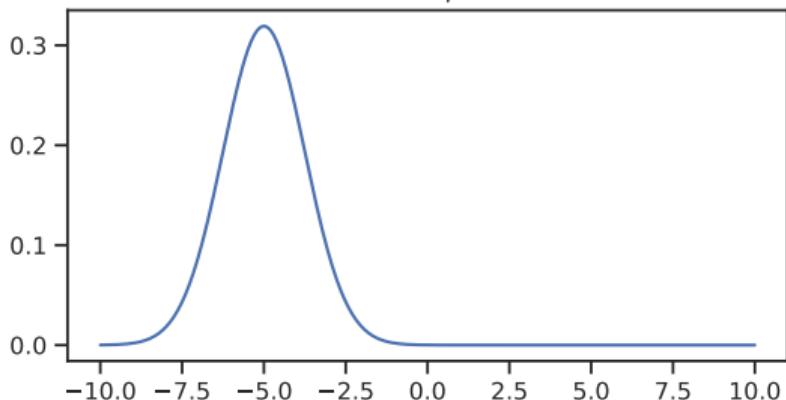
- **Sample space \mathcal{X} :** where the data live (usually \mathbb{N} , \mathbb{R} or \mathbb{R}^N).
- **Statistical model \mathcal{M} :** the set of probability distributions on \mathcal{X} indexed by a set \mathcal{T} :

$$\mathcal{M} = \{P_\theta \mid \theta \in \mathcal{T}\}. \quad (1)$$

- **Parameter θ and parameter space \mathcal{T} .**

Examples:

Gaussian with mean μ and variance σ^2



The Statistical Model

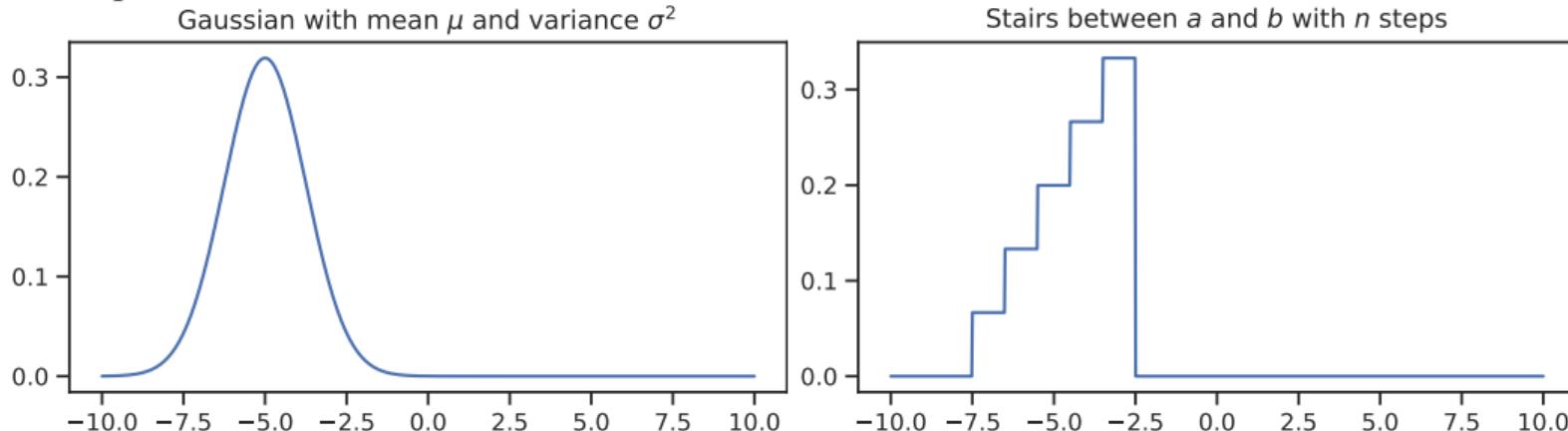
Notation

- **Sample space \mathcal{X} :** where the data live (usually \mathbb{N} , \mathbb{R} or \mathbb{R}^N).
- **Statistical model \mathcal{M} :** the set of probability distributions on \mathcal{X} indexed by a set \mathcal{T} :

$$\mathcal{M} = \{P_\theta \mid \theta \in \mathcal{T}\}. \quad (1)$$

- **Parameter θ and parameter space \mathcal{T} .**

Examples:



The Statistical Model

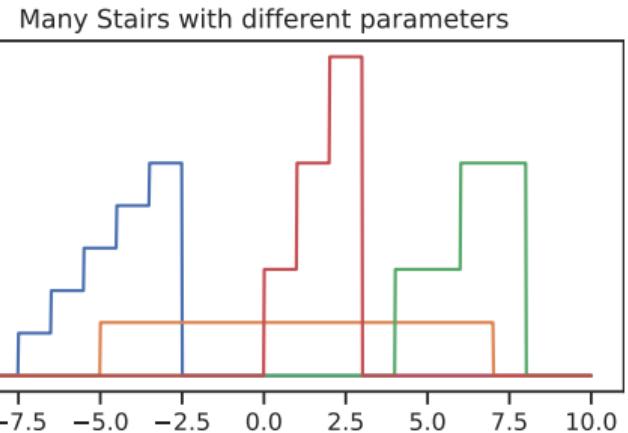
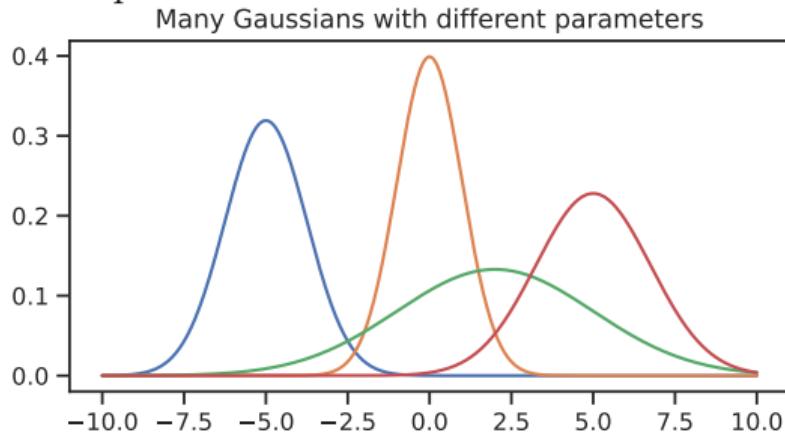
Notation

- **Sample space \mathcal{X} :** where the data live (usually \mathbb{N} , \mathbb{R} or \mathbb{R}^N).
- **Statistical model \mathcal{M} :** the set of probability distributions on \mathcal{X} indexed by a set \mathcal{T} :

$$\mathcal{M} = \{P_\theta \mid \theta \in \mathcal{T}\}. \quad (1)$$

- **Parameter θ and parameter space \mathcal{T} .**

Examples:

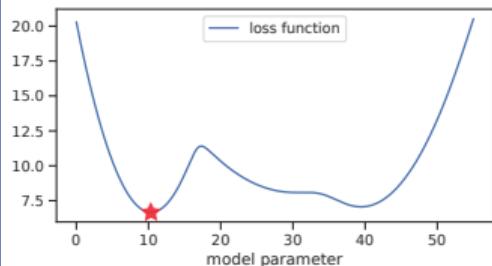


Finding the the parameter

(and thus the model)

Alternative 1: best model

- i) Define the model space
- ii) Gather data
- iii) Give each model a (fit) score:

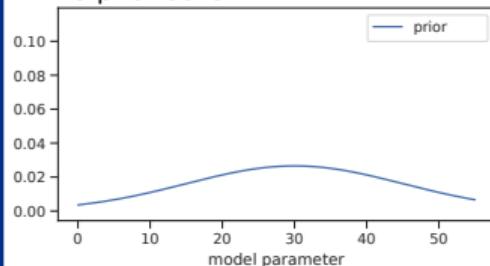


- iv) Chose the model with the best score ★

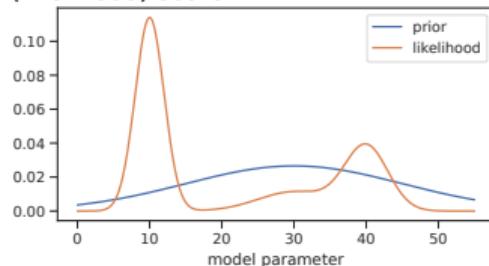
Predictions: "plug in" the chosen model into the predictor

Alternative 2: posterior belief

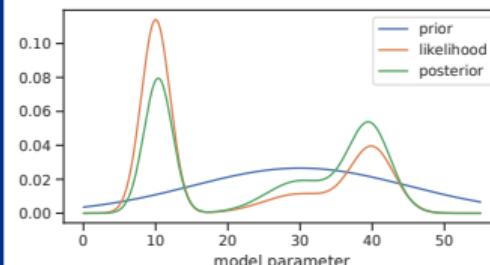
- i) Define the model space and equip it with a prior belief



- ii) Gather data and give each model a (likelihood) score



- iii) Update your prior belief into a posterior one



Predictions: predict with each model and average predictions wrt to posterior belief

Bayesian inference: a one-slide cheat sheet

Model specification: Before seeing any data, define some parameters as random variables (θ) and some as deterministic values (λ), hierarchically:

$$\lambda \in \Lambda \quad \& \quad \theta \sim p_\lambda(\theta) \quad \Rightarrow \quad X|\theta, \lambda \sim p_\lambda(x|\theta). \quad (2)$$

Now, in the light of data:

- **Set λ (model comparison):** e.g., by maximising the likelihood averages

$$p_\lambda(\text{data}) = \int p_\lambda(\text{data}|\theta)p_\lambda(\theta)d\theta. \quad (3)$$

- **Bayesian update:** compute the posterior over θ

$$p_\lambda(\theta) \rightarrow p_\lambda(\theta|\text{data}) = \frac{p_\lambda(\text{data}|\theta)p_\lambda(\theta)}{p_\lambda(\text{data})}. \quad (4)$$

- **Prediction:** via posterior averages

$$p_\lambda(x|\text{data}) = \int p_\lambda(x|\theta)p_\lambda(\theta|\text{data})d\theta. \quad (5)$$

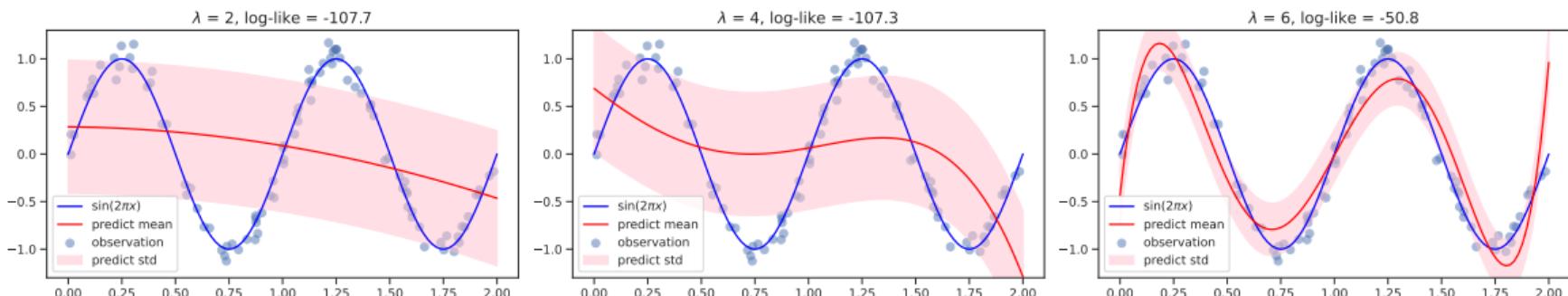
(given λ) the data is *summarised* in $p_\lambda(\theta|\text{data}) \Rightarrow$ the posterior replaces the data for prediction.

Example: Bayesian linear regression (1)

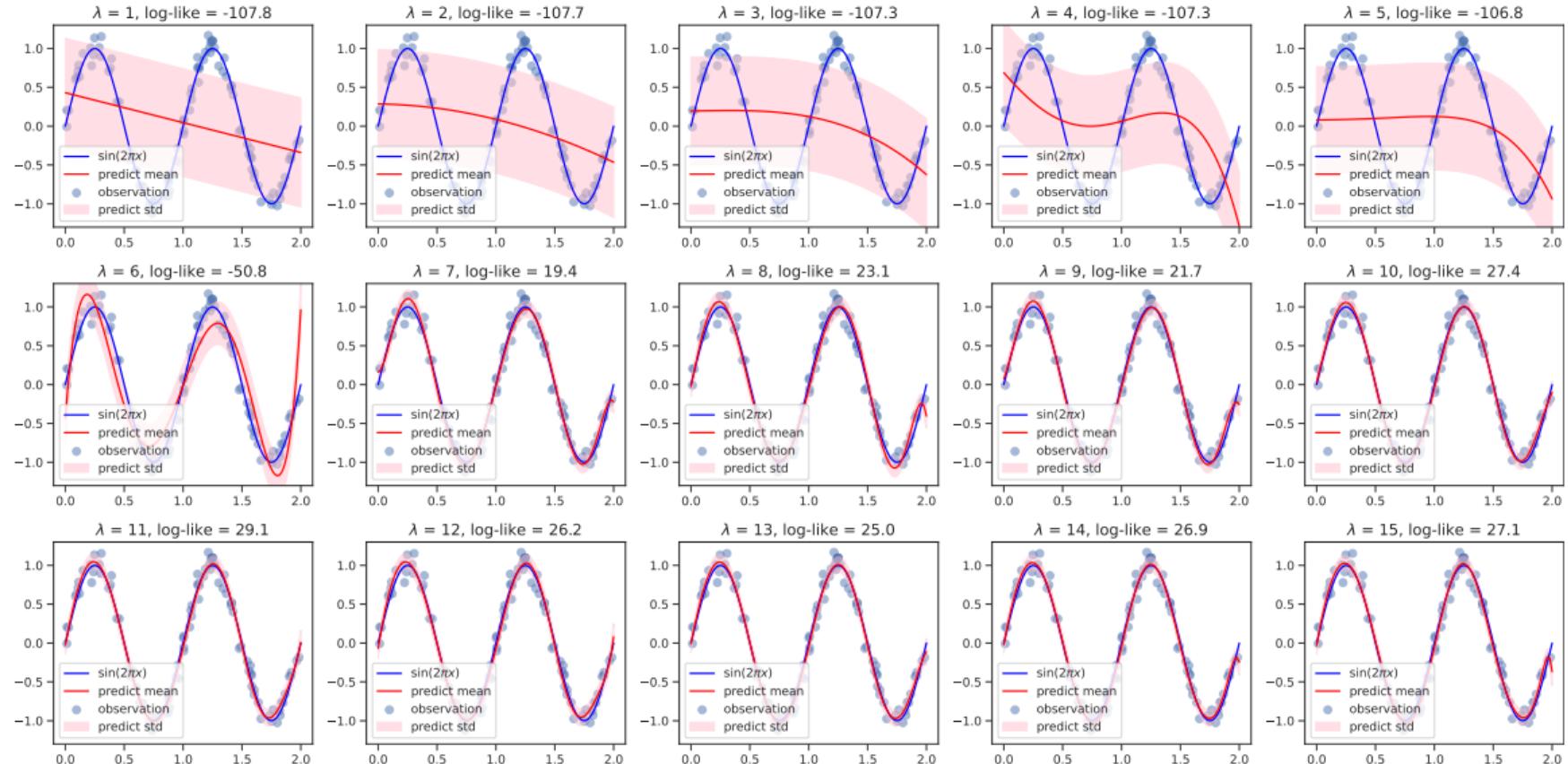
Let us consider $\lambda \in \{1, 2, 3, 4, \dots\}$, $\theta \sim \mathcal{N}(0, \sigma_\theta^2 \mathbf{I})$, $\sigma_\epsilon^2 \sim \text{Inv-Gamma}(a, b)$ and

$$X | \theta, \lambda \sim \mathcal{N}(x; \theta_0 + \theta_1 t + \theta_2 t^2 + \theta_3 t^3 + \theta_4 t^4 + \dots + \theta_\lambda t^\lambda, \sigma_\epsilon^2). \quad (6)$$

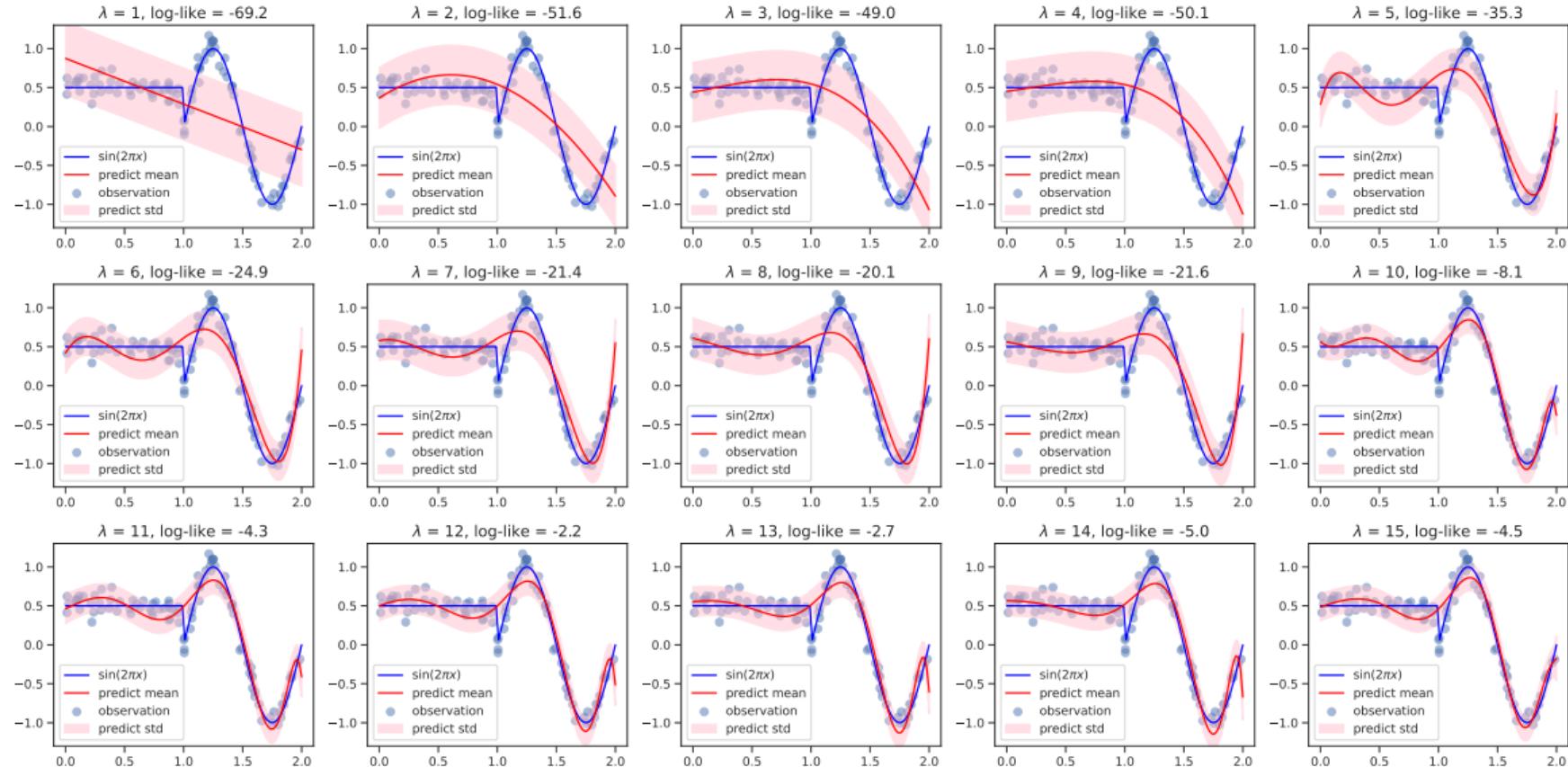
This prior is conjugate: the posterior over (θ, σ^2) is also a Normal-Inverse-Gamma. However, for the order λ , we evaluate each case:



Bayesian linear regression: more



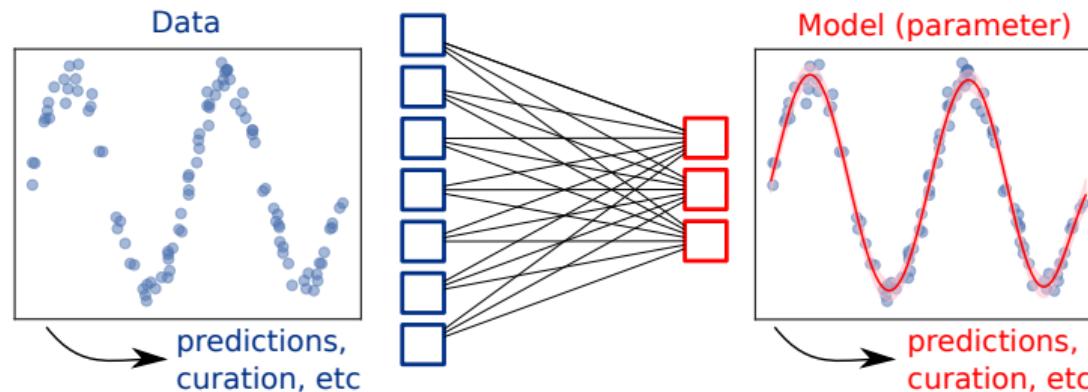
Bayesian linear regression: more (different dataset)



The cardinality of the parameter

Parameters (usually) live in a space where we can search, such as \mathbb{R}^N (or a space bijective with \mathbb{R}^N). Its dimensionality N can be understood in two (related) ways:

the **complexity** of the model & the amount of **information** extracted from the data



Thus it is natural to wonder:

- how *large* should the parameter be?
- can we bypass that question and use **infinite-dimensional** parameters?
- if so, how can we learn infinite parameters with a finite amount of data?

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

Why Gaussians? i) Central Limit theorem

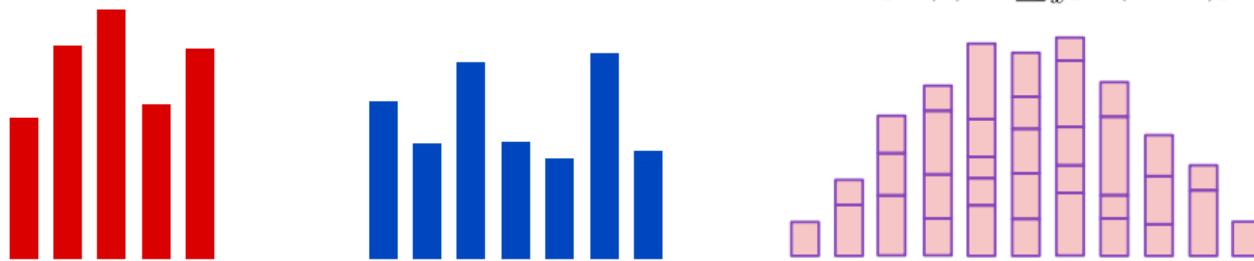
Let us assume the phenomenon we want to model is made up of multiple sources of uncertainty, which are independent and possibly of different distributions, added together.

$$X \sim p_X(x)$$

$$Y \sim p_Y(y)$$

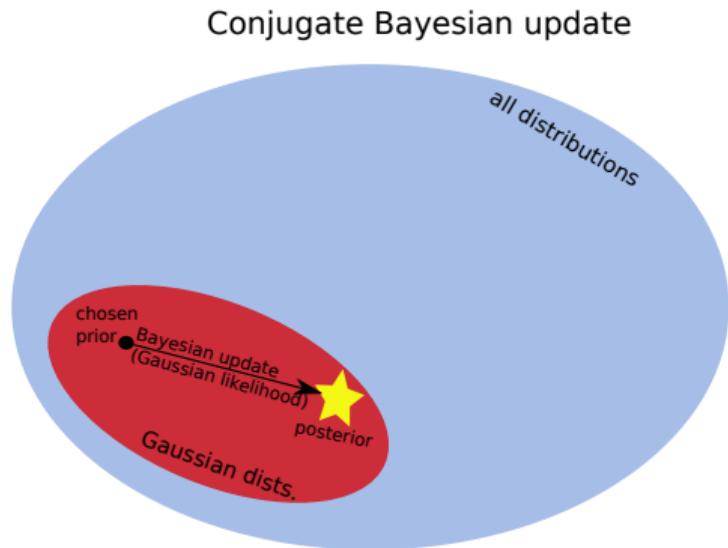
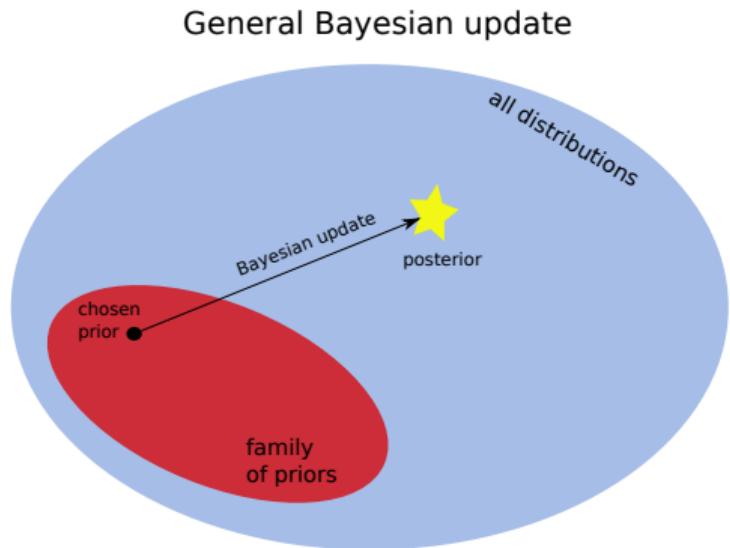
$$Z = X + Y$$

$$Z \sim p_Z(z) = \sum_x p_Y(z - x)p_X(x)$$



The sum of two RVs with a density given by the **convolution** of the densities of the original RVs. The Gaussian distribution is the fixed-point of this operation (convolution and scaling), therefore, the more sources of uncertainty we add together the *more Gaussian* the result **irrespective of the original distributions**.

Why Gaussians? ii) Conjugacy



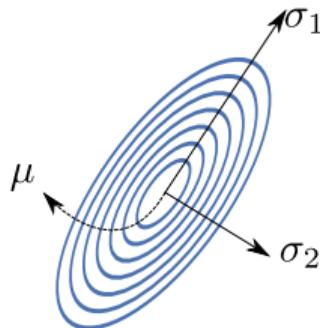
It is desired that the prior and posterior belong to the same family for the following reasons:

- the dimension of the parameter does not change (needed for numerical implementation)
- the properties of the model can be easily compared before/after data

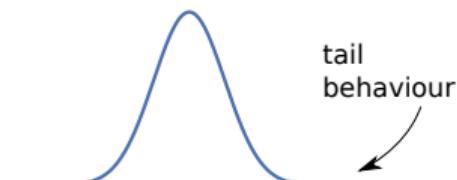
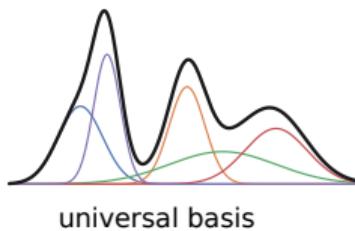
The Gaussian law

Some properties

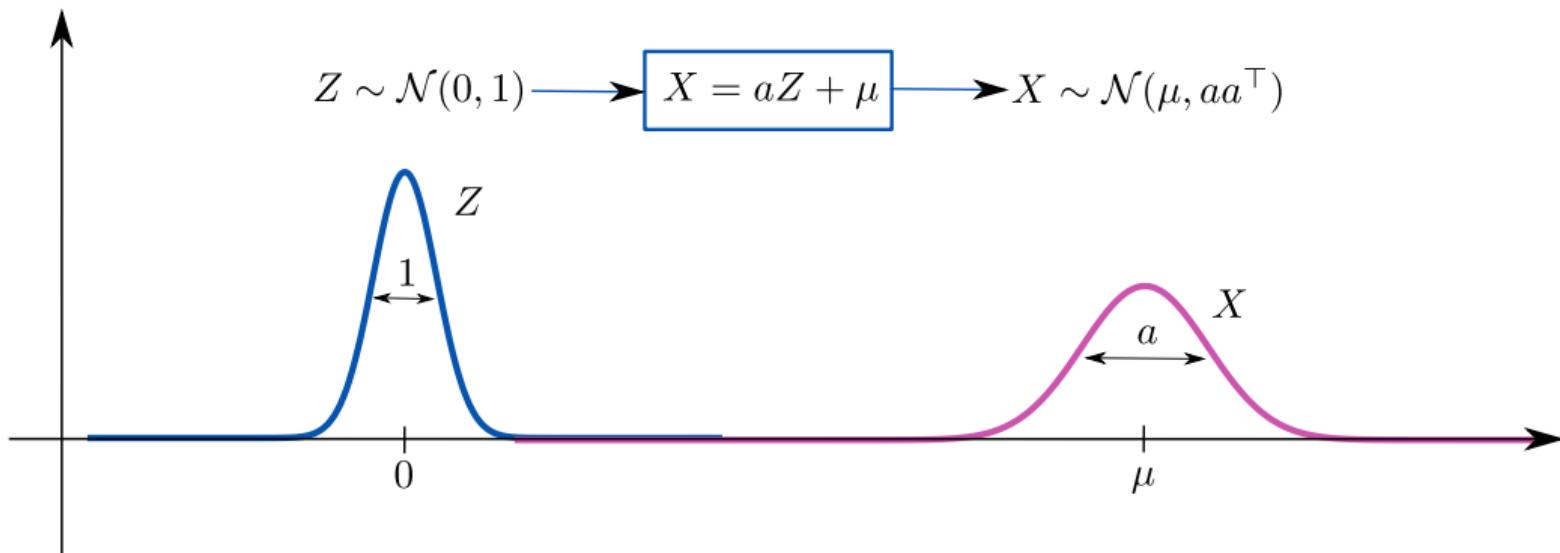
- central limit theorem
- conjugacy
- infinite support
- infinitely differentiable
- maximum entropy (fixed variance)
- convex loss (linear regression)
- conditional expectation is linear
- equivalence to least squares
- uncorrelatedness \Leftrightarrow independence
- universal basis



$$X \sim \mathcal{N}(\mu, \sigma^2) \Leftrightarrow p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

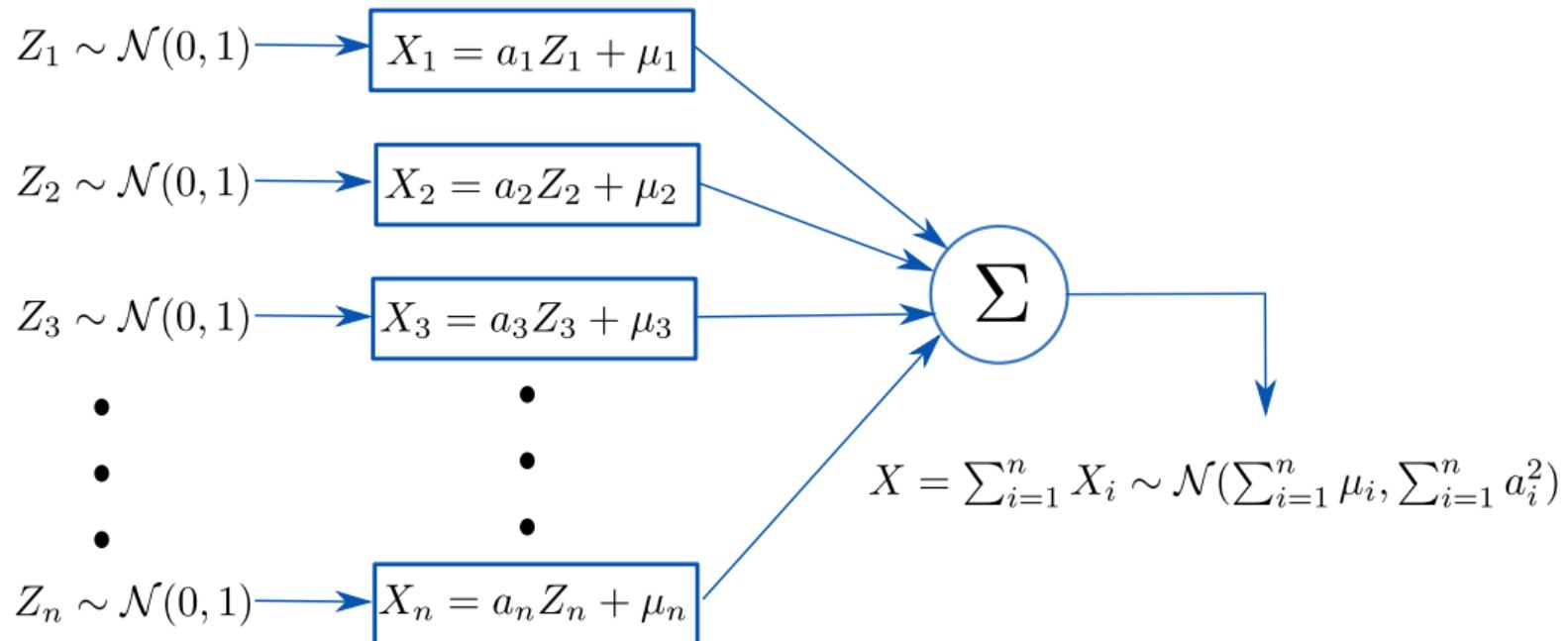


The Gaussian law (linear transformation)



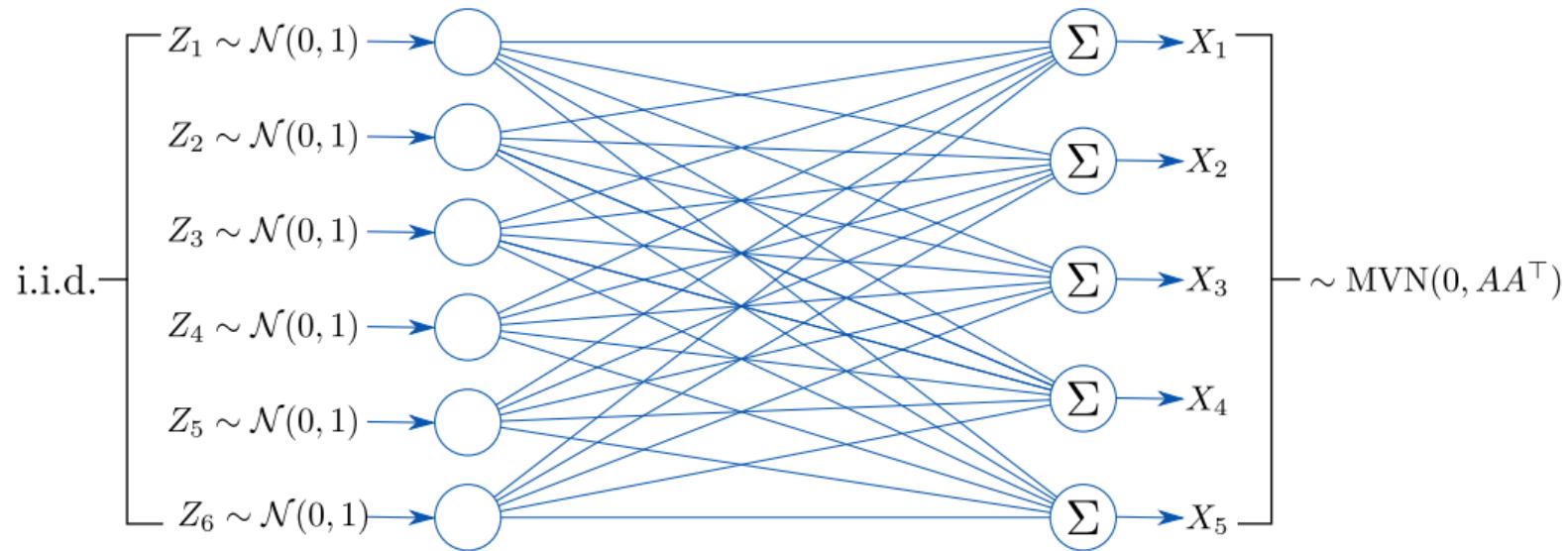
- Closed under linear transformations
- Any Gaussian RV can be produced from a linear transformation of $\mathcal{N}(0, 1)$

The Gaussian law (linear combinations)



(we have assumed the RVs to be independent but the result holds for the correlated case)

The Gaussian law (multivariate case)

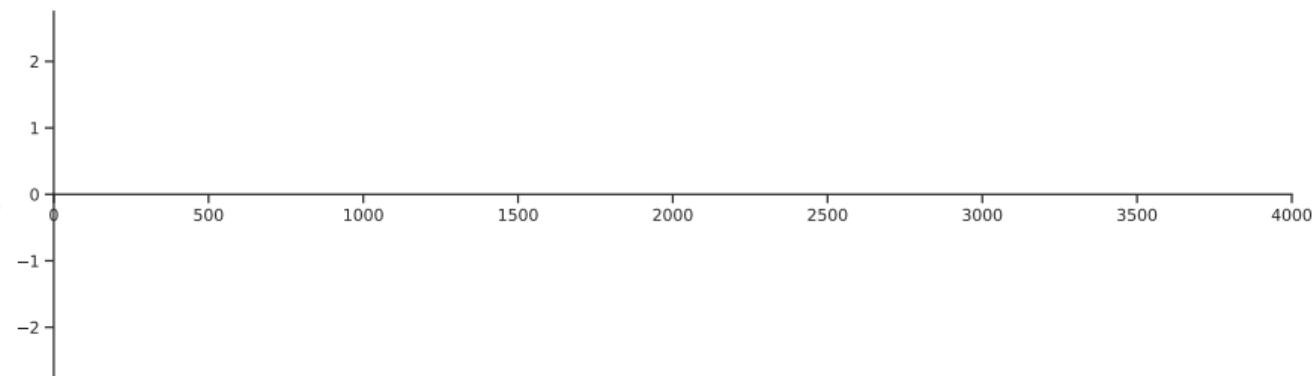
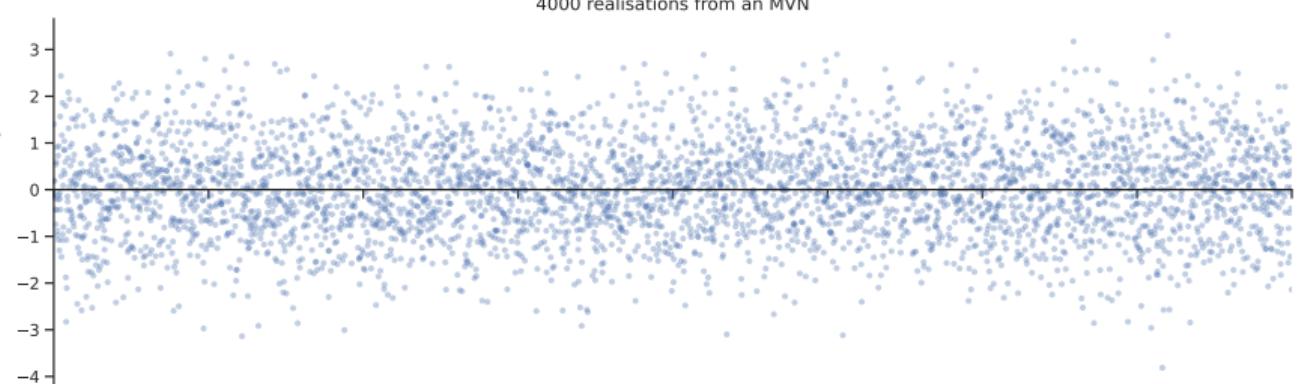


$$X \sim \text{MVN}(\mu, \Sigma) \Leftrightarrow p(x) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp \left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right)$$

MVN: The Multivariate Normal Distribution

A numerical illustration

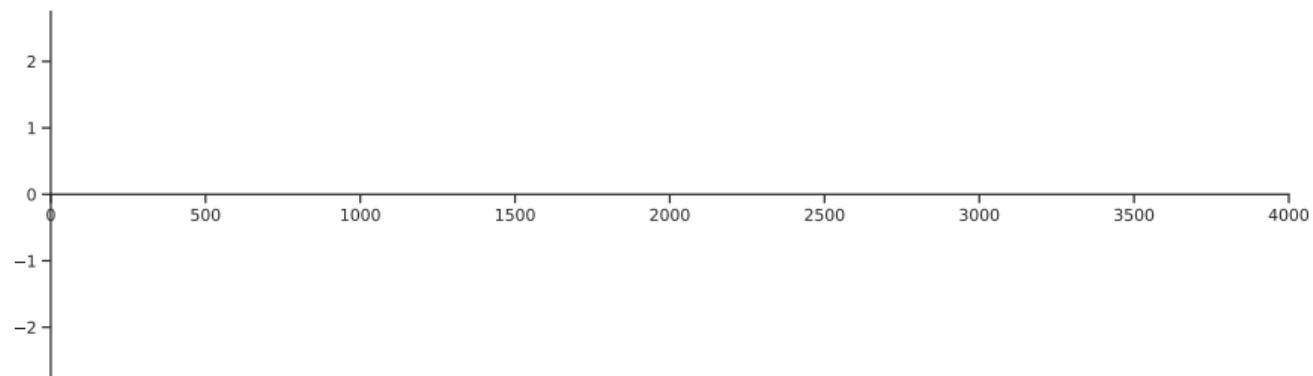
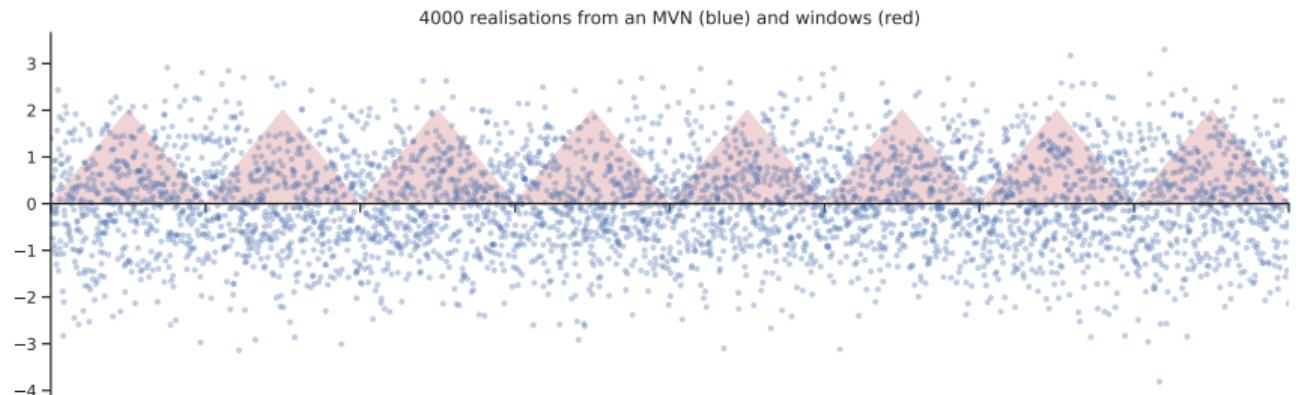
- 4000 realisations of $Z \sim \mathcal{N}(0, 1)$, in **blue**, against their indices
- a set of *mixing weights* (or windows), in **red**
- compute the inner product between each **window** and the **realisations** it intersects.
Plot the result below.
- notice the effect of overlapping windows



MVN: The Multivariate Normal Distribution

A numerical illustration

- 4000 realisations of $Z \sim \mathcal{N}(0, 1)$, in **blue**, against their indices
- a set of *mixing weights* (or windows), in **red**
- compute the inner product between each **window** and the **realisations** it intersects.
Plot the result below.
- notice the effect of overlapping windows



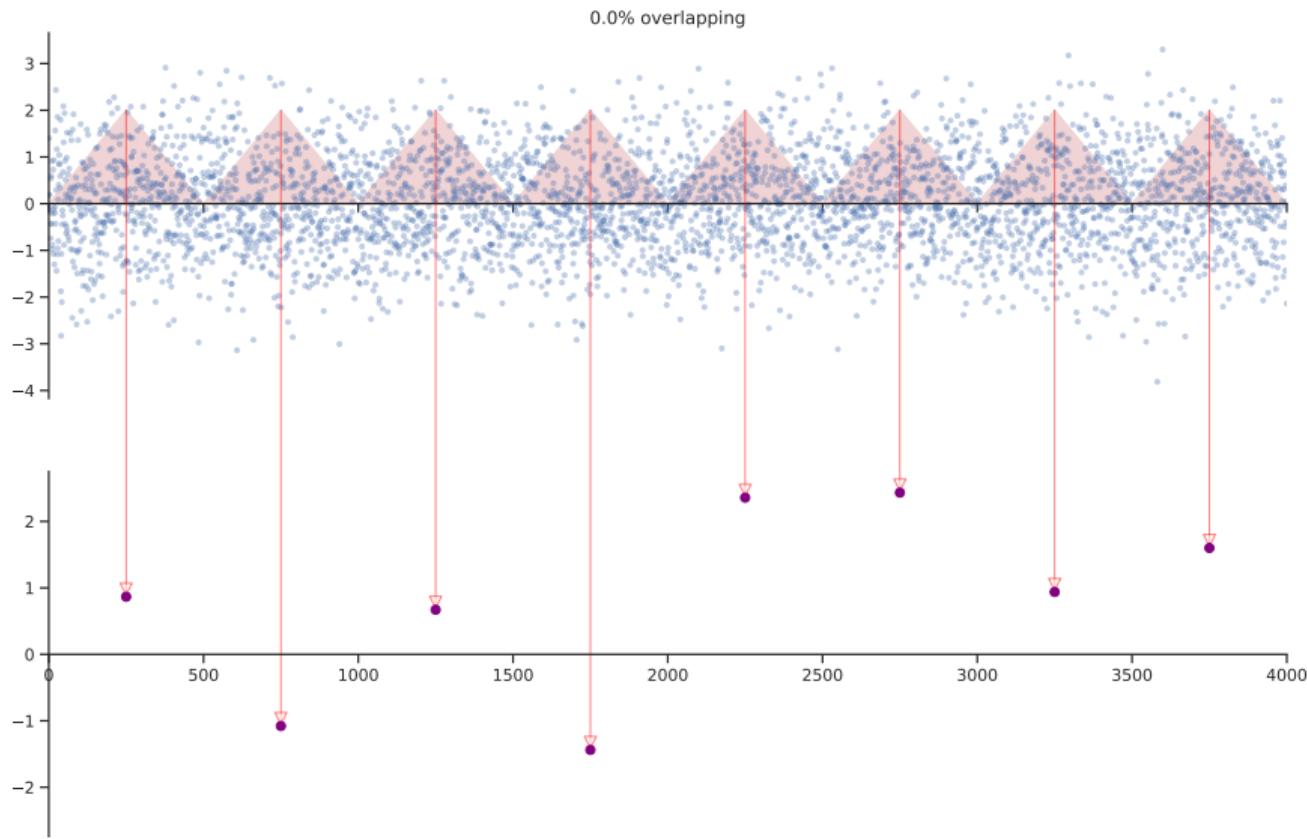
MVN: The Multivariate Normal Distribution

A numerical illustration

- 4000 realisations of $Z \sim \mathcal{N}(0, 1)$, in **blue**, against their indices
- a set of *mixing weights* (or windows), in **red**
- compute the inner product between each **window** and the **realisations** it intersects.

Plot the result below.

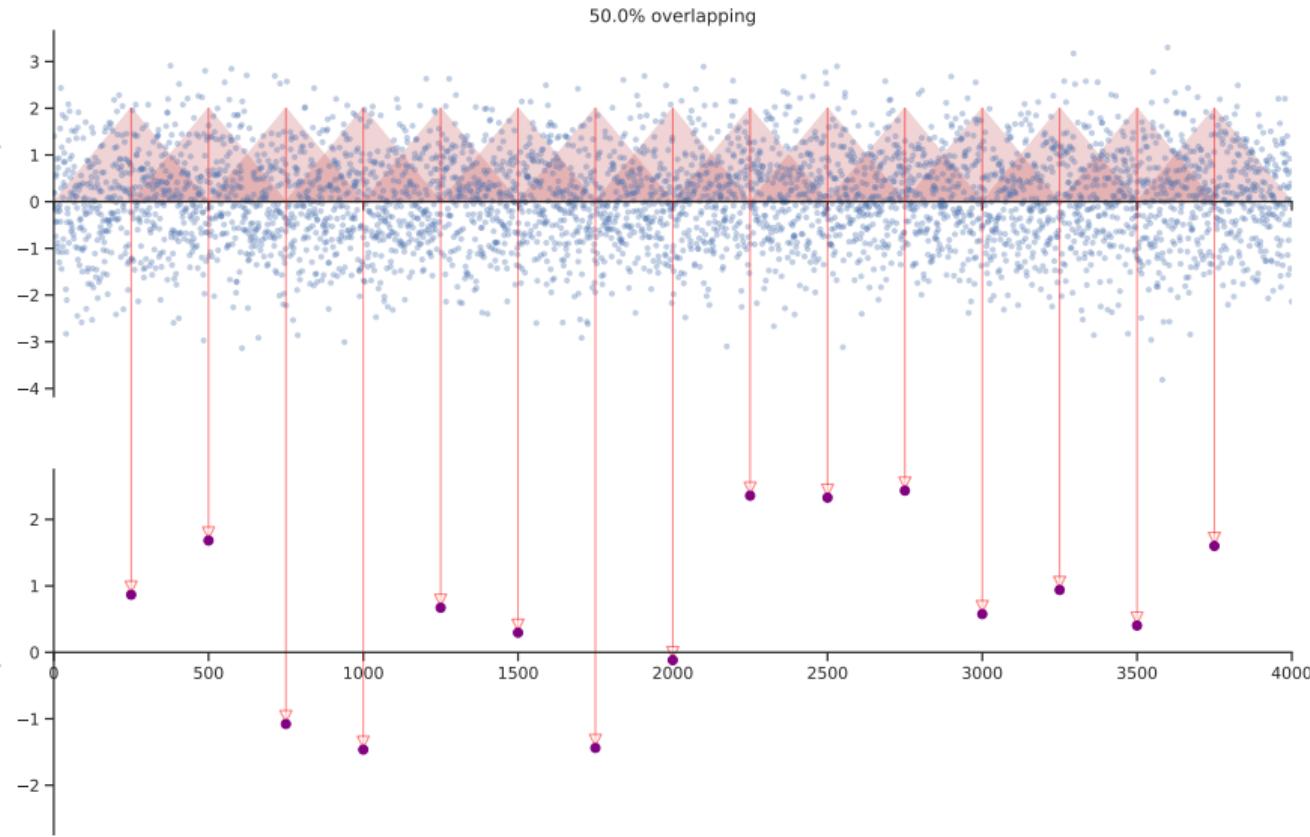
- notice the effect of overlapping windows



MVN: The Multivariate Normal Distribution

A numerical illustration

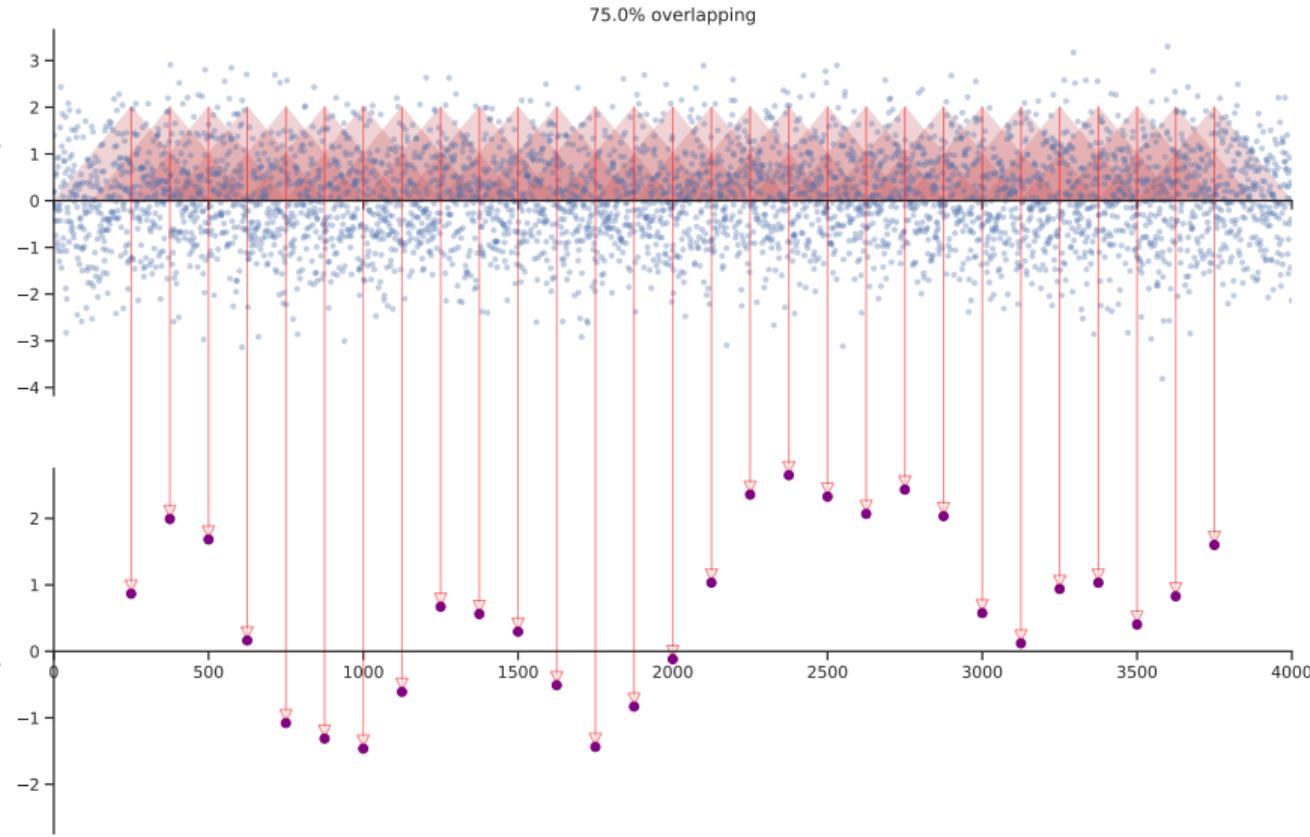
- 4000 realisations of $Z \sim \mathcal{N}(0, 1)$, in **blue**, against their indices
 - a set of *mixing weights* (or windows), in **red**
 - compute the inner product between each **window** and the **realisations** it intersects.
- Plot the result below.
- notice the effect of overlapping windows



MVN: The Multivariate Normal Distribution

A numerical illustration

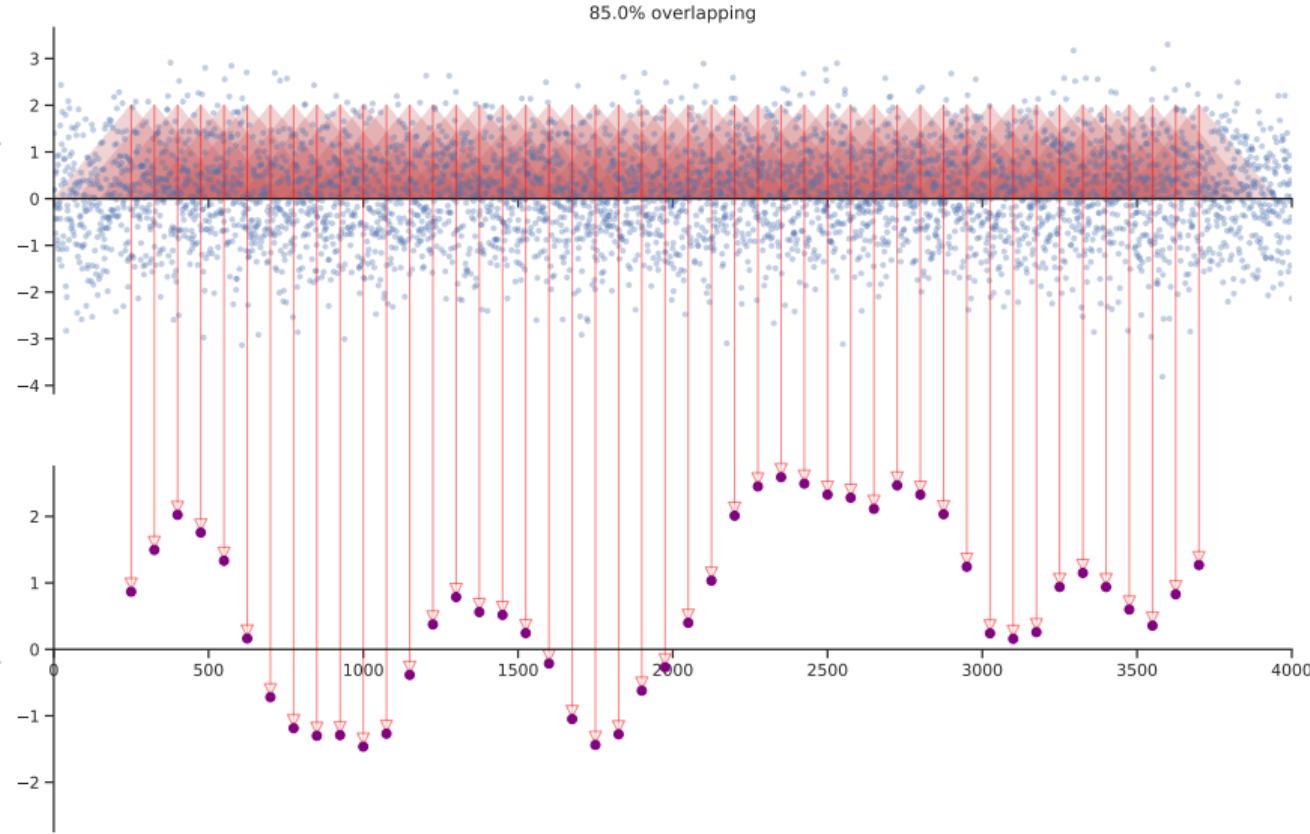
- 4000 realisations of $Z \sim \mathcal{N}(0, 1)$, in **blue**, against their indices
 - a set of *mixing weights* (or windows), in **red**
 - compute the inner product between each **window** and the **realisations** it intersects.
- Plot the result below.
- notice the effect of overlapping windows



MVN: The Multivariate Normal Distribution

A numerical illustration

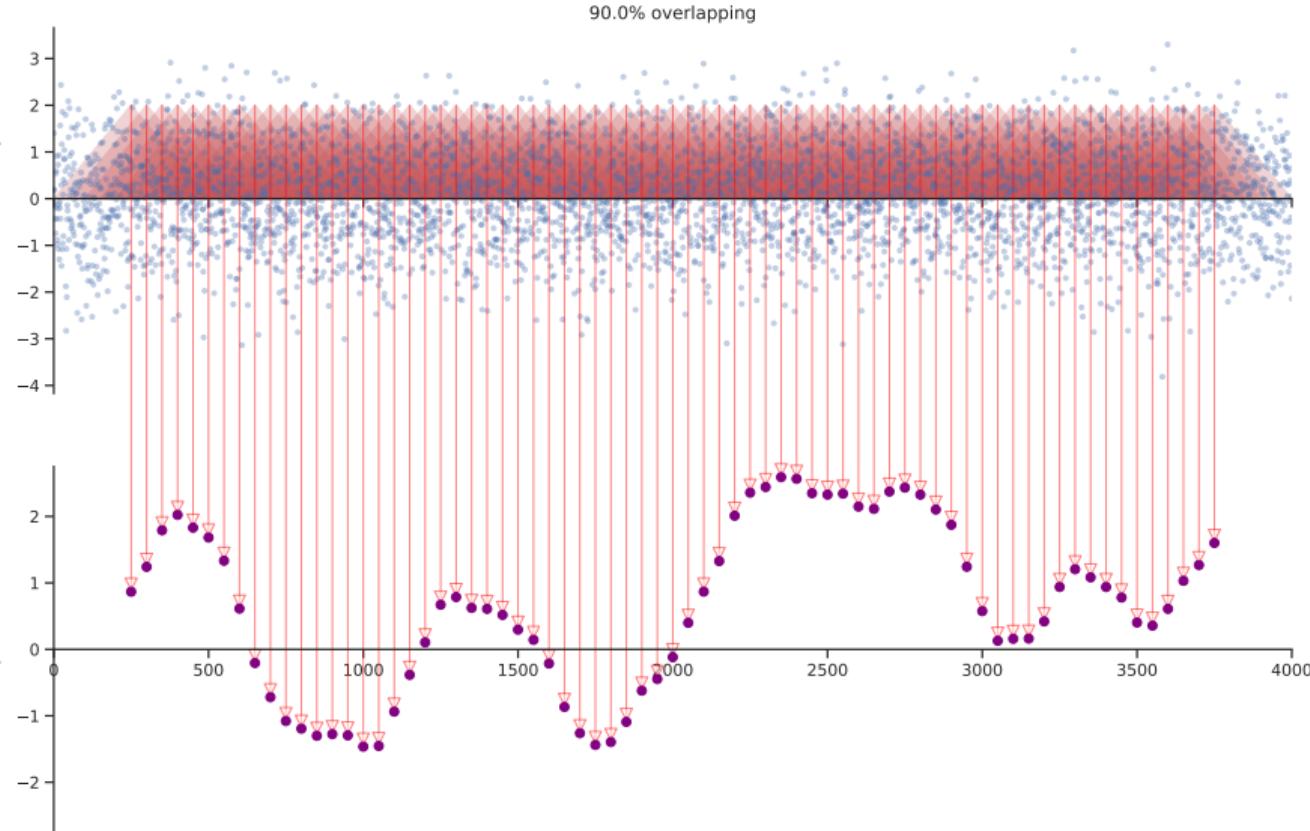
- 4000 realisations of $Z \sim \mathcal{N}(0, 1)$, in **blue**, against their indices
 - a set of *mixing weights* (or windows), in **red**
 - compute the inner product between each **window** and the **realisations** it intersects.
- Plot the result below.
- notice the effect of overlapping windows



MVN: The Multivariate Normal Distribution

A numerical illustration

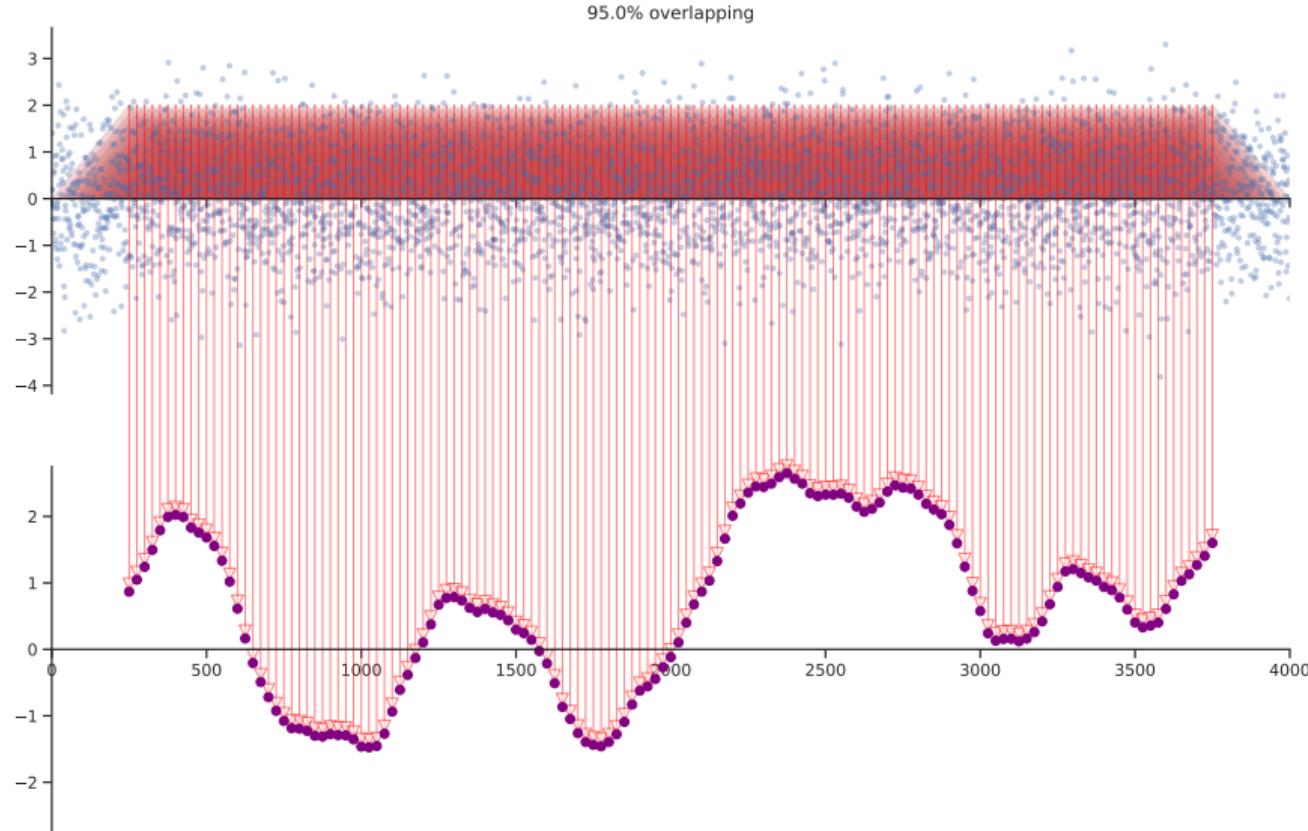
- 4000 realisations of $Z \sim \mathcal{N}(0, 1)$, in **blue**, against their indices
 - a set of *mixing weights* (or windows), in **red**
 - compute the inner product between each **window** and the **realisations** it intersects.
- Plot the result below.
- notice the effect of overlapping windows



MVN: The Multivariate Normal Distribution

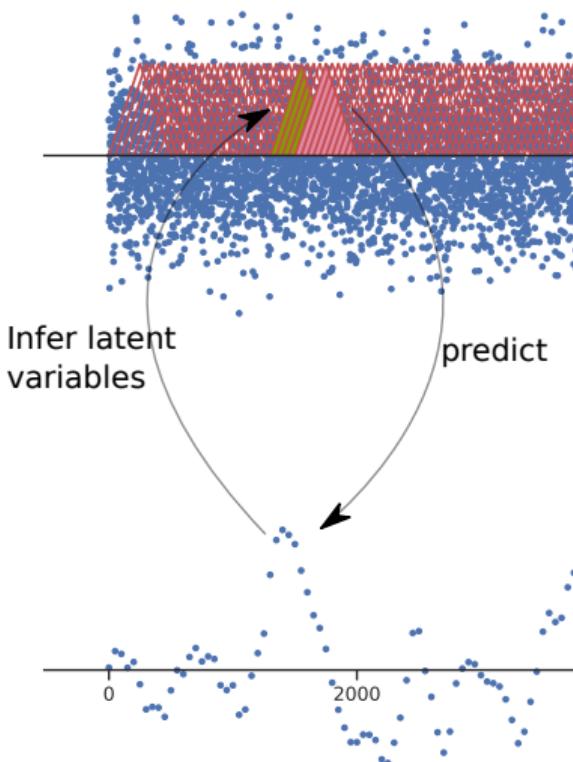
A numerical illustration

- 4000 realisations of $Z \sim \mathcal{N}(0, 1)$, in **blue**, against their indices
- a set of *mixing weights* (or windows), in **red**
- compute the inner product between each **window** and the **realisations** it intersects.
Plot the result below.
- notice the effect of overlapping windows



Observations from the above illustration

Correlation and Inference



- The induced structure in X is a consequence of the shared latent variables Z
- How many latent variables should we consider? Intuitively, more than the observations X
- We can model missing data as an inverse problem:
observations → latent variables → predictions
- The induced law is:

$$X \sim \text{MVN}(\mu, \Sigma^2) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right)$$

where the mean μ and covariance matrix Σ are uniquely determined by the mixing weights and the law of Z .

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

The Covariance: detaching from the latent variables

Is this a kernel trick?

Recall that if $X = A^\top Z$, with Z being a collection of IID- $\mathcal{N}(0, 1)$, then

$$p_X(x) = \text{MVN}(x; \mu, \Sigma^2) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right), \quad (7)$$

where

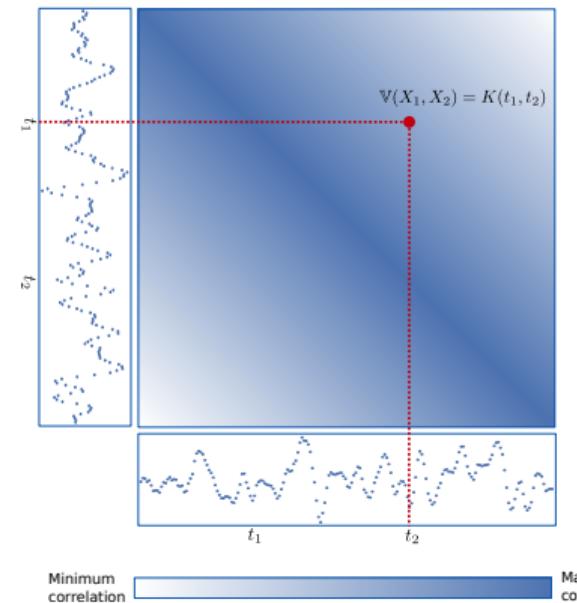
- $\mu = 0$ — we assume zero mean
- $\Sigma = A^\top A$ — we need more latent variables than observations, since $\text{rank}(\Sigma) = \text{rank}(A)$.

The advantage of building Σ through A was that A was unrestricted. Designing Σ , however, needs

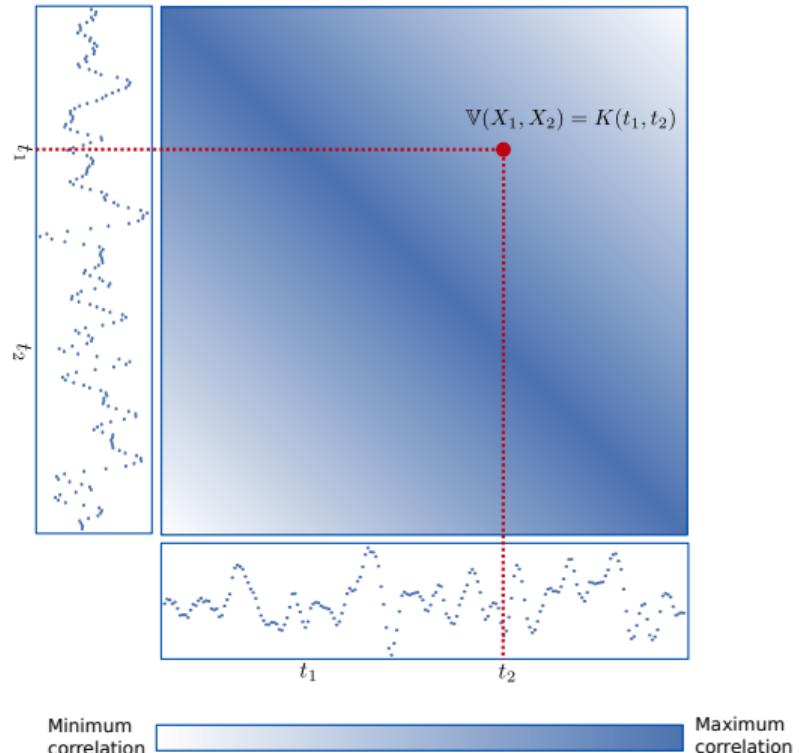
- Symmetry: $\Sigma = \Sigma^\top$
- Positive semidefinitenes: $v^\top \Sigma v \geq 0, \forall v \in \mathbb{R}^n$

⇒ we *only* need to choose $\frac{1}{2}n(n + 1)$ values

(we bypassed the latent variables)

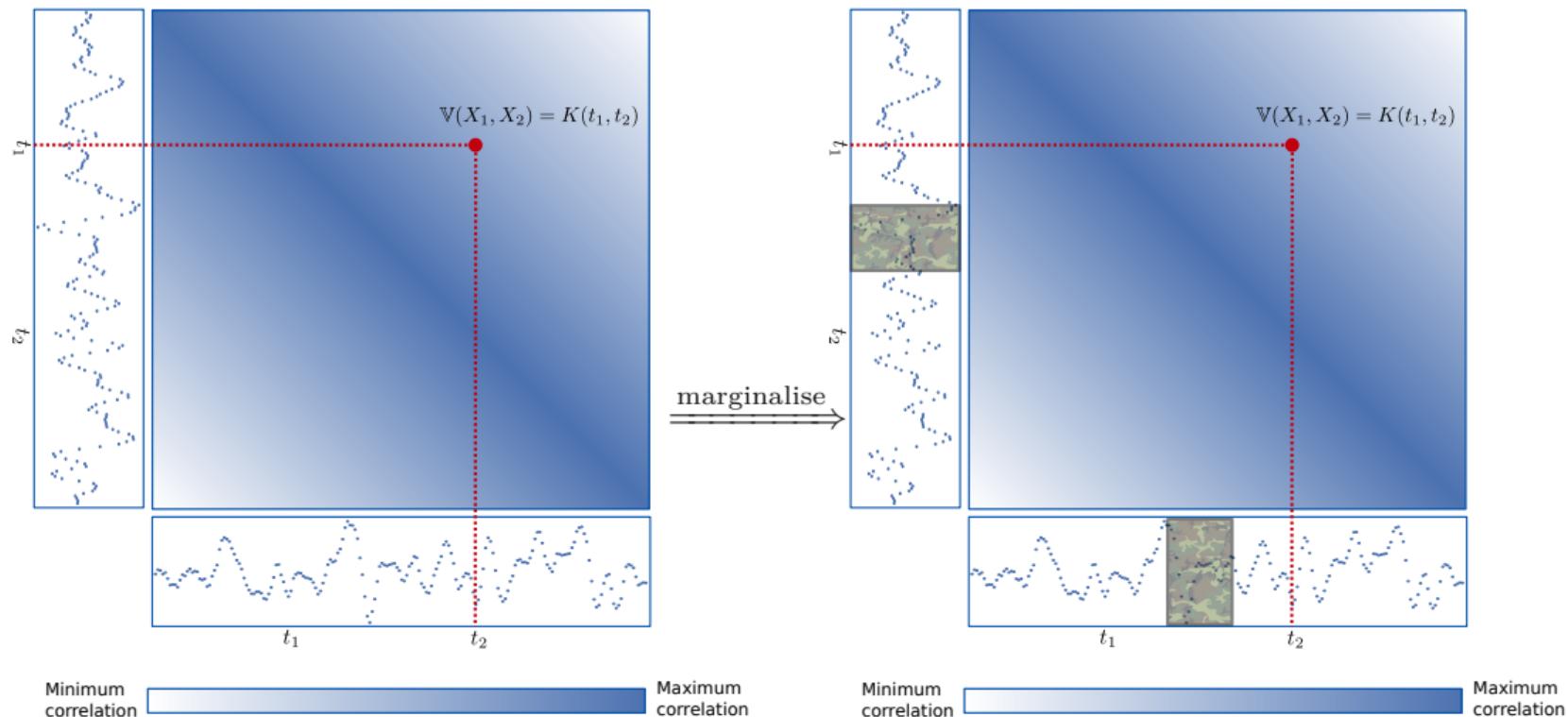


Some properties of the Covariance

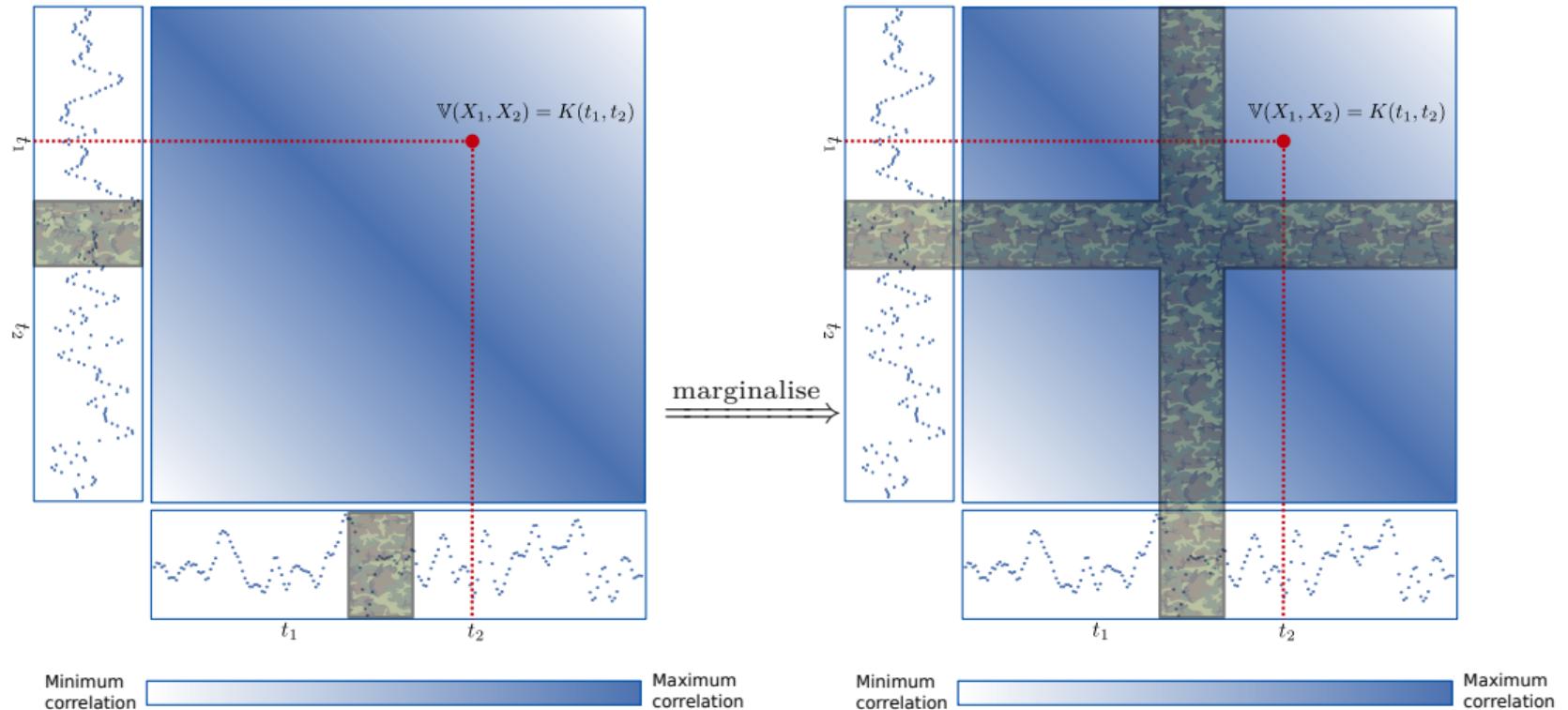


- symmetry: $\mathbb{V}[t_1, t_2] = \mathbb{V}[t_2, t_1]$
- positive semidefiniteness implies diagonal dominance: *an RV cannot co-vary more with another RV than itself*
- Structural dependence among Gaussians are second-order only: the covariance fully determines the relationship among Gaussians and thus uncorrelatedness is equivalent to independence
- **Gaussians will be Gaussians**, no matter how many of them we take

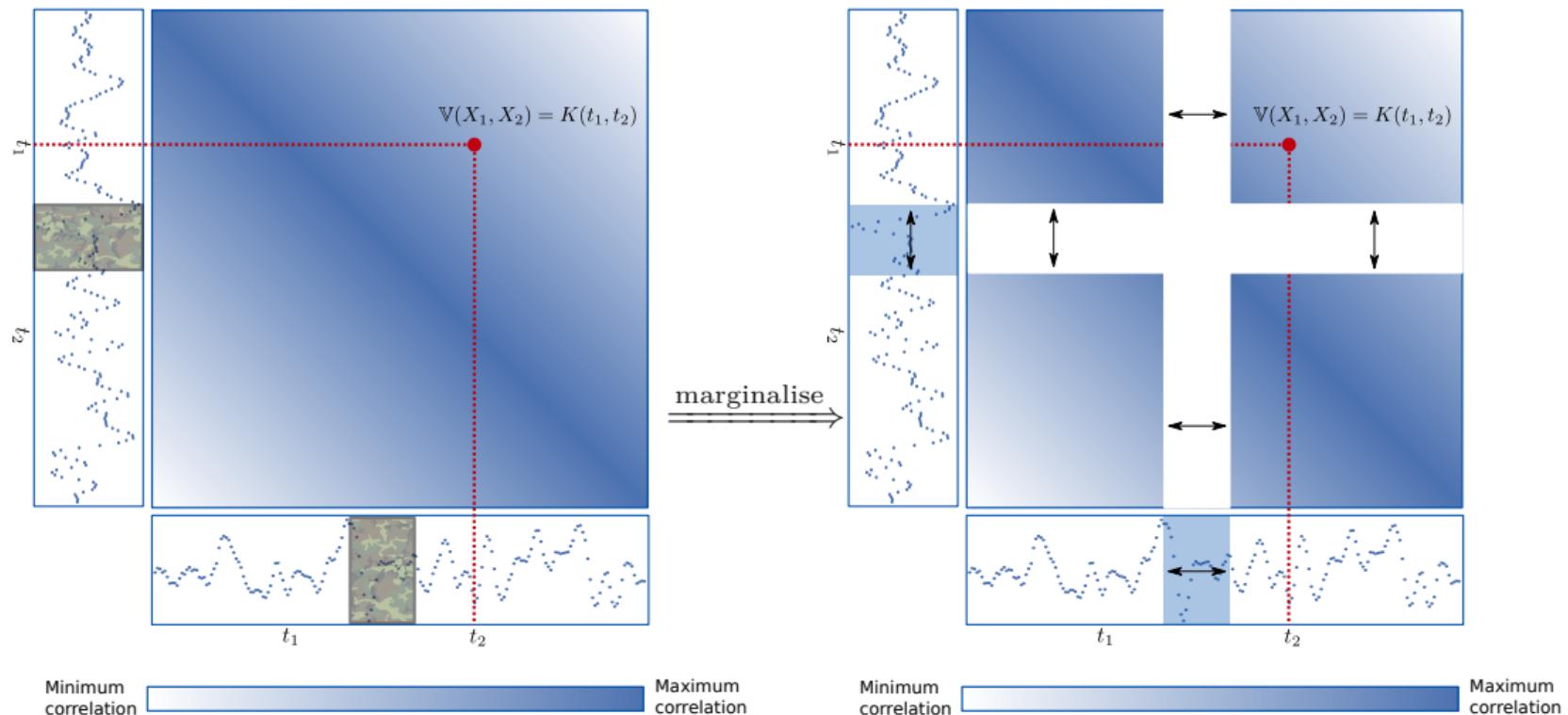
The Marginalisation property of the Gaussians



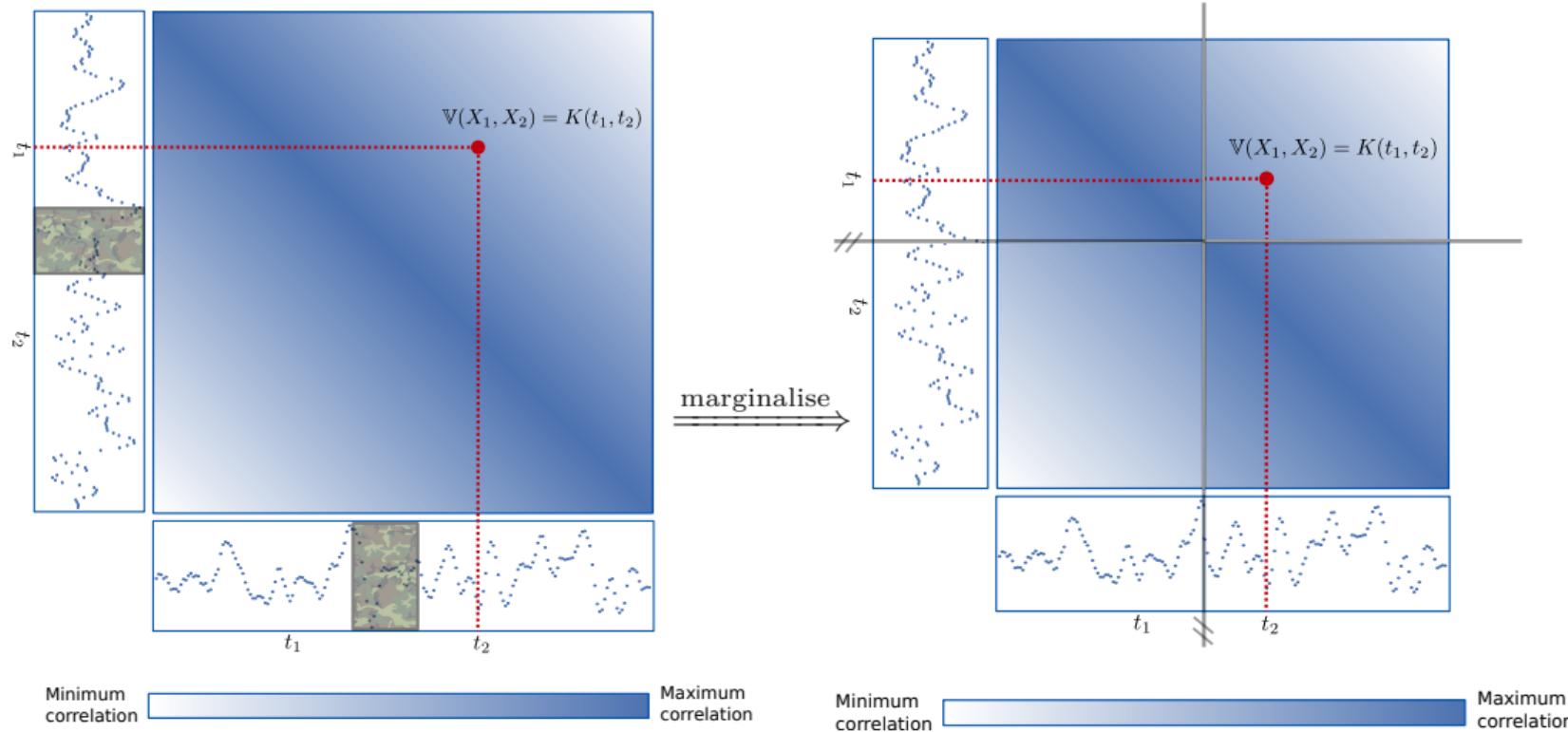
The Marginalisation property of the Gaussians



The Marginalisation property of the Gaussians



The Marginalisation property of the Gaussians



Marginalisation: removing RVs

Another way to understand the closure of the Gaussians under marginalisation is to present it as a linear operation (recall that linear combinations of Gaussians are also Gaussian).

If $y \in \mathbb{R}^m$ is a marginalisation of $x \in \mathbb{R}^n$ ($m \leq n$):

$$[y]_i = [x]_{\lambda(i)}, \quad (8)$$

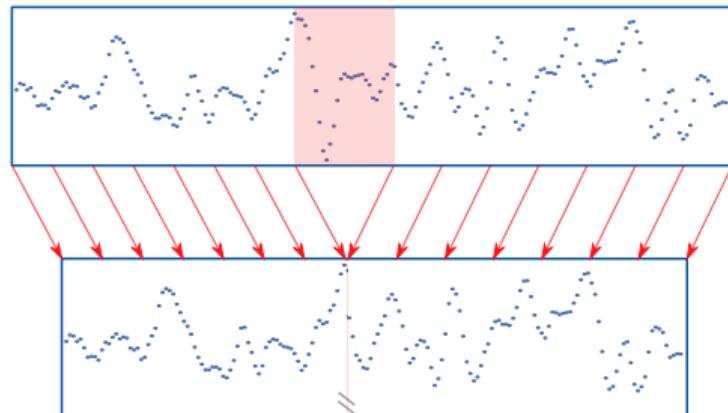
where $\lambda : \{1, 2, \dots, m\} \rightarrow \{1, 2, \dots, n\}$ is an injective non-surjective function. Then, we can write

$$y = Mx, \quad M \in \mathbb{R}^{m \times n} \quad (9)$$

where $[M]_{\lambda(i), i} = 1$ and zero elsewhere. Therefore, $X \sim \text{MVN}(\mu, \Sigma)$ implies

$$Y = MX \sim \text{MVN}(M\mu, \mathbf{M}\Sigma\mathbf{M}^\top) \quad (10)$$

where $\mathbf{M}\Sigma\mathbf{M}^\top$ is the marginalisation of Σ .



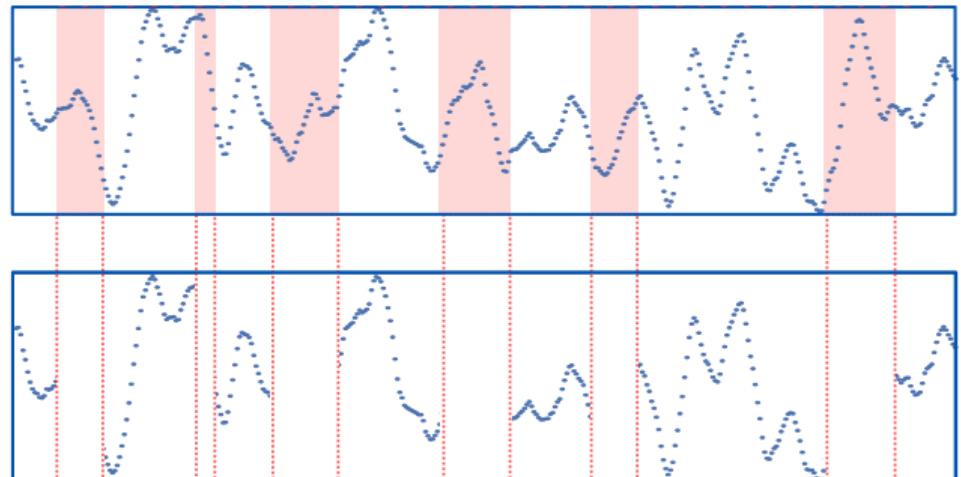
Marginalisation: where are these Gaussian RVs coming from?

- We can marginalise various subsets of data - *still Gaussian*

- **Idea:** The MVN vector X might be itself the marginalisation from another, longer, vector

- What if the *original* process was a continuous or *dense* collection of RVs?

- That continuous-time stochastic process does exist (due to the Kolmogorov's extension theorem)



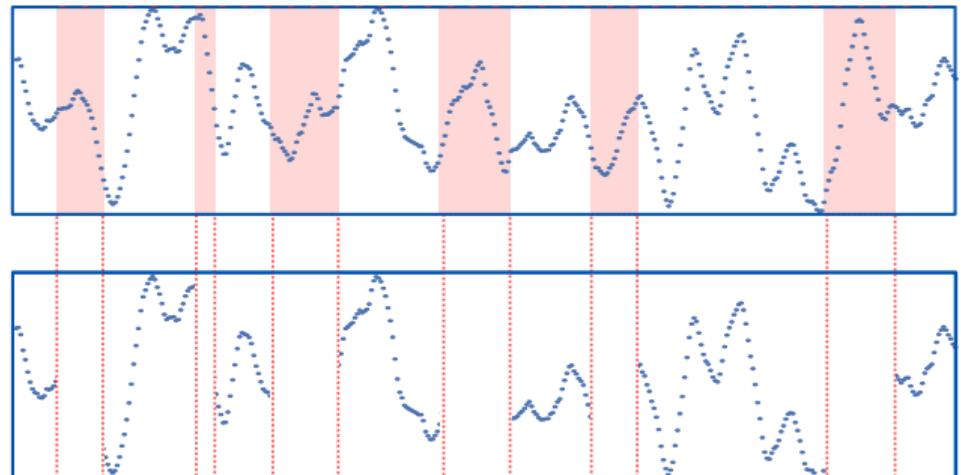
Marginalisation: where are these Gaussian RVs coming from?

- We can marginalise various subsets of data - *still Gaussian*

- **Idea:** The MVN vector X might be itself the marginalisation from another, longer, vector

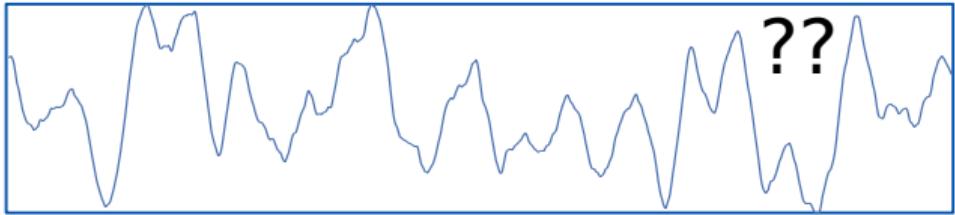
- What if the *original* process was a continuous or *dense* collection of RVs?

- That continuous-time stochastic process does exist (due to the Kolmogorov's extension theorem)

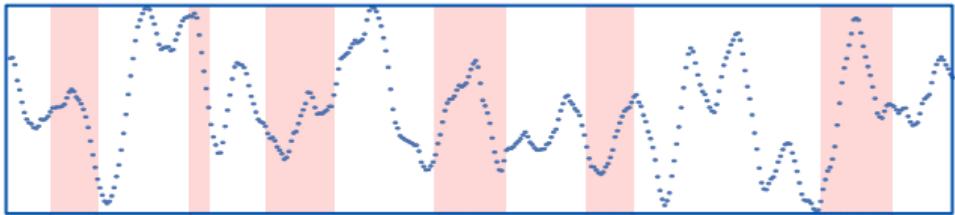


Marginalisation: where are these Gaussian RVs coming from?

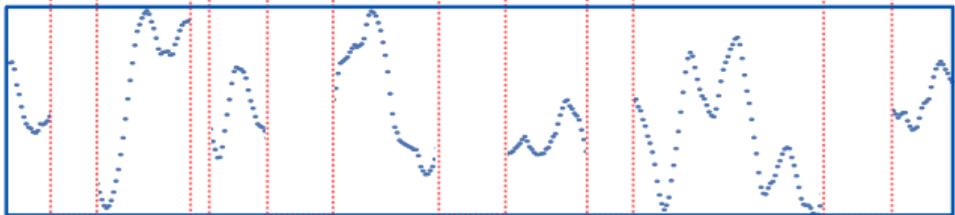
- We can marginalise various subsets of data - *still Gaussian*



- **Idea:** The MVN vector X might be itself the marginalisation from another, longer, vector



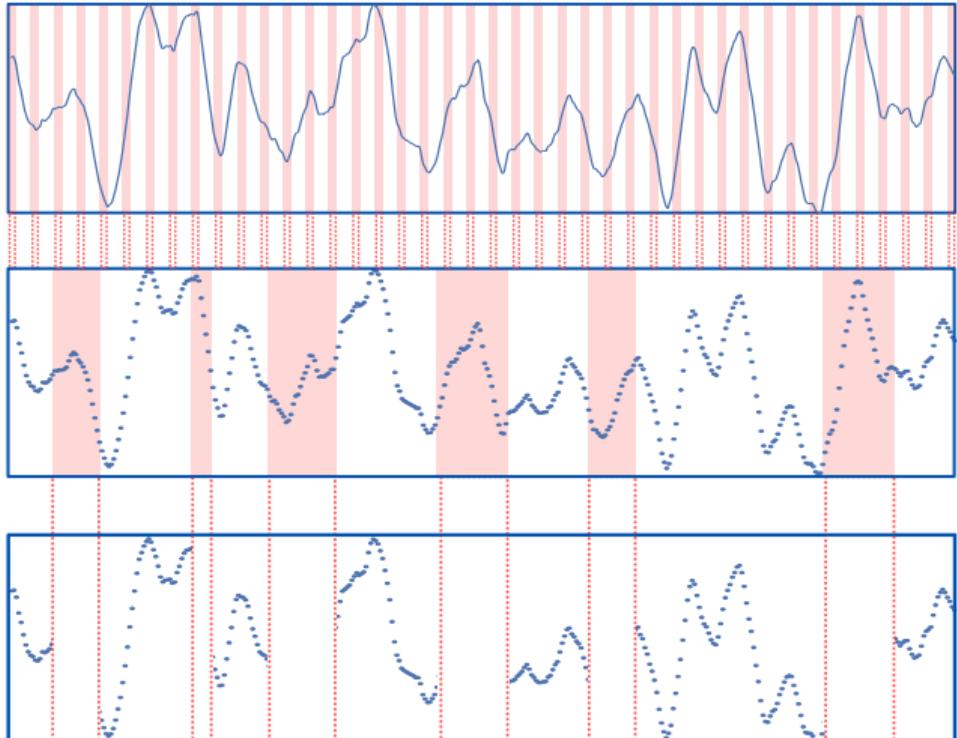
- What if the *original* process was a continuous or *dense* collection of RVs?



- That continuous-time stochastic process does exist (due to the Kolmogorov's extension theorem)

Marginalisation: where are these Gaussian RVs coming from?

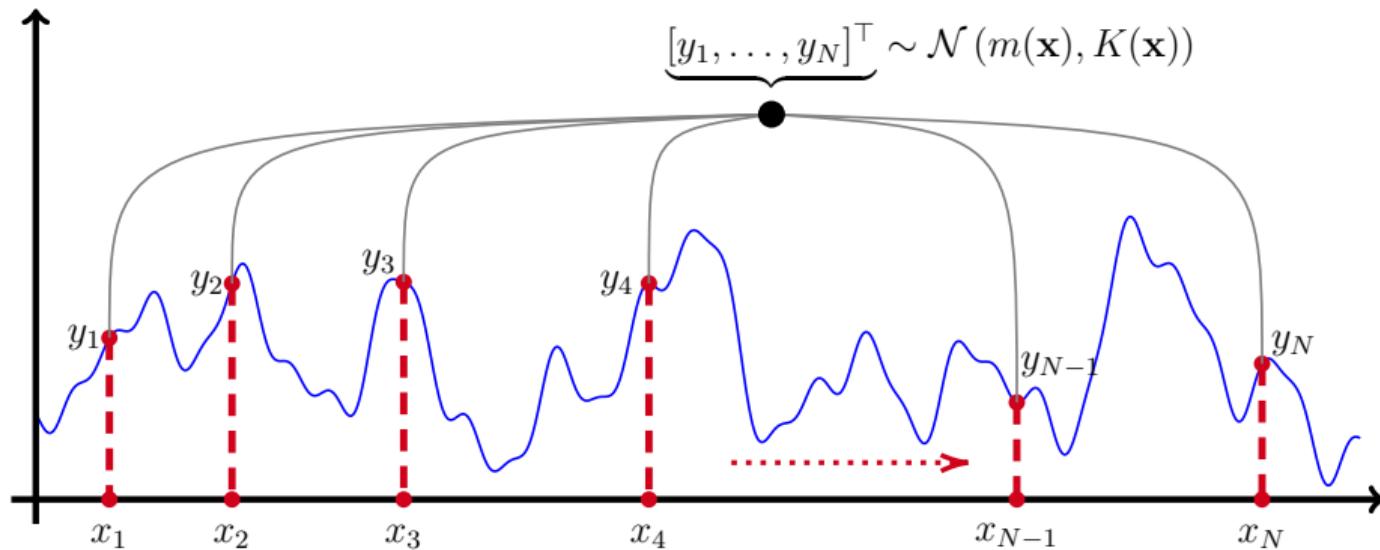
- We can marginalise various subsets of data - *still Gaussian*
- **Idea:** The MVN vector X might be itself the marginalisation from another, longer, vector
- What if the *original* process was a continuous or *dense* collection of RVs?
- That continuous-time stochastic process does exist (due to the Kolmogorov's extension theorem)



Enter the Gaussian Process

A nonparametric generative model for continuous functions

Definition: A GP is a stochastic process such that any finite collection of values follows a multivariate normal distribution.

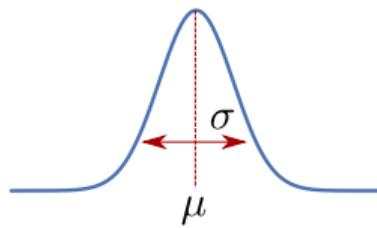


(change of) Notation: $f \sim \mathcal{GP}(\mu, K) \iff f(\mathbf{x}) \sim \mathcal{N}(m(\mathbf{x}), K(\mathbf{x}, \mathbf{x})), \forall \mathbf{x} \in \mathcal{X}^n, n \in \mathbb{N}$

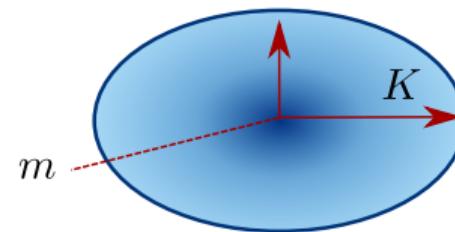
GPs - what are they?

- GPs are nonparametric **prior distributions** over the space of functions.
- These distributions are defined by a **mean** (function) and a **covariance** (kernel)
- GPs are used for **Bayesian inference**: it is possible to update such belief using data (observations of the function)
- They are **nonparametric**: The samples of the GP cannot be represented by a finite-dimensional parameters (but rather by all their *infinite* values)

Gaussian over the real line



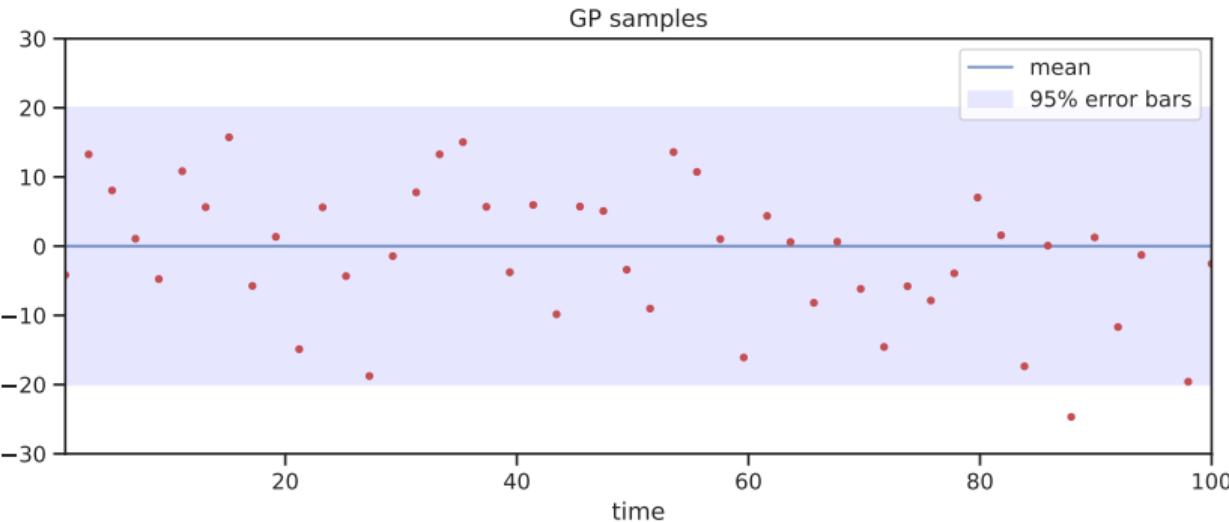
Gaussian over function space



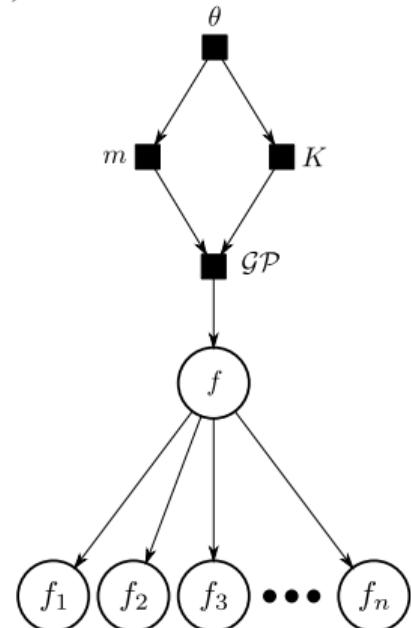
Hyperparameters: $m : x \mapsto m(x)$ belongs to the space of functions and $K : (x, x') \mapsto K(x, x')$ is a positive-semidefinite symmetric function. We refer all the parameters of these functions as θ .

Sampling from the GP prior

- The (complete) sample of a GP $f \sim \mathcal{GP}$ is infinite dimensional
- However, using the marginalisation property we can draw (finite) *parts* of the GP sample
- Set up a grid and sample the *marginalisation* of f in that grid (an MVN)

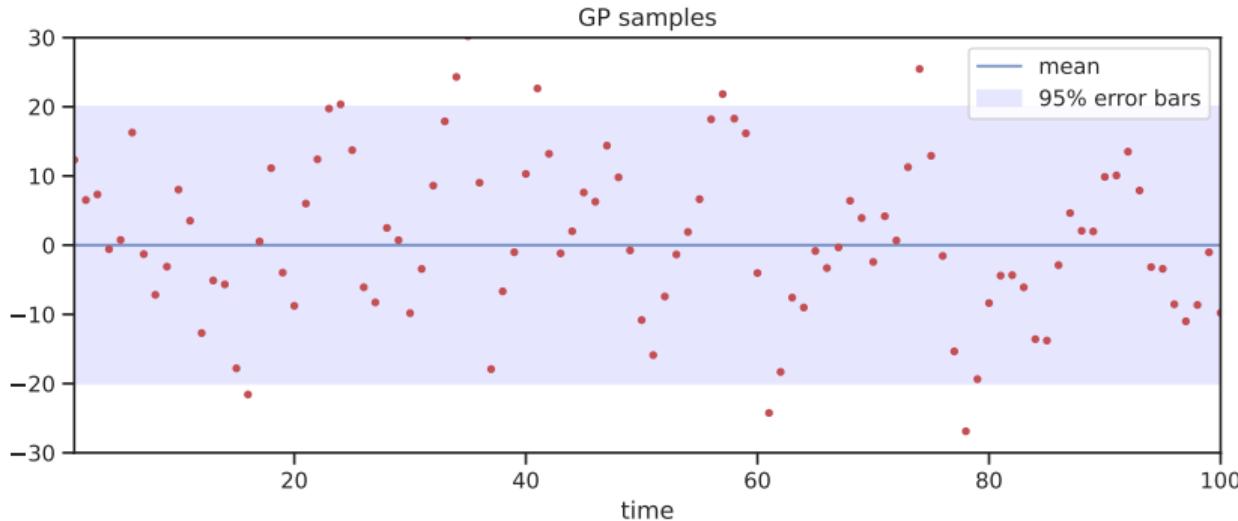


- Trick: use `linestyle = '-'` to make the grid look continuous :)

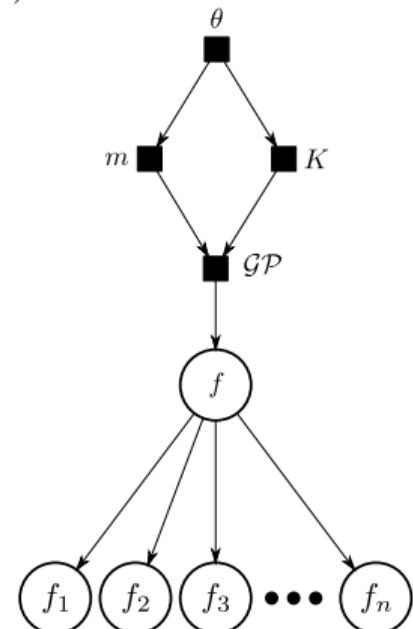


Sampling from the GP prior

- The (complete) sample of a GP $f \sim \mathcal{GP}$ is infinite dimensional
- However, using the marginalisation property we can draw (finite) *parts* of the GP sample
- Set up a grid and sample the *marginalisation* of f in that grid (an MVN)

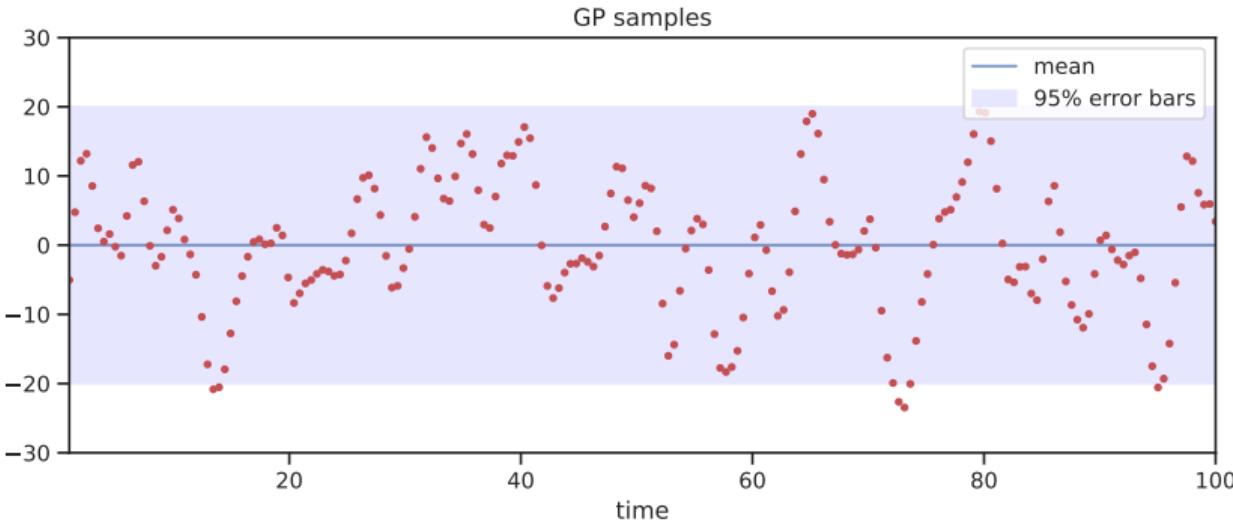


- Trick: use `linestyle = '-'` to make the grid look continuous :)

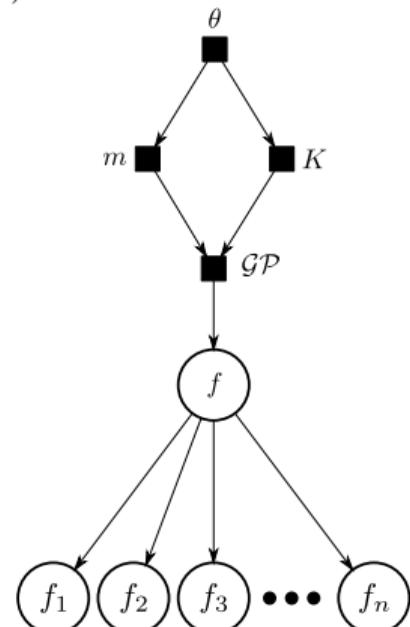


Sampling from the GP prior

- The (complete) sample of a GP $f \sim \mathcal{GP}$ is infinite dimensional
- However, using the marginalisation property we can draw (finite) *parts* of the GP sample
- Set up a grid and sample the *marginalisation* of f in that grid (an MVN)

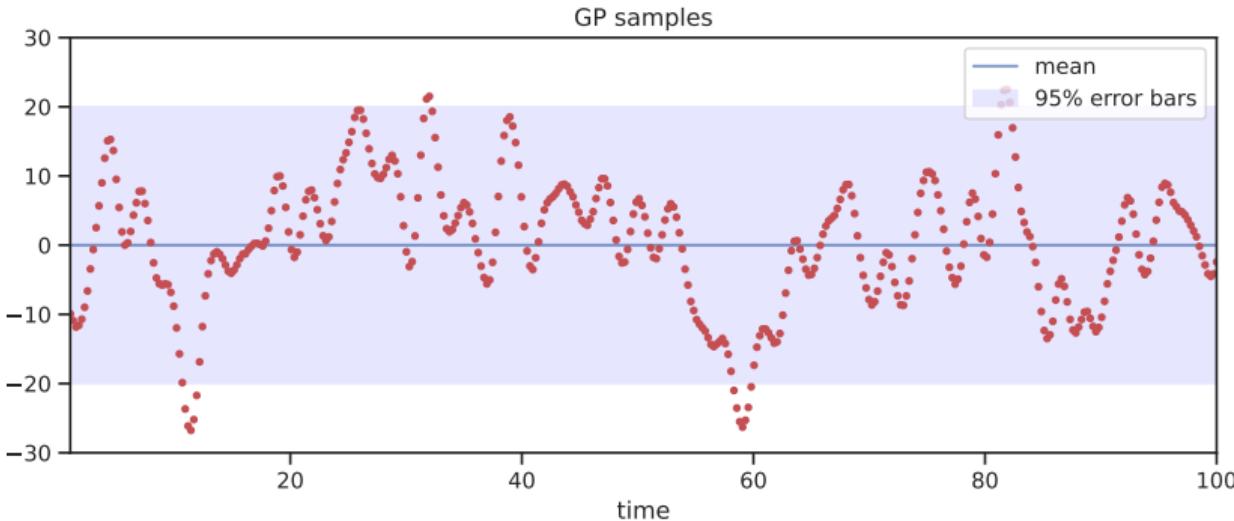


- Trick: use `linestyle = '-'` to make the grid look continuous :)

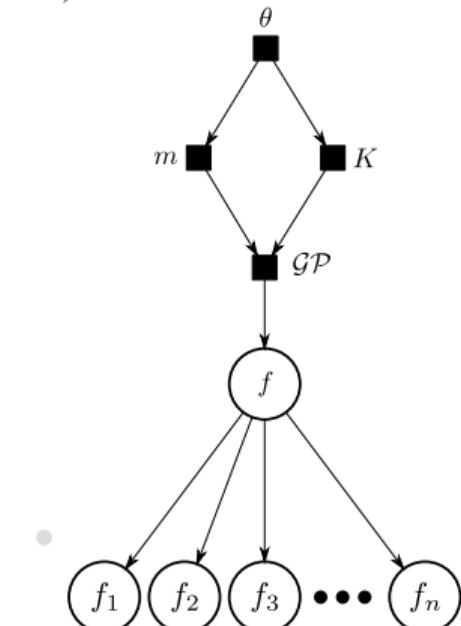


Sampling from the GP prior

- The (complete) sample of a GP $f \sim \mathcal{GP}$ is infinite dimensional
- However, using the marginalisation property we can draw (finite) *parts* of the GP sample
- Set up a grid and sample the *marginalisation* of f in that grid (an MVN)

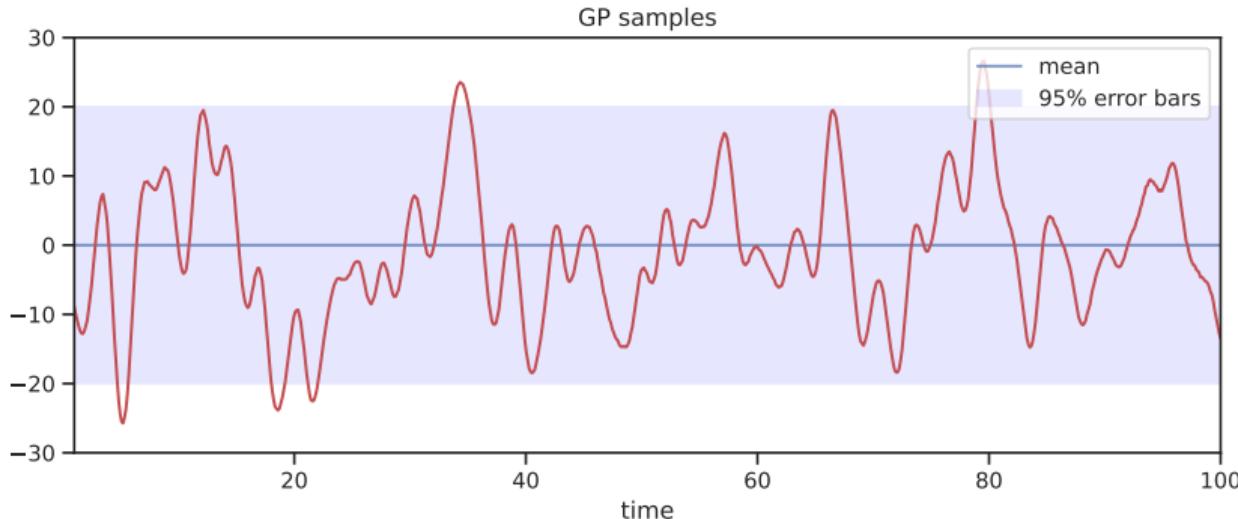


Trick: use `linestyle = '--'` to make the grid look continuous :)

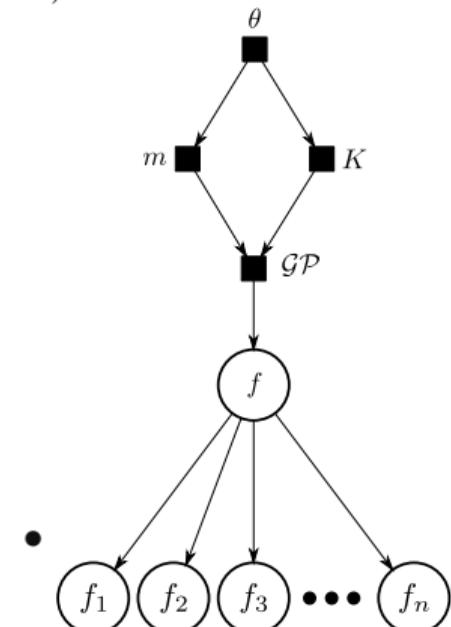


Sampling from the GP prior

- The (complete) sample of a GP $f \sim \mathcal{GP}$ is infinite dimensional
- However, using the marginalisation property we can draw (finite) *parts* of the GP sample
- Set up a grid and sample the *marginalisation* of f in that grid (an MVN)

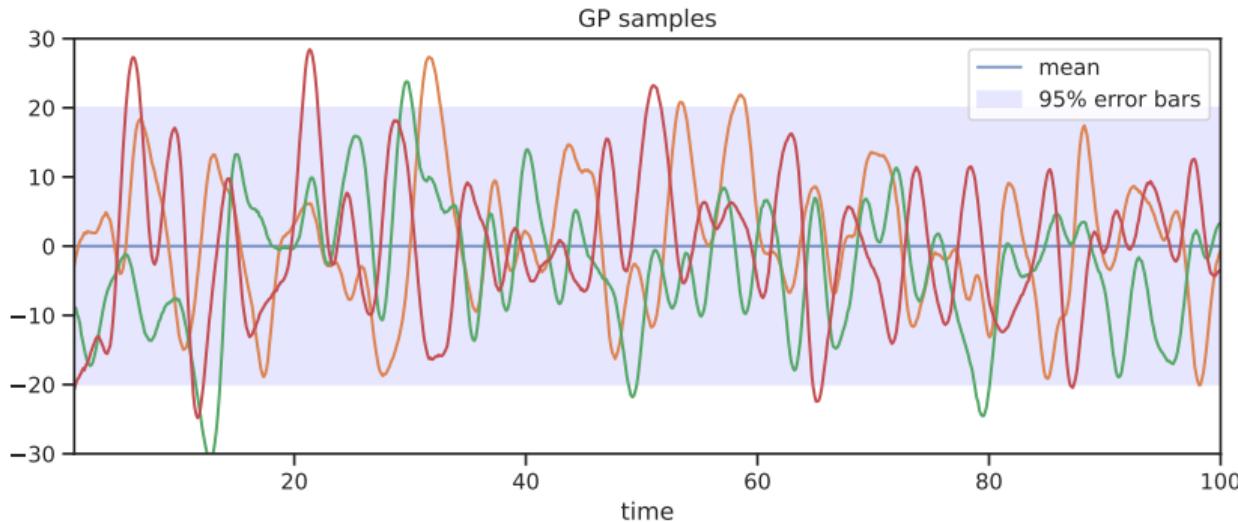


Trick: use `linestyle = '-'` to make the grid look continuous :)

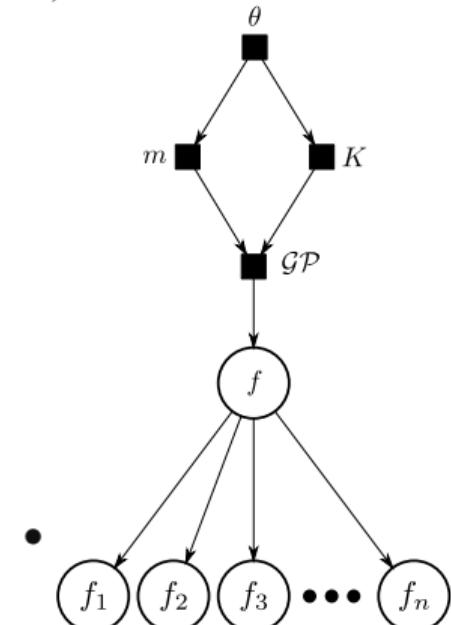


Sampling from the GP prior

- The (complete) sample of a GP $f \sim \mathcal{GP}$ is infinite dimensional
- However, using the marginalisation property we can draw (finite) *parts* of the GP sample
- Set up a grid and sample the *marginalisation* of f in that grid (an MVN)



Trick: use `linestyle = '-'` to make the grid look continuous :)

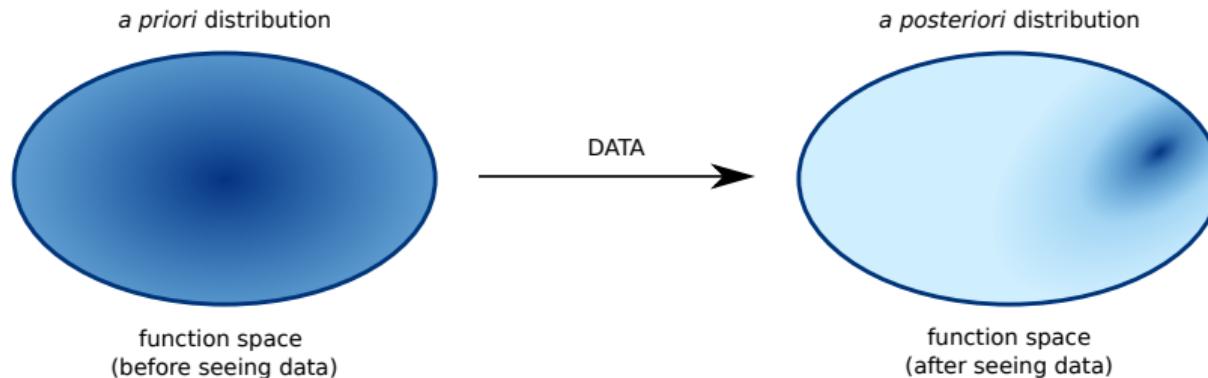


Incorporating data

Gaussians are also closed under conditioning. For data \mathbf{x} and \mathbf{y} , the posterior GP is:

$$f \sim \mathcal{GP}(\mu, K) \longrightarrow f|\mathbf{x}, \mathbf{y} \sim \mathcal{GP}\left(\mu + K_* \mathbf{K}^{-1} (\mathbf{y} - \mu(\mathbf{x})), K - K_* \mathbf{K}^{-1} K_*\right), \quad (11)$$

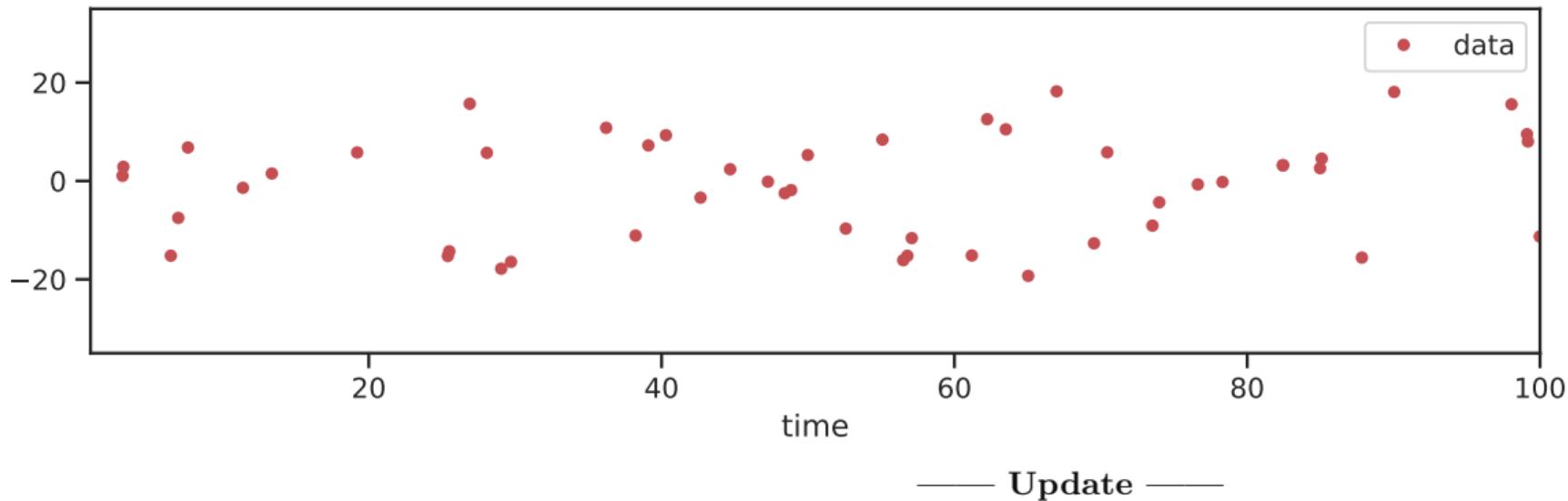
where $K_* = K(\cdot, \mathbf{x})$ and $\mathbf{K} = K(\mathbf{x}, \mathbf{x})$.



In the posterior GP, the mean and variance become *parametrised by the data*. This further reveals the **nonparametric** feature of the model: The more data we see, the more complex the model becomes.

Conditioning to observations (predictions)

Posterior



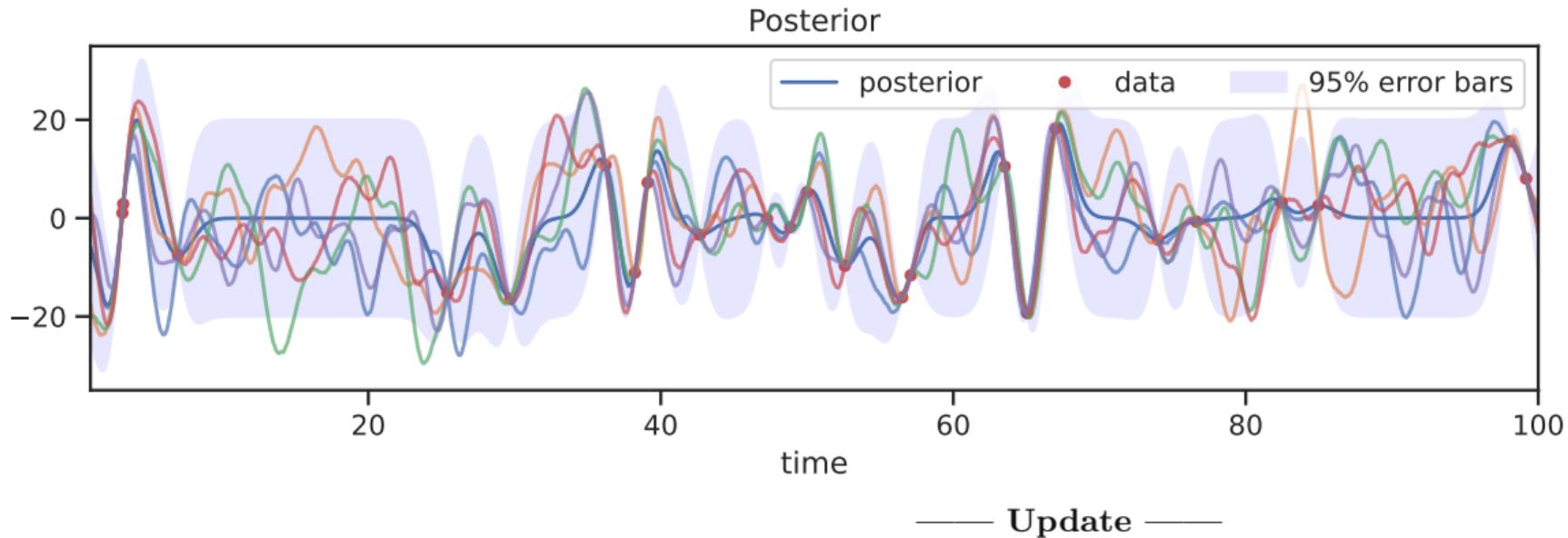
$$\mu \mapsto \mu + K_* \mathbf{K}^{-1} (\mathbf{y} - \mu(\mathbf{x})) \quad (12)$$

→ prediction is a linear combination of observed values

$$K \mapsto K - K_* \mathbf{K}^{-1} K_* \quad (13)$$

→ variance is reduced based on observed locations

Conditioning to observations (predictions)



$$\mu \mapsto \mu + K_* \mathbf{K}^{-1} (\mathbf{y} - \mu(\mathbf{x})) \quad (12)$$

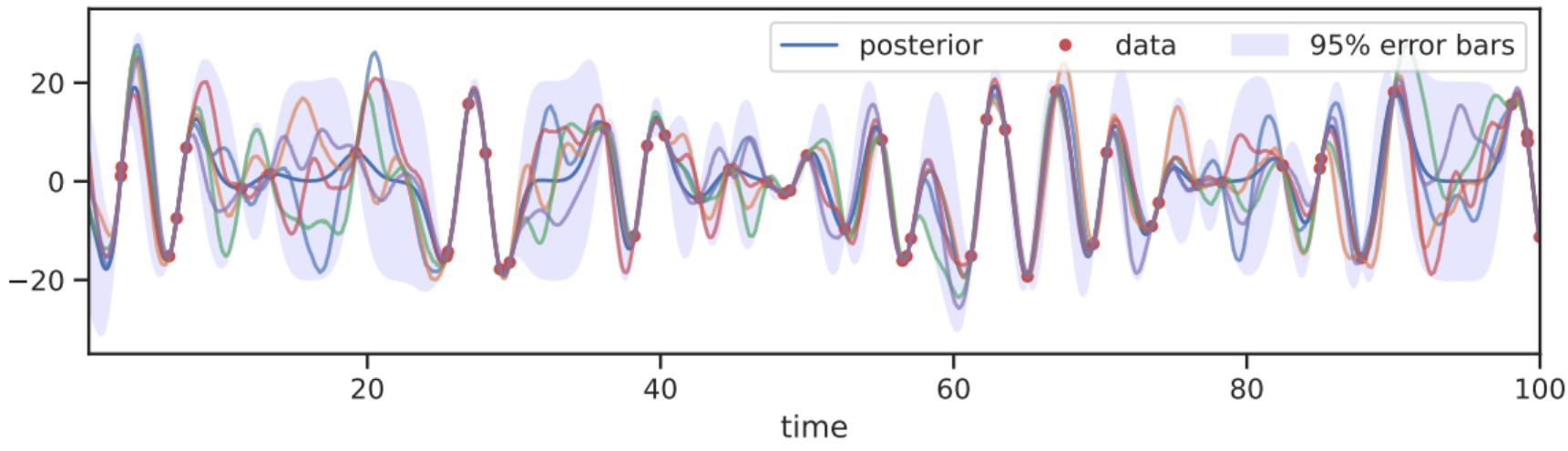
→ prediction is a linear combination of observed values

$$K \mapsto K - K_* \mathbf{K}^{-1} K_* \quad (13)$$

→ variance is reduced based on observed locations

Conditioning to observations (predictions)

Posterior



— Update —



$$\mu \mapsto \mu + K_* \mathbf{K}^{-1} (\mathbf{y} - \mu(\mathbf{x})) \quad (12)$$

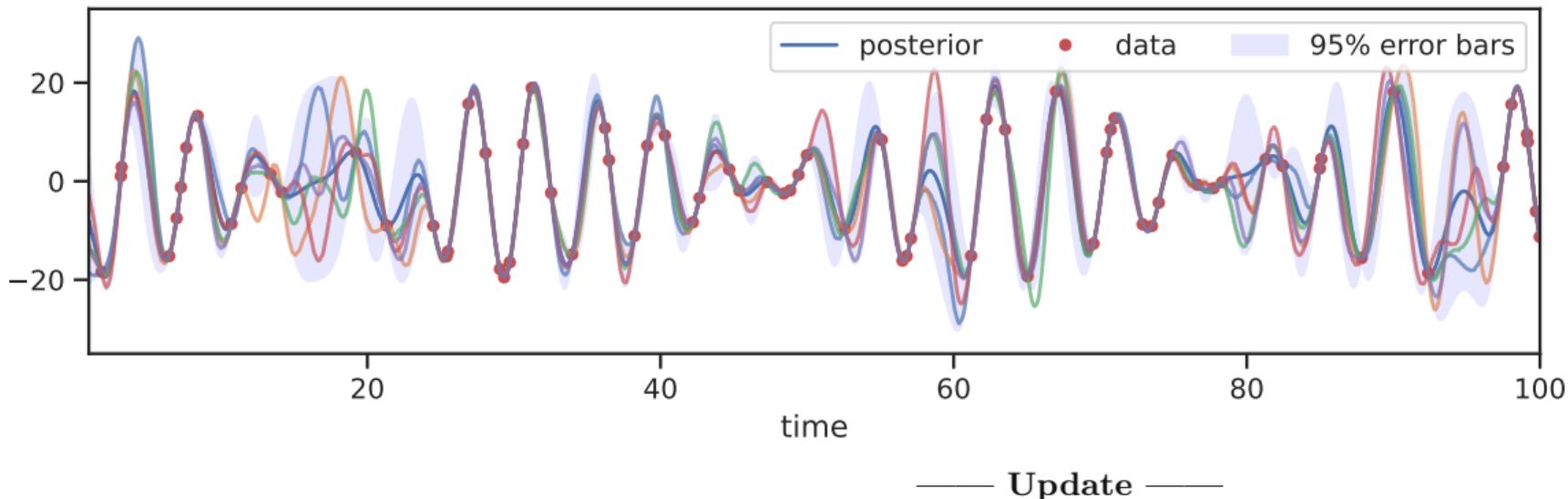
→ prediction is a linear combination of observed values

$$K \mapsto K - K_* \mathbf{K}^{-1} K_* \quad (13)$$

→ variance is reduced based on observed locations

Conditioning to observations (predictions)

Posterior



$$\mu \mapsto \mu + K_* \mathbf{K}^{-1} (\mathbf{y} - \mu(\mathbf{x})) \quad (12)$$

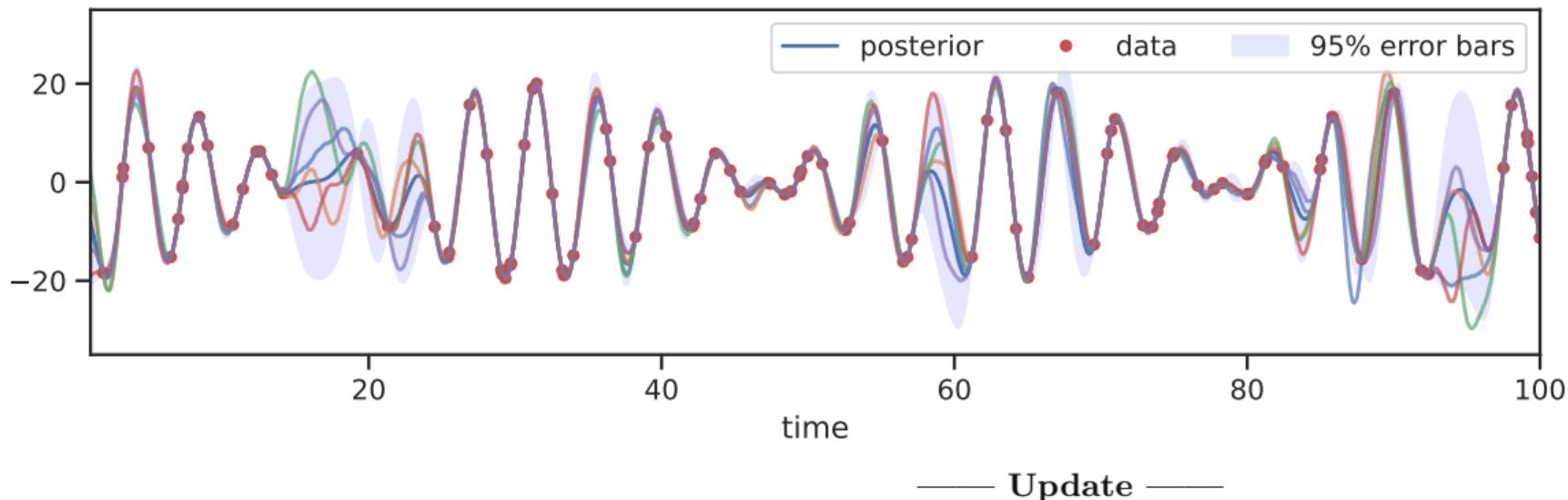
→ prediction is a linear combination of observed values

$$K \mapsto K - K_* \mathbf{K}^{-1} K_* \quad (13)$$

→ variance is reduced based on observed locations

Conditioning to observations (predictions)

Posterior



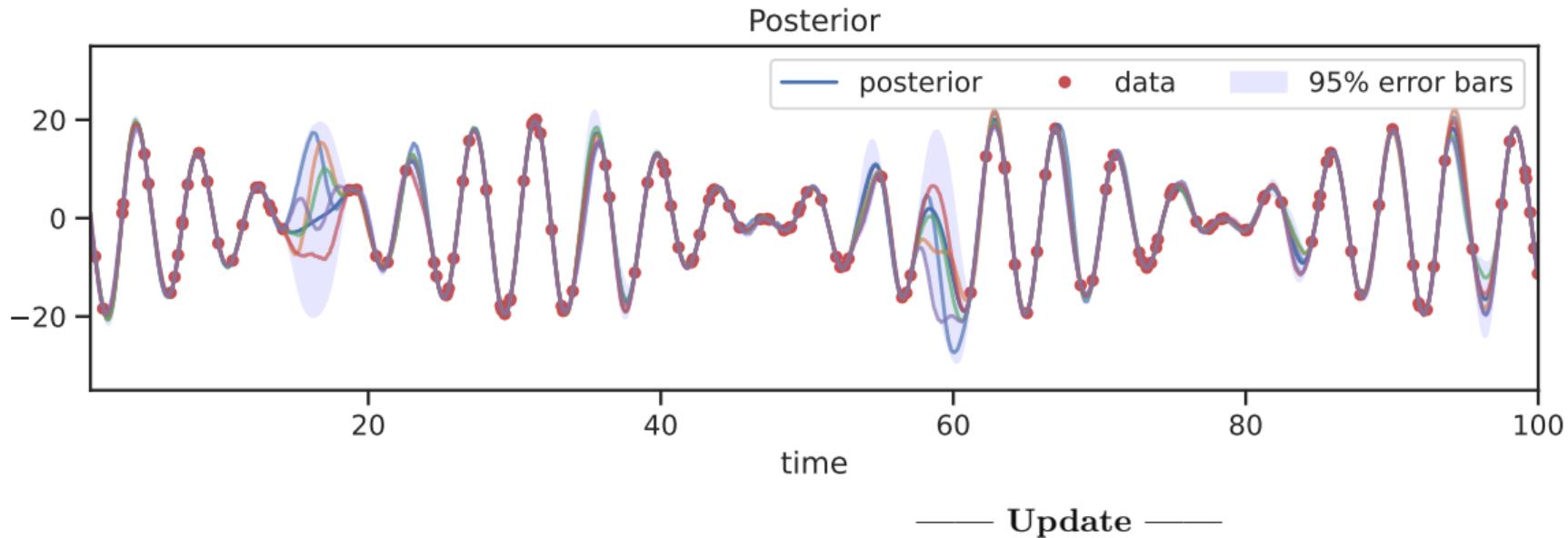
$$\mu \mapsto \mu + K_* \mathbf{K}^{-1} (\mathbf{y} - \mu(\mathbf{x})) \quad (12)$$

→ prediction is a linear combination of observed values

$$K \mapsto K - K_* \mathbf{K}^{-1} K_* \quad (13)$$

→ variance is reduced based on observed locations

Conditioning to observations (predictions)



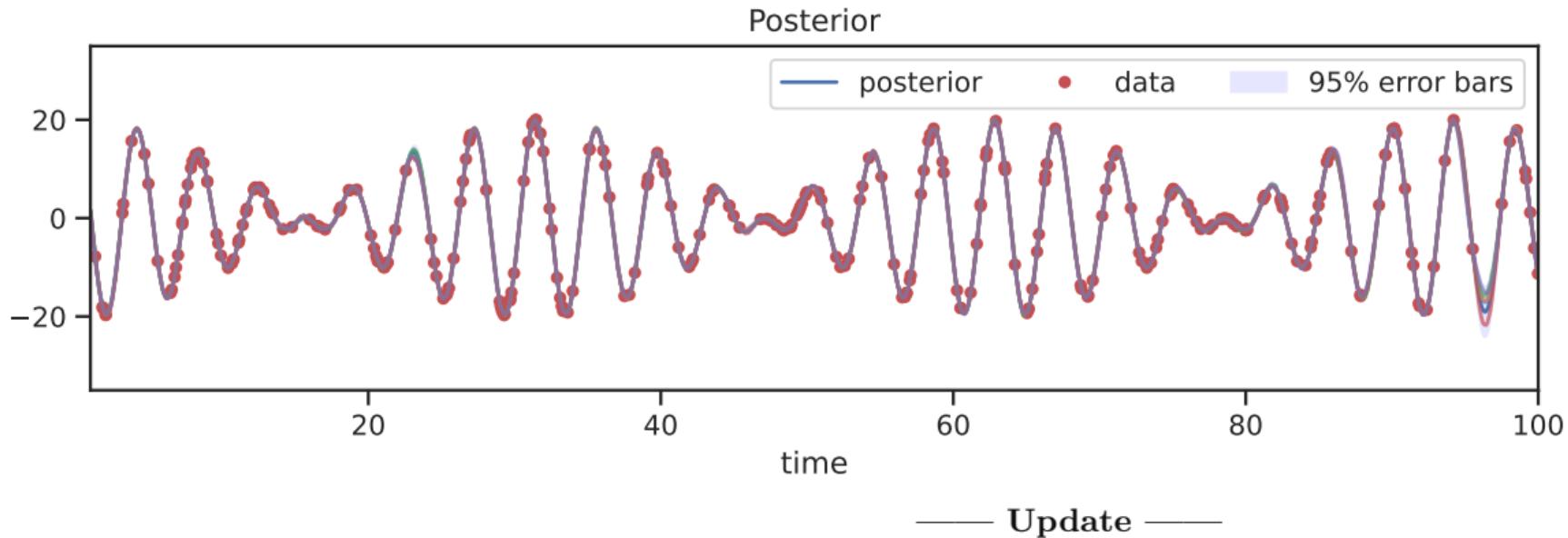
$$\mu \mapsto \mu + K_* \mathbf{K}^{-1} (\mathbf{y} - \mu(\mathbf{x})) \quad (12)$$

→ prediction is a linear combination of observed values

$$K \mapsto K - K_* \mathbf{K}^{-1} K_* \quad (13)$$

→ variance is reduced based on observed locations

Conditioning to observations (predictions)



$$\mu \mapsto \mu + K_* \mathbf{K}^{-1} (\mathbf{y} - \mu(\mathbf{x})) \quad (12)$$

→ prediction is a linear combination of observed values

$$K \mapsto K - K_* \mathbf{K}^{-1} K_* \quad (13)$$

→ variance is reduced based on observed locations

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

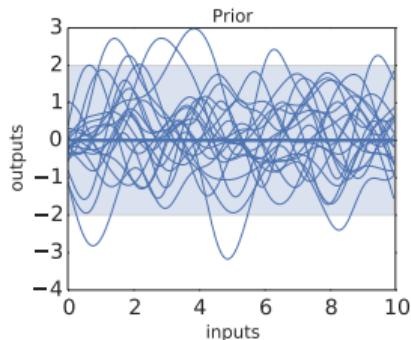
From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

Conditioning to real-world observations

Posterior predictive: $p(f(\mathbf{x})|\mathbf{x}, f(\mathbf{x}), \mathbf{x}) = \mathcal{N}\left(f(\mathbf{x}); m_{\mathbf{x}|\mathbf{x}}, \sigma_{\mathbf{x}|\mathbf{x}}^2\right)$ — (blue means new, red means observed)



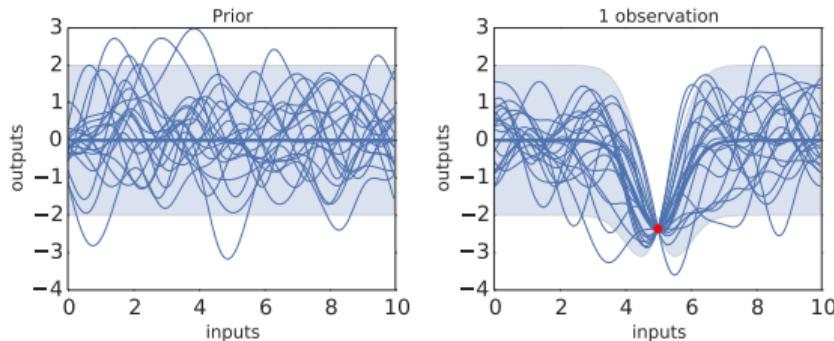
Generally, observations do not correspond to samples from a Gaussian process (or at least not one with the chosen kernel)



one needs to define an observation model that considers noise-corrupted GP observations

Conditioning to real-world observations

Posterior predictive: $p(f(\mathbf{x})|\mathbf{x}, f(\mathbf{x}), \mathbf{x}) = \mathcal{N}\left(f(\mathbf{x}); m_{\mathbf{x}|\mathbf{x}}, \sigma_{\mathbf{x}|\mathbf{x}}^2\right)$ — (blue means new, red means observed)



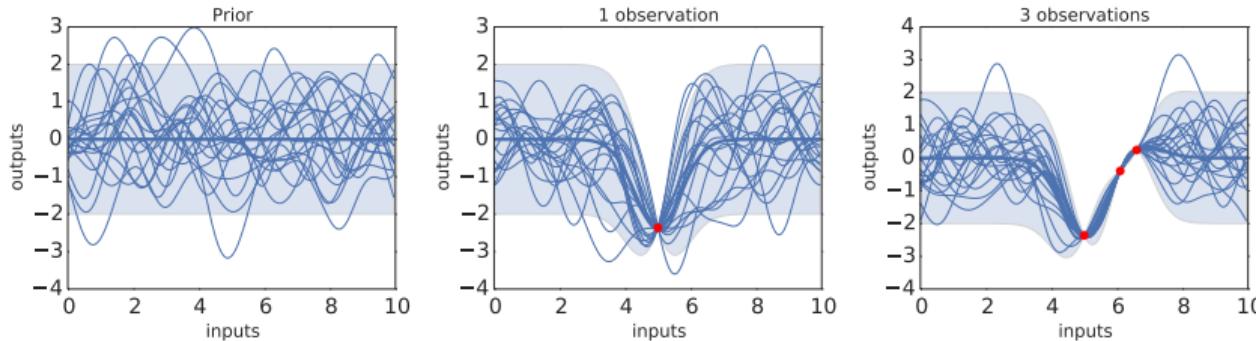
Generally, observations do not correspond to samples from a Gaussian process (or at least not one with the chosen kernel)



one needs to define an observation model that considers noise-corrupted GP observations

Conditioning to real-world observations

Posterior predictive: $p(f(\mathbf{x})|\mathbf{x}, f(\mathbf{x}), \mathbf{x}) = \mathcal{N}\left(f(\mathbf{x}); m_{\mathbf{x}|\mathbf{x}}, \sigma_{\mathbf{x}|\mathbf{x}}^2\right)$ — (blue means new, red means observed)



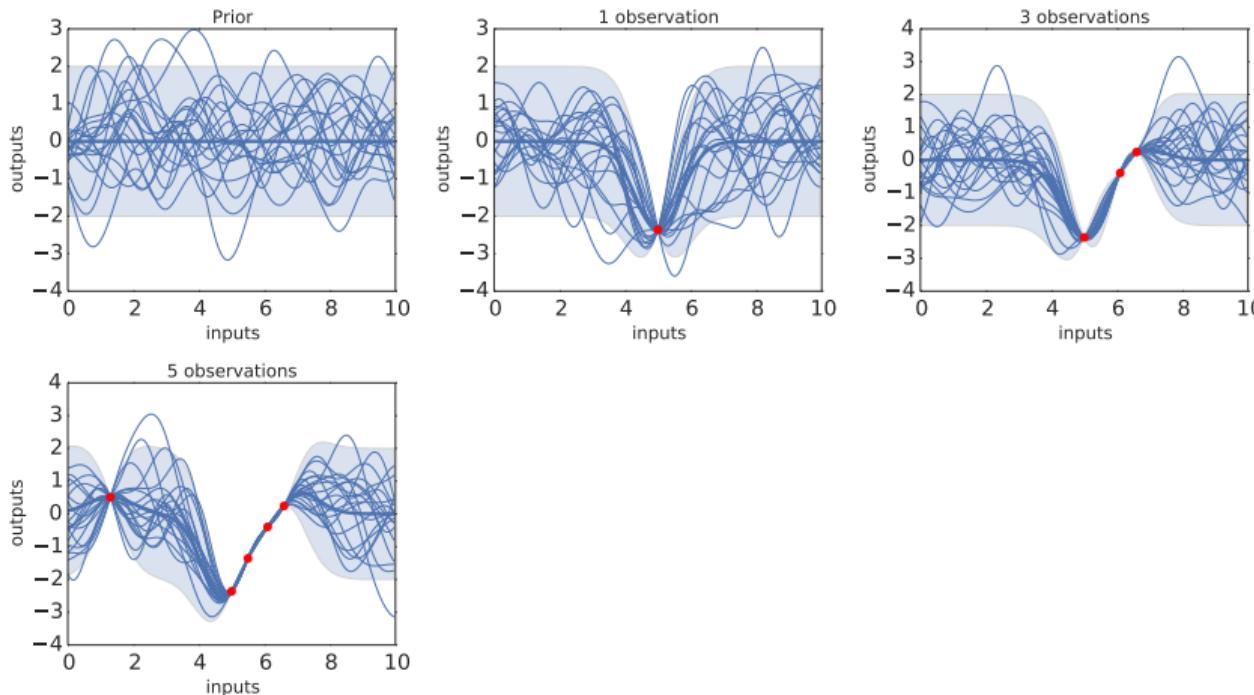
Generally, observations do not correspond to samples from a Gaussian process (or at least not one with the chosen kernel)



one needs to define an observation model that considers noise-corrupted GP observations

Conditioning to real-world observations

Posterior predictive: $p(f(\mathbf{x})|\mathbf{x}, f(\mathbf{x}), \mathbf{x}) = \mathcal{N}\left(f(\mathbf{x}); m_{\mathbf{x}|\mathbf{x}}, \sigma_{\mathbf{x}|\mathbf{x}}^2\right)$ — (blue means new, red means observed)



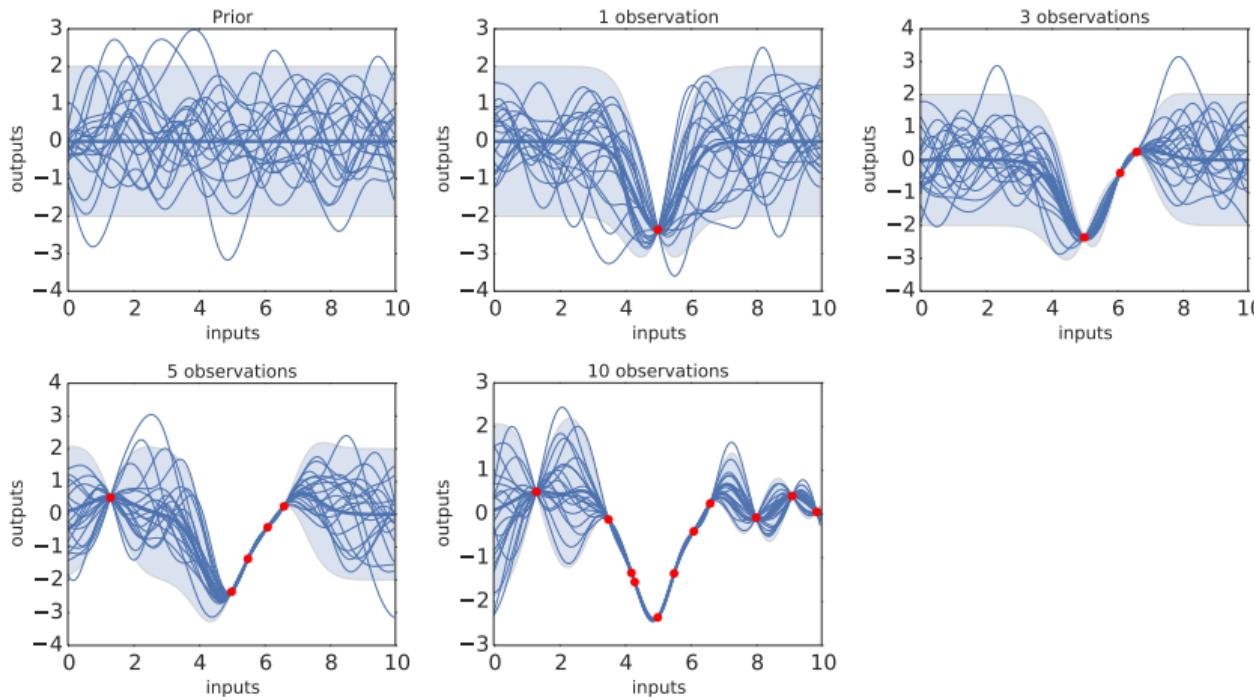
Generally, observations do not correspond to samples from a Gaussian process (or at least not one with the chosen kernel)



one needs to define an observation model that considers noise-corrupted GP observations

Conditioning to real-world observations

Posterior predictive: $p(f(\mathbf{x})|\mathbf{x}, f(\mathbf{x}), \mathbf{x}) = \mathcal{N}\left(f(\mathbf{x}); m_{\mathbf{x}|\mathbf{x}}, \sigma^2_{\mathbf{x}|\mathbf{x}}\right)$ — (blue means new, red means observed)



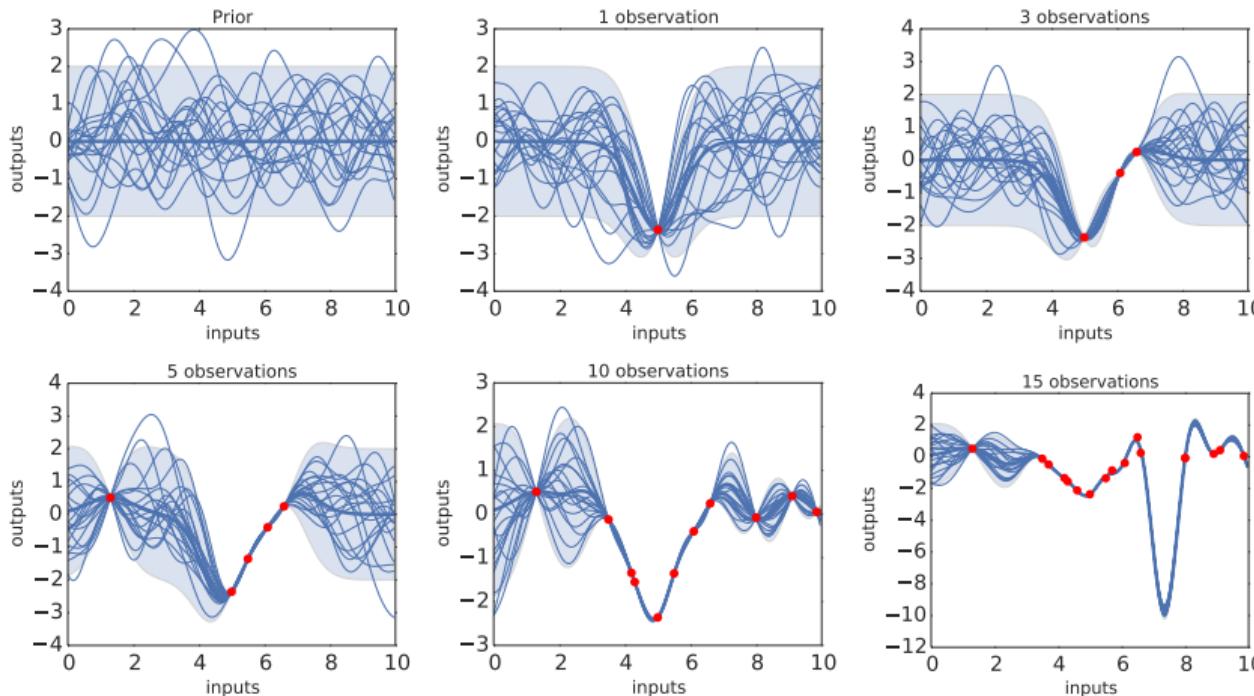
Generally, observations do not correspond to samples from a Gaussian process (or at least not one with the chosen kernel)



one needs to define an observation model that considers noise-corrupted GP observations

Conditioning to real-world observations

Posterior predictive: $p(f(\mathbf{x})|\mathbf{x}, f(\mathbf{x}), \mathbf{x}) = \mathcal{N}\left(f(\mathbf{x}); m_{\mathbf{x}|\mathbf{x}}, \sigma^2_{\mathbf{x}|\mathbf{x}}\right)$ — (blue means new, red means observed)



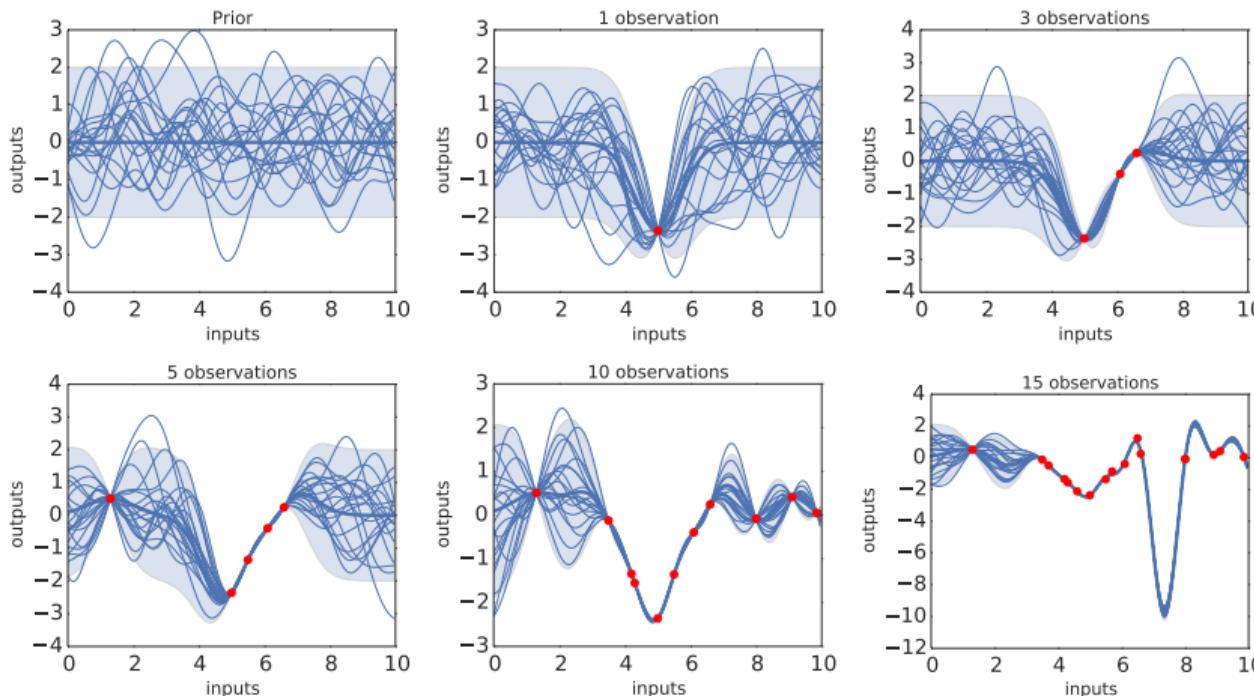
Generally, observations do not correspond to samples from a Gaussian process (or at least not one with the chosen kernel)



one needs to define an observation model that considers noise-corrupted GP observations

Conditioning to real-world observations

Posterior predictive: $p(f(\mathbf{x})|\mathbf{x}, f(\mathbf{x}), \mathbf{x}) = \mathcal{N}\left(f(\mathbf{x}); m_{\mathbf{x}|\mathbf{x}}, \sigma^2_{\mathbf{x}|\mathbf{x}}\right)$ — (blue means new, red means observed)



Generally, observations do not correspond to samples from a Gaussian process (or at least not one with the chosen kernel)

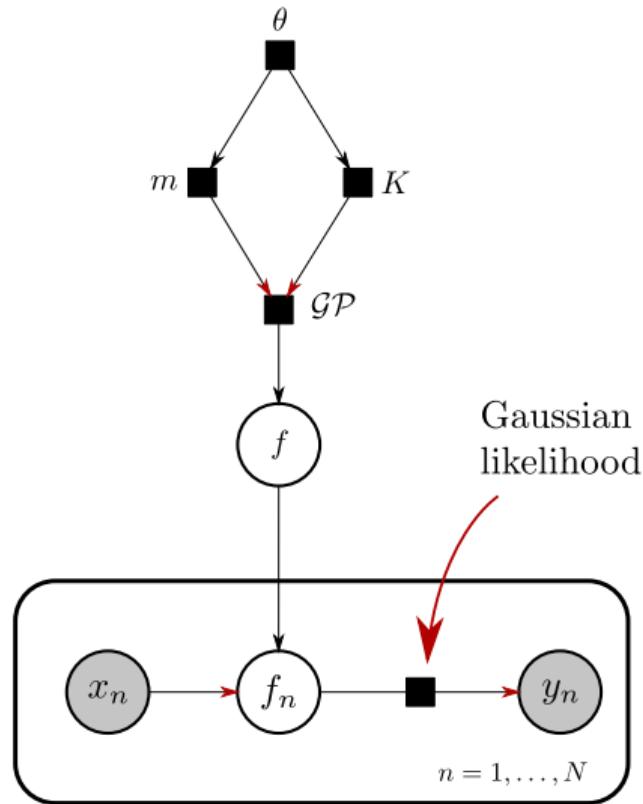


one needs to define an observation model that considers noise-corrupted GP observations

The Gaussian likelihood

linking real-world data to the GP generative model

- Modelling data (strictly) as a GP with a smooth kernel may result in mismodelling.
 - In practice, we use a likelihood function to connect the GP to the observations



Observation model with noise

To account for model mismatch, define $f(\cdot) \sim \mathcal{GP}(m(\cdot), K(\cdot, \cdot))$, with an observation

$$y = f(\mathbf{x}) + \eta, \quad \eta \sim \mathcal{N}(0, \sigma_{\text{noise}}^2). \quad (14)$$

The marginal distribution of the observation is **also a GP**:

$$y(\cdot) \sim \mathcal{GP}\left(m(\cdot), K(\cdot, \cdot) + I\sigma_{\text{noise}}^2\right). \quad (15)$$

Posterior is $\mathcal{GP}(\mu + K_\star \mathbf{K}^{-1}(\mathbf{y} - \mu(\mathbf{x})), K - K_\star \mathbf{K}^{-1} K_\star)$ — noise acts as a regulariser

Observation model with noise

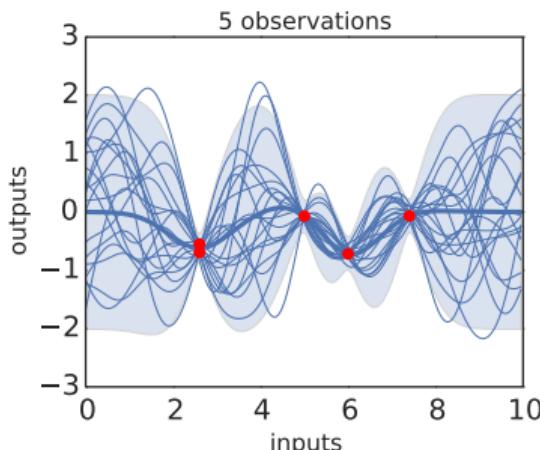
To account for model mismatch, define $f(\cdot) \sim \mathcal{GP}(m(\cdot), K(\cdot, \cdot))$, with an observation

$$y = f(\mathbf{x}) + \eta, \quad \eta \sim \mathcal{N}(0, \sigma_{\text{noise}}^2). \quad (14)$$

The marginal distribution of the observation is **also a GP**:

$$y(\cdot) \sim \mathcal{GP}\left(m(\cdot), K(\cdot, \cdot) + I\sigma_{\text{noise}}^2\right). \quad (15)$$

Posterior is $\mathcal{GP}(\mu + K_\star \mathbf{K}^{-1}(\mathbf{y} - \mu(\mathbf{x})), K - K_\star \mathbf{K}^{-1} K_\star)$ \rightarrow noise acts as a regulariser



Observation model with noise

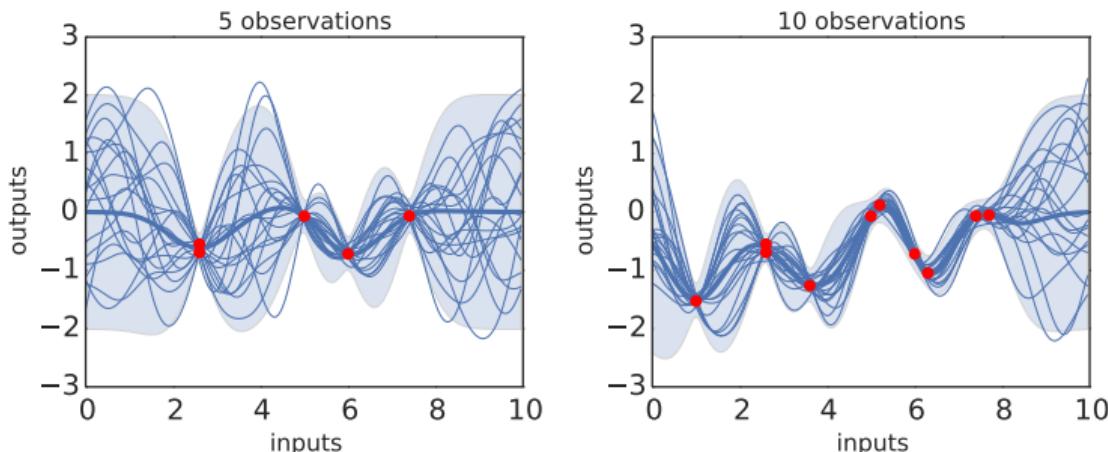
To account for model mismatch, define $f(\cdot) \sim \mathcal{GP}(m(\cdot), K(\cdot, \cdot))$, with an observation

$$y = f(\mathbf{x}) + \eta, \quad \eta \sim \mathcal{N}(0, \sigma_{\text{noise}}^2). \quad (14)$$

The marginal distribution of the observation is **also a GP**:

$$y(\cdot) \sim \mathcal{GP}\left(m(\cdot), K(\cdot, \cdot) + I\sigma_{\text{noise}}^2\right). \quad (15)$$

Posterior is $\mathcal{GP}\left(\mu + K_* \mathbf{K}^{-1}(\mathbf{y} - \mu(\mathbf{x})), K - K_* \mathbf{K}^{-1} K_*\right) \rightarrow$ noise acts as a regulariser



Observation model with noise

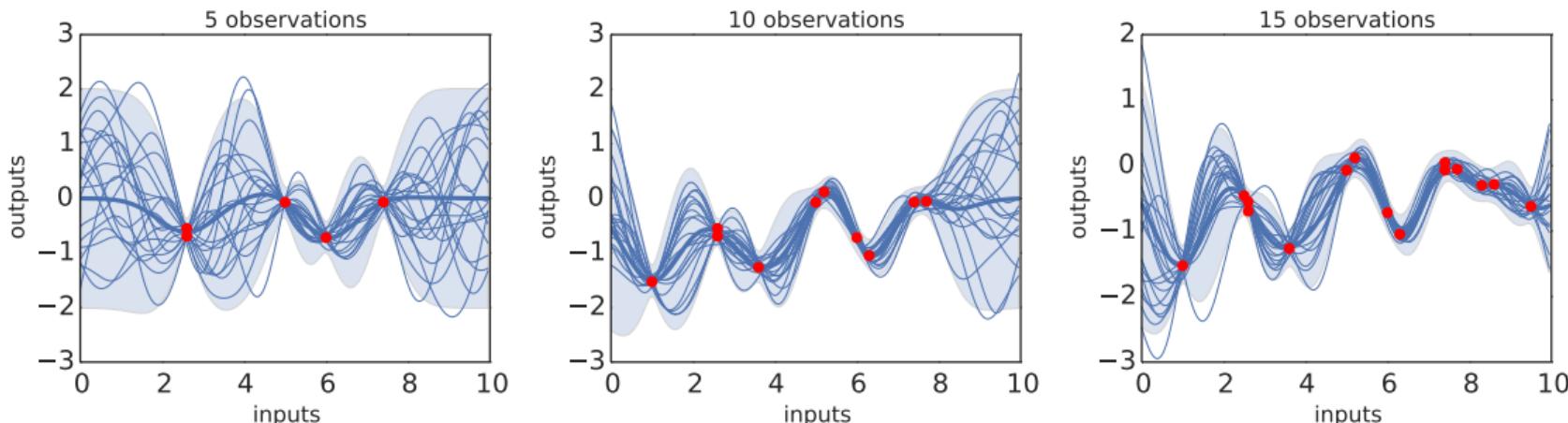
To account for model mismatch, define $f(\cdot) \sim \mathcal{GP}(m(\cdot), K(\cdot, \cdot))$, with an observation

$$y = f(\mathbf{x}) + \eta, \quad \eta \sim \mathcal{N}(0, \sigma_{\text{noise}}^2). \quad (14)$$

The marginal distribution of the observation is **also a GP**:

$$y(\cdot) \sim \mathcal{GP}\left(m(\cdot), K(\cdot, \cdot) + I\sigma_{\text{noise}}^2\right). \quad (15)$$

Posterior is $\mathcal{GP}(\mu + K_* \mathbf{K}^{-1}(\mathbf{y} - \mu(\mathbf{x})), K - K_* \mathbf{K}^{-1} K_*)$ \rightarrow noise acts as a regulariser



Learning: fitting the GP using data and maximum likelihood

- **Data:** $\mathbf{x} = \{x_i\}_{i=1:T} \subset \mathbb{R}^N$, $\mathbf{y} = \{y_i\}_{i=1:T} \subset \mathbb{R}$,
- **Model:** $y \sim \mathcal{GP}(m, K)$, where $K(x) = K_f(x) + \delta(x)\sigma_{\text{noise}}^2$ is the covariance of the (noisy) y

The marginal likelihood is given by:

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{(2\pi|\mathbf{K}|)^{T/2}} \exp\left(-\frac{1}{2}(\mathbf{y} - m(\mathbf{x}))^\top \mathbf{K}^{-1}(\mathbf{y} - m(\mathbf{x}))\right), \quad (16)$$

and the more optimisation-friendly **negative log-likelihood** NLL = $-\log p(\mathbf{y}|\mathbf{x})$ is:

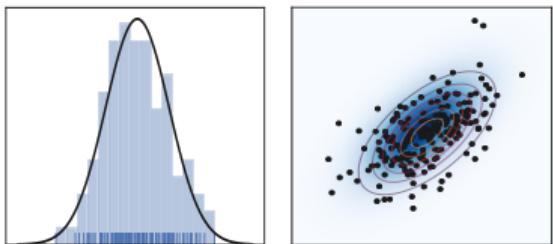
$$\text{NLL} = \frac{1}{2} \log |\mathbf{K}| + \frac{1}{2}(\mathbf{y} - m(\mathbf{x}))^\top \mathbf{K}^{-1}(\mathbf{y} - m(\mathbf{x})) + \frac{T}{2} \log 2\pi, \quad (17)$$

where $\mathbf{K} = K_y(\mathbf{x}, \mathbf{x})$ and $\mathbf{m} = m(\mathbf{x})$.

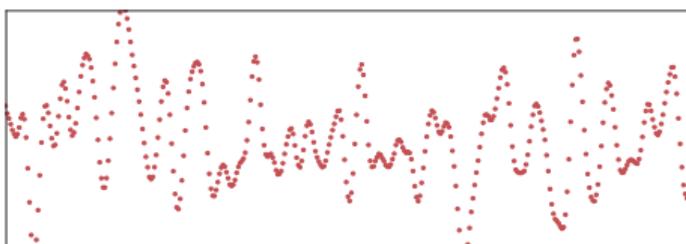
Learning a GP requires minimisation of the NLL with respect to the functions $K(\cdot, \cdot)$ and $m(\cdot)$. We avoid optimising on infinite-dimensional spaces by *making some compromises*.

Fitting Gaussians

Finite Gaussians: # samples >> parameter dimension



Gaussian process: # samples << 1 << param dimension

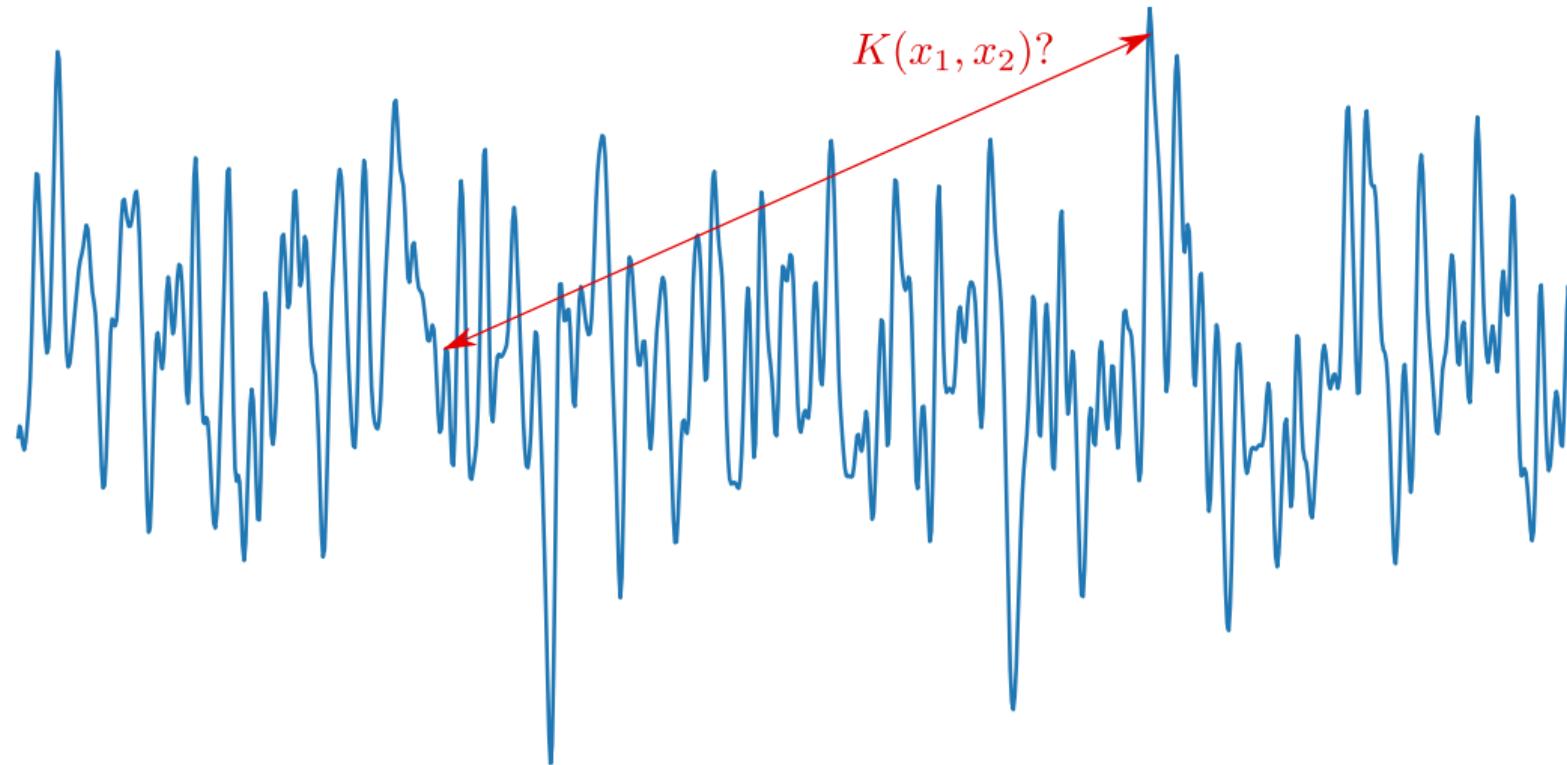


Learning (finite) Gaussians requires several observations: to fit an n -dimensional MVN we set $n + n(n + 1)/2$ parameters and thus require more data than that. This is why learning with parametric models is usually an **overdetermined** problem (e.g., linear regression).

With GPs, however, we don't even have a **single observation of the model**, but only have parts of it; therefore, learning a GP is an **underdetermined** problem.

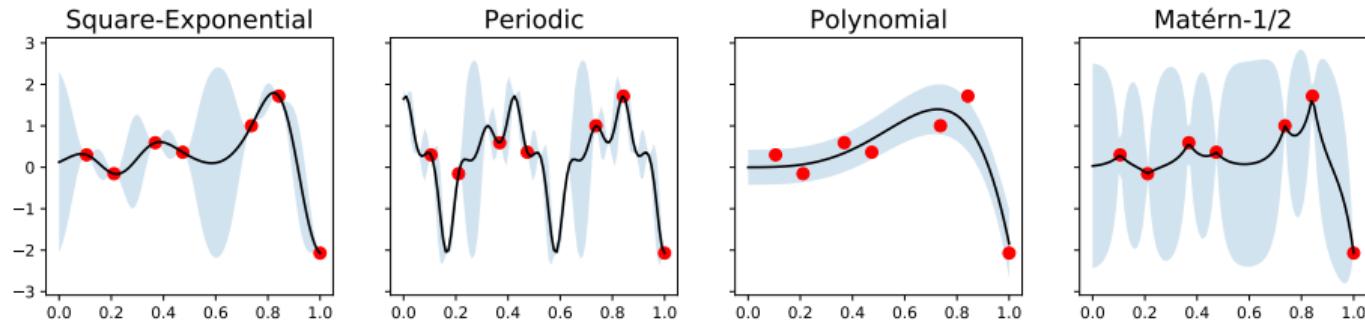
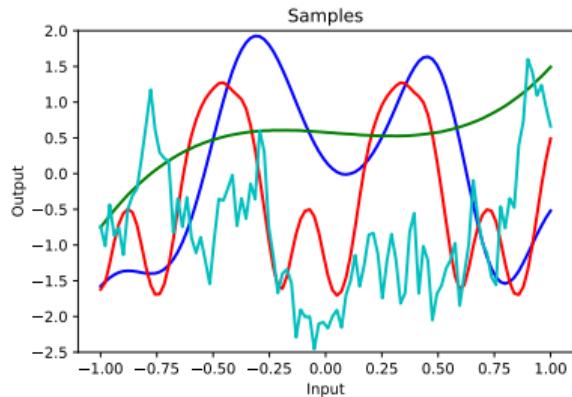
Thus, to learn a GP we need to *tie* the values of the mean and covariance to reduce their degrees of freedom. In practice, rather than freely assigning a covariance to each pair (x, y) we choose a function $(x, y) \mapsto K_\theta(x, y)$ and inspect alternatives for θ .

How to parametrise the kernel?



Different kernels, different predictions

- Square-Exponential: $k_{\text{SE}}(x, x') = \sigma^2 \exp\left(-\frac{(x-x')^2}{2l^2}\right)$
- Periodic: $k_{\text{Per}}(x, x') = \sigma^2 \exp\left(-\frac{2 \sin(\pi|x-x'|/p)}{l^2}\right)$
- Polynomial: $k_{\text{Pol}}(x, x') = (x^\top x' + c)^d$
- Matérn-1/2: $k_{\text{M}}(x, x') = \sigma^2 \exp\left(-\frac{|x-x'|}{2l^2}\right)$



Stationary kernels

Definition: A *stationary kernel* is function of the difference of its inputs

$$K(x, x') = K(x - x'). \quad (18)$$

Observation 1: GPs with stationary kernels are *invariant under translations*

$$K(x, x') = K(x - \delta, x' - \delta), \forall \delta \in \mathbb{R}. \quad (19)$$

Observation 2: The assumption of stationarity dramatically reduces kernel design.

Stationary kernels

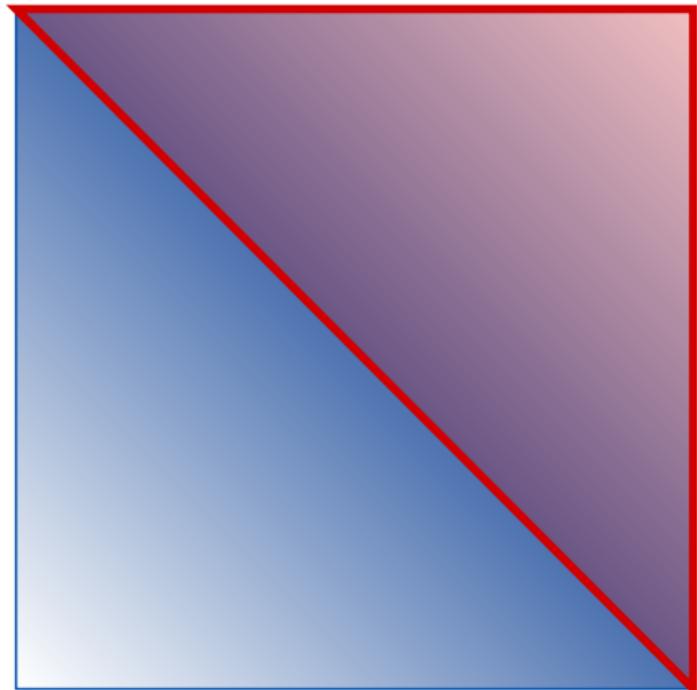
Definition: A *stationary kernel* is function of the difference of its inputs

$$K(x, x') = K(x - x'). \quad (18)$$

Observation 1: GPs with stationary kernels are *invariant under translations*

$$K(x, x') = K(x - \delta, x' - \delta), \forall \delta \in \mathbb{R}. \quad (19)$$

Observation 2: The assumption of stationarity dramatically reduces kernel design.



Stationary kernels

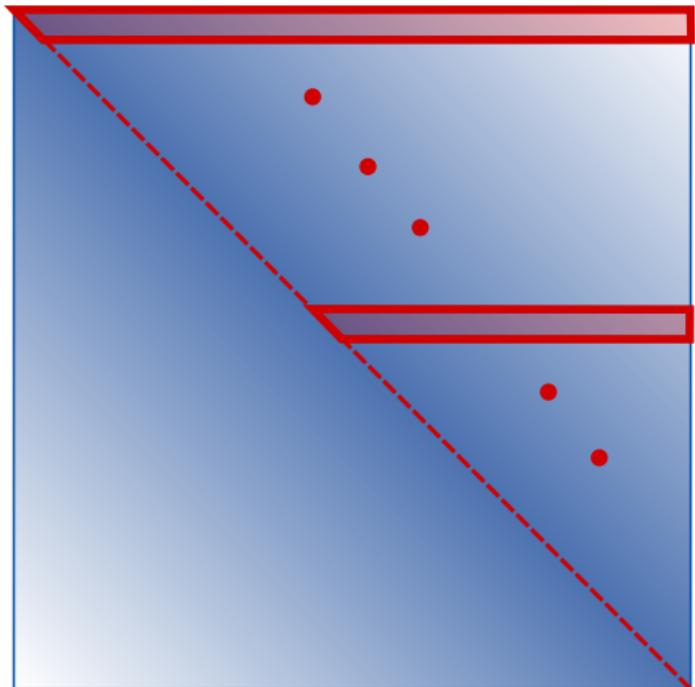
Definition: A *stationary kernel* is function of the difference of its inputs

$$K(x, x') = K(x - x'). \quad (18)$$

Observation 1: GPs with stationary kernels are *invariant under translations*

$$K(x, x') = K(x - \delta, x' - \delta), \forall \delta \in \mathbb{R}. \quad (19)$$

Observation 2: The assumption of stationarity dramatically reduces kernel design.



Stationary kernels: Spectral Mixture

Theorem (Bochner): A complex-valued function k on \mathbb{R}^d is the covariance function of a weakly-stationary mean-square-continuous stochastic process on \mathbb{R}^d if and only if it admits

$$k(x, x') = \int_{\mathbb{R}^d} e^{i\omega^\top (x-x')} S(\omega) d\omega, \quad (20)$$

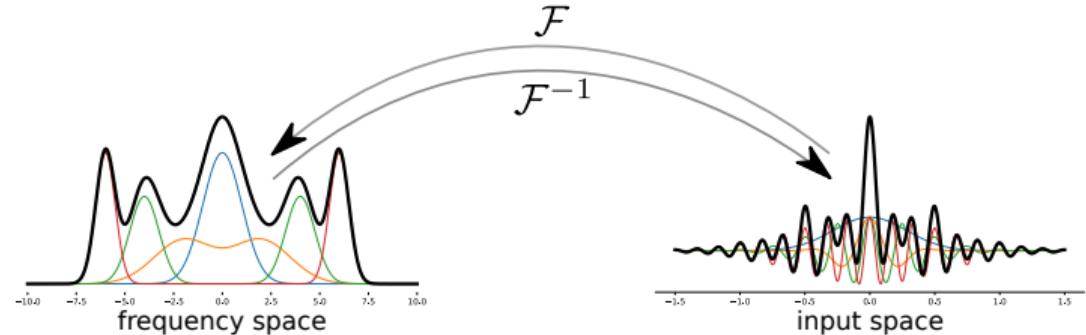
where $S(\omega)$ is a non-negative bounded function on \mathbb{R}^d , called the power spectral density.

Corollary: k and S are one-to-one \Rightarrow covariances can be designed in the frequency domain.

Spectral Mixture:

parametrise the power spectrum as a Gaussian mixture.

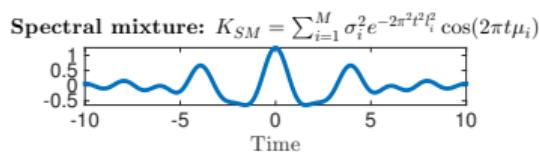
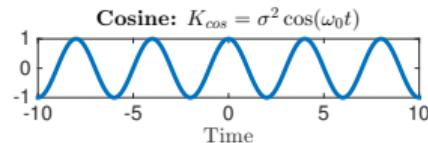
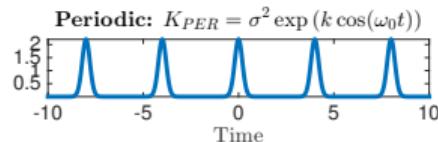
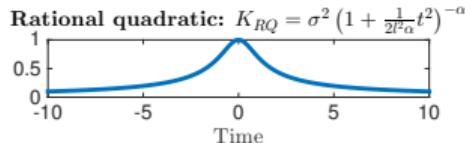
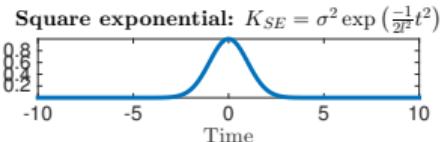
Wilson & Adams, GP kernels for pattern discovery and extrapolation. ICML, 2013



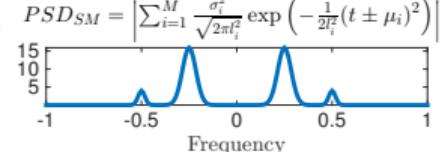
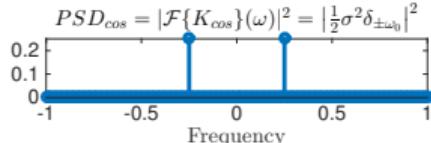
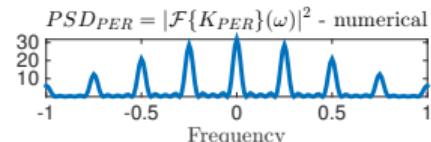
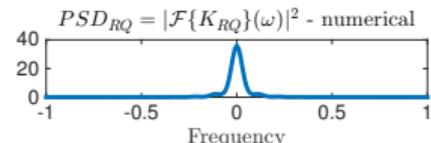
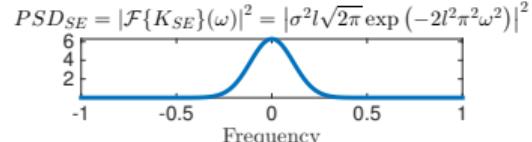
$$S(\xi) = \sum_{q=1}^Q \sigma_q^2 \frac{1}{\sqrt{2\pi l_q^2}} \exp\left(-\frac{1}{2l_q^2}(\xi \pm \mu_q)^2\right)$$

$$K(x) = \sum_{q=1}^Q \sigma_q^2 \exp(-2\pi^2 l_q^2 x^2) \cos(2\pi \mu_q x)$$

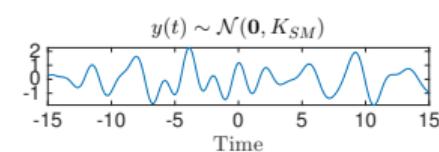
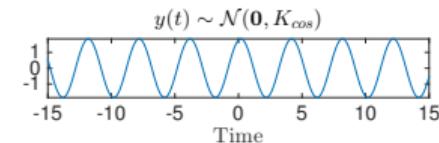
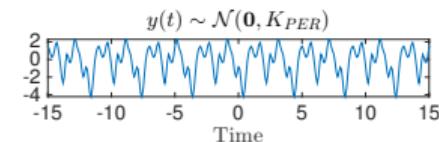
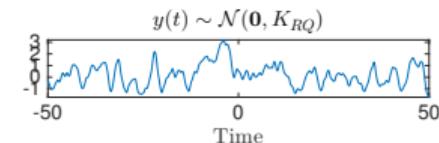
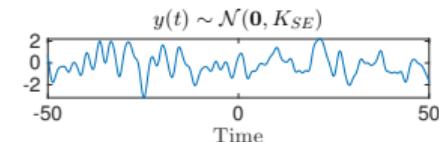
Kernels



Power spectral density



Sample function



After choosing the kernel

we optimise and predict

Training: $m, K = \arg \max p(\mathbf{y}|\mathbf{x})$ - **Data:** $\{(\mathbf{x}, \mathbf{y})\}$

$$\log p(\mathbf{y}|\mathbf{x}) = \log \mathcal{N}(\mathbf{y}; m(\mathbf{x}), K(\mathbf{x}, \mathbf{x})) \quad (21)$$

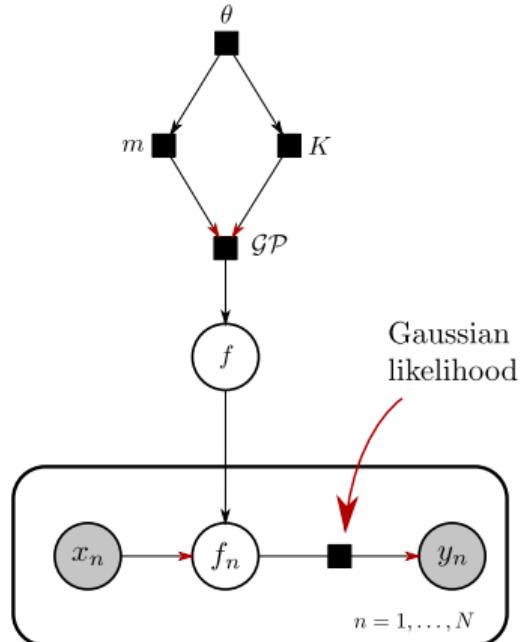
$$\propto -\log |K(\mathbf{x}, \mathbf{x})| - (\mathbf{y} - m(\mathbf{x}))^\top K(\mathbf{x}, \mathbf{x})^{-1} (\mathbf{y} - m(\mathbf{x})) \quad (22)$$

Prediction: $y_\star | x_\star, \mathbf{x}, \mathbf{y} \sim \mathcal{N}(m_p(x_\star), K_p(x_\star))$

$$m_p(x_\star) = K(x_\star, \mathbf{x}) K(\mathbf{x}, \mathbf{x})^{-1} (\mathbf{y} - m(\mathbf{x})) \quad (23)$$

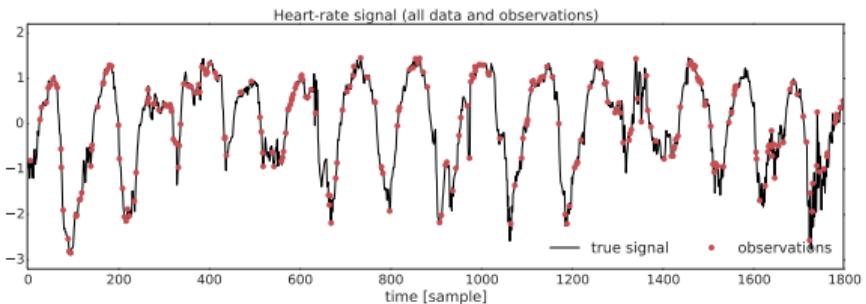
$$K_p(x_\star) = K(x_\star, \mathbf{x}) K(\mathbf{x}, \mathbf{x})^{-1} K(\mathbf{x}, x_\star) \quad (24)$$

how are predictions reported?



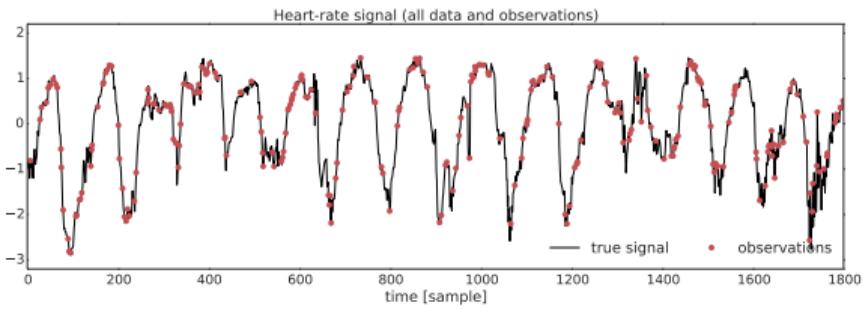
Example: Heart-rate signal

- 1800 data points
- 270 observations
- reconstruction+filtering



Example: Heart-rate signal

- 1800 data points
- 270 observations
- reconstruction+filtering

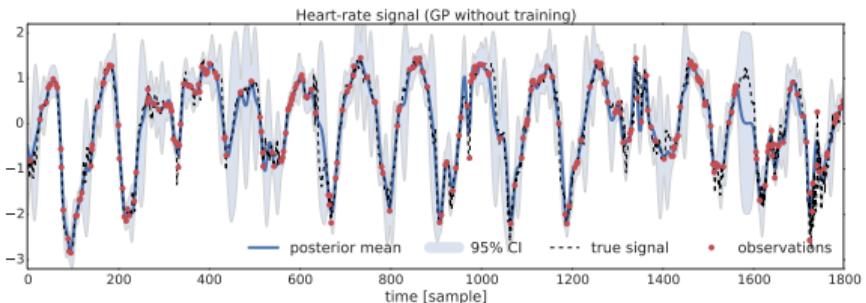


Model (not trained)

$$f(\cdot) \sim \mathcal{GP}(0, K), \quad \eta(t) \sim \mathcal{N}(0, \sigma_{\text{noise}}^2)$$

$$y(t) = f(t) + \eta(t)$$

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



Example: Heart-rate signal

- 1800 data points
- 270 observations
- reconstruction+filtering

Model (not trained)

$$f(\cdot) \sim \mathcal{GP}(0, K), \quad \eta(t) \sim \mathcal{N}(0, \sigma_{\text{noise}}^2)$$

$$y(t) = f(t) + \eta(t)$$

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

Model (trained)

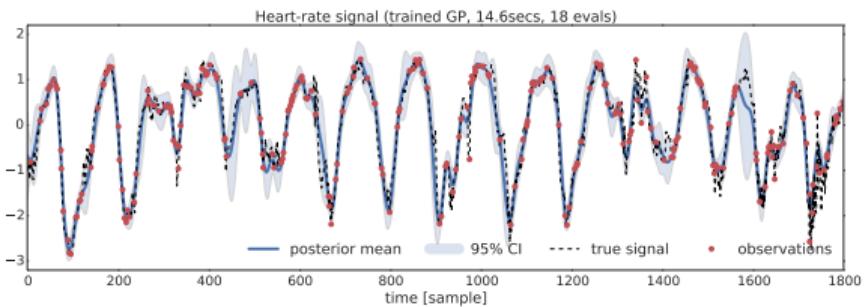
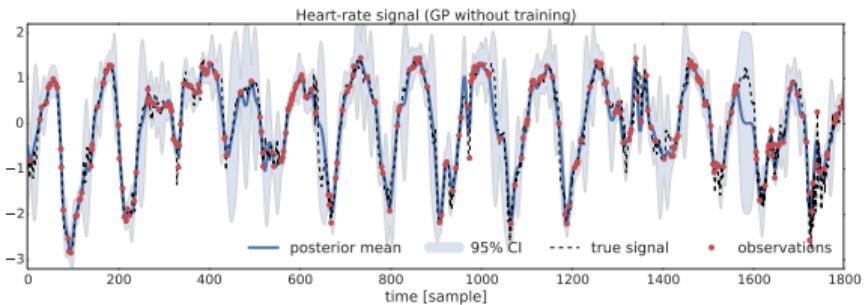
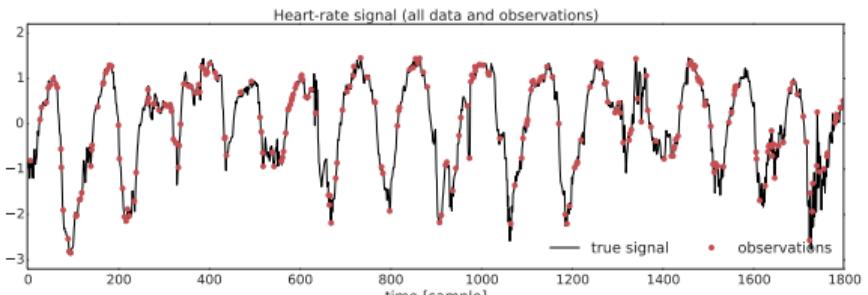
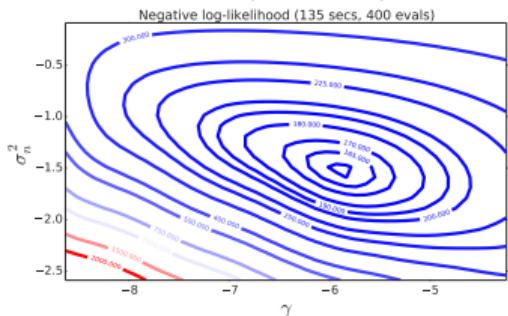


Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

GPlite - A Python package

- **Objective:** Show how to produce a minimal GP toolbox (< 200 lines)
- **Paradigm:** Object oriented, we define objects and their methods
- **Modules:**
 - constructor (initialisation)
 - sampler
 - data loader
 - posterior computation
 - likelihood & training
 - various plot functions
- **Demo**

Constructor

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.optimize import minimize
4 from utils import *
5
6 class gp_lite:
7
8     def __init__(self, kernel = 'SE'):
9         self.kernel = kernel
10
11    def init_hypers(self, case = 'canonical', theta = None):
12        self.x = None
13        self.y = None
14        # for the SE kernel
15        if case == 'canonical':
16            self.sigma = 10
17            self.gamma = 1/2
18            self.mu = 0
19            self.sigma_n = 0.1
20        elif case == 'manual':
21            self.sigma = theta[0]
22            self.gamma = theta[1]
23            self.mu = theta[2]
24            self.sigma_n = theta[3]
```

Use:

```
1 gp = gp_lite()
2 gp.init_hypers()
```

Data loading & sampler

```
1 def load(self, x, y):
2     self.Nx = len(x)
3     self.x = x
4     self.y = y
5
6 def sample(self, how_many=1):
7     samples = MVN(self.mean, self.cov, size=how_many)
8     self.samples = samples.T
9     return self.samples
```

Use:

```
1 gp.load(x_sample,y_sample)
2 gp.plot_data()
3 gp.sample(how_many=1)
```

Posterior computation

```
1 def compute_posterior(self, times):
2     if np.size(times) == 1:
3         self.N = dimension
4         self.time = np.linspace(1,100,dimension)
5     else:
6         self.time = where
7         self.N = len(where)
8
9     cov_grid = Spec_Mix(self.time,self.time, self.gamma, self.mu, self.sigma) +
10    1e-5*np.eye(self.N) + self.sigma_n**2*np.eye(self.N)
11
12    cov_obs = Spec_Mix(self.x,self.x, self.gamma, self.mu, self.sigma) +
13    1e-5*np.eye(self.Nx) + self.sigma_n**2*np.eye(self.Nx)
14
15    cov_star = Spec_Mix(self.time,self.x, self.gamma, self.mu, self.sigma)
16
17    self.mean = np.squeeze(cov_star@np.linalg.solve(cov_obs, self.y))
18
19    self.cov = cov_grid - (cov_star@np.linalg.solve(cov_obs, cov_star.T))
```

Use:

```
1 gp.compute_posterior(times=1000)
```

Likelihood & training

```
1 def nlogp(self, hypers):
2     sigma = np.exp(hypers[0])
3     gamma = np.exp(hypers[1])
4     mu = np.exp(hypers[2])
5     sigma_n = np.exp(hypers[3])
6
7     Y = self.y
8     Gram = Spec_Mix(self.x, self.x, gamma, mu, sigma)
9     K = Gram + sigma_n**2*np.eye(self.Nx) + 1e-5*np.eye(self.Nx)
10    (sign, logdet) = np.linalg.slogdet(K)
11    return 0.5*(Y.T@np.linalg.solve(K, Y) + logdet + self.Nx*np.log(2*np.pi))
12
13 def train(self, flag = 'quiet'):
14     hypers0 = np.array([np.log(self.sigma), np.log(self.gamma), np.log(self.mu),
15                         np.log(self.sigma_n)])
16     res = minimize(self.nlogp, hypers0, args=(), method='L-BFGS-B', jac =
17                     self.dnlogp, options={'maxiter': 500, 'disp': True})
18
19     self.sigma = np.exp(res.x[0])
20     self.gamma = np.exp(res.x[1])
21     self.mu = np.exp(res.x[2])
22     self.sigma_n = np.exp(res.x[3])
```

Use:

```
1 gp.train()
```

“Live” Demo :)

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

Beyond Gaussian likelihood

- Gaussian likelihood is suitable only when Gaussian observations are expected:

$$\begin{aligned}\mathbf{y} &= \mathbf{f} + \boldsymbol{\epsilon}, \\ p(\mathbf{f}|\mathbf{X}) &\sim \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_y^2 \mathbf{I}), \\ p(\mathbf{y}|\mathbf{f}) &= \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}).\end{aligned}$$

- For non-Gaussian data, we choose other likelihood $p(\mathbf{y}|\mathbf{f}, \phi)$ parameterised by ϕ .
 - Heavy-tailed regression: e.g. Student-t or Laplace likelihoods;
 - Non-negative regression: e.g. gamma or exponential likelihoods;
 - Discrete regression: e.g. Poisson likelihood;
 - Binary classification: e.g. Bernoulli likelihood;
 - Multiclass classification: e.g. categorical likelihood.
- **Problem:** The tractable integrals of GP models require a Gaussian likelihood.

Beyond Gaussian likelihood

- Gaussian likelihood is suitable only when Gaussian observations are expected:

$$\begin{aligned}\mathbf{y} &= \mathbf{f} + \boldsymbol{\epsilon}, \\ p(\mathbf{f}|\mathbf{X}) &\sim \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_y^2 \mathbf{I}), \\ p(\mathbf{y}|\mathbf{f}) &= \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}).\end{aligned}$$

- For non-Gaussian data, we choose other likelihood $p(\mathbf{y}|\mathbf{f}, \boldsymbol{\phi})$ parameterised by $\boldsymbol{\phi}$.
 - Heavy-tailed regression: e.g. **Student-t** or **Laplace** likelihoods;
 - Non-negative regression: e.g. **gamma** or **exponential** likelihoods;
 - Discrete regression: e.g. **Poisson** likelihood;
 - Binary classification: e.g. **Bernoulli** likelihood;
 - Multiclass classification: e.g. **categorical** likelihood.
- Problem: The tractable integrals of GP models require a Gaussian likelihood.

Beyond Gaussian likelihood

- Gaussian likelihood is suitable only when Gaussian observations are expected:

$$\begin{aligned}\mathbf{y} &= \mathbf{f} + \boldsymbol{\epsilon}, \\ p(\mathbf{f}|\mathbf{X}) &\sim \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_y^2 \mathbf{I}), \\ p(\mathbf{y}|\mathbf{f}) &= \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_y^2 \mathbf{I}).\end{aligned}$$

- For non-Gaussian data, we choose other likelihood $p(\mathbf{y}|\mathbf{f}, \boldsymbol{\phi})$ parameterised by $\boldsymbol{\phi}$.
 - Heavy-tailed regression: e.g. **Student-t** or **Laplace** likelihoods;
 - Non-negative regression: e.g. **gamma** or **exponential** likelihoods;
 - Discrete regression: e.g. **Poisson** likelihood;
 - Binary classification: e.g. **Bernoulli** likelihood;
 - Multiclass classification: e.g. **categorical** likelihood.
- **Problem:** The tractable integrals of GP models require a Gaussian likelihood.

Beyond Gaussian likelihood

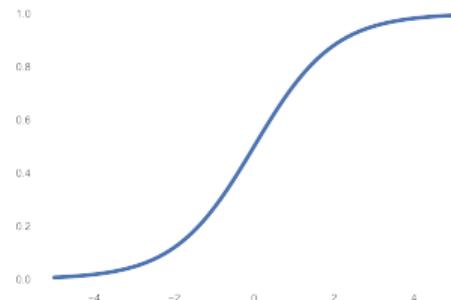
- Consider a **binary classification task** with Bernoulli observations:

$$p(y) = \mathcal{B}(y|a) = a^y(1 - a^{1-y}), \quad a = p(y=1).$$

- A Bernoulli likelihood $p(y|f)$ is obtained with $a = \sigma(f) = p(y=1|f)$, where $\sigma(f) \in [0, 1], \forall f \in \mathbb{R}$, is a **sigmoid function**:

$$p(y|f) = \sigma(f)^y(1 - \sigma(f))^{1-y},$$

$$p(\mathbf{y}|f) = \prod_{i=1}^N \sigma(f_i)^{y_i}(1 - \sigma(f_i))^{1-y_i}.$$



- $\sigma(f) = a$ is the **response (or transformation) function**.
- $\sigma^{-1}(a) = f$ is the **link function**, which must be monotonic.
- $f = \theta^\top x$ results in **logistic regression**, an example of a generalized linear model.

Beyond Gaussian likelihood

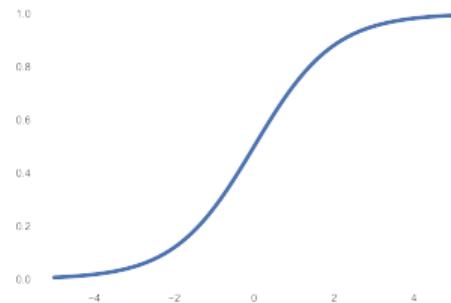
- Consider a **binary classification task** with Bernoulli observations:

$$p(y) = \mathcal{B}(y|a) = a^y(1 - a^{1-y}), \quad a = p(y=1).$$

- A Bernoulli likelihood $p(y|f)$ is obtained with $a = \sigma(f) = p(y=1|f)$, where $\sigma(f) \in [0, 1], \forall f \in \mathbb{R}$, is a **sigmoid function**:

$$p(y|f) = \sigma(f)^y(1 - \sigma(f)^{1-y}),$$

$$p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^N \sigma(f_i)^{y_i}(1 - \sigma(f_i)^{1-y_i}).$$



- $\sigma(f) = a$ is the **response** (or transformation) function.
- $\sigma^{-1}(a) = f$ is the **link function**, which must be monotonic.
- $f = \theta^\top x$ results in **logistic regression**, an example of a generalized linear model.

Beyond Gaussian likelihood

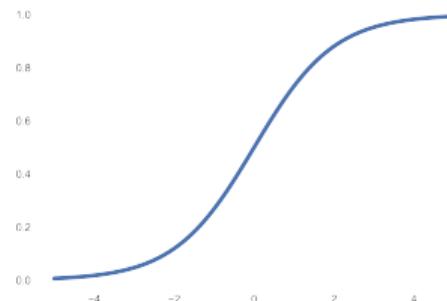
- Consider a **binary classification task** with Bernoulli observations:

$$p(y) = \mathcal{B}(y|a) = a^y(1 - a^{1-y}), \quad a = p(y=1).$$

- A Bernoulli likelihood $p(y|f)$ is obtained with $a = \sigma(f) = p(y=1|f)$, where $\sigma(f) \in [0, 1], \forall f \in \mathbb{R}$, is a **sigmoid function**:

$$p(y|f) = \sigma(f)^y(1 - \sigma(f)^{1-y}),$$

$$p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^N \sigma(f_i)^{y_i}(1 - \sigma(f_i)^{1-y_i}).$$



- $\sigma(f) = a$ is the **response (or transformation) function**.
- $\sigma^{-1}(a) = f$ is the **link function**, which must be monotonic.
- $f = \theta^\top x$ results in **logistic regression**, an example of a **generalized linear model**.

Beyond Gaussian likelihood

- GP models for binary classification can be similarly obtained:

$$p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^N \sigma(f_i)^{y_i} (1 - \sigma(f_i)^{1-y_i}), \quad \mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f).$$

- \mathbf{f} is a **nuisance variable** and we are interested in predictions in the observed space:

$$\underbrace{p(f_*|x_*, \mathbf{X}, \mathbf{y})}_{\text{latent space}} = \int p(f_*|\mathbf{f}, x_*, \mathbf{X}) p(\mathbf{f}|\mathbf{X}, \mathbf{y}) d\mathbf{f},$$

$$\underbrace{p(y_*|x_*, \mathbf{X}, \mathbf{y})}_{\text{observed space}} = \int p(y_*|f_*) p(f_*|x_*, \mathbf{X}, \mathbf{y}) df_*.$$

- For a binary classification task with Bernoulli likelihood:

$$p(y_*|x_*, \mathbf{X}, \mathbf{y}) = \int \overbrace{\sigma(f_*)^{y_*} (1 - \sigma(f_*)^{1-y_*})}^{p(y_*|f_*)} \overbrace{p(f_*|\mathbf{f}, x_*, \mathbf{X})}^{\text{cond. Gaussian}} \overbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}^{\text{posterior}} d\mathbf{f} df_*,$$

$$p(y_* = 1|x_*, \mathbf{X}, \mathbf{y}) = \int \sigma(f_*) p(f_*|\mathbf{f}, x_*, \mathbf{X}) p(\mathbf{f}|\mathbf{X}, \mathbf{y}) d\mathbf{f} df_*.$$

Beyond Gaussian likelihood

- GP models for binary classification can be similarly obtained:

$$p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^N \sigma(f_i)^{y_i} (1 - \sigma(f_i)^{1-y_i}), \quad \mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f).$$

- \mathbf{f} is a **nuisance variable** and we are interested in predictions in the observed space:

$$\underbrace{p(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y})}_{\text{latent space}} = \int p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f}|\mathbf{X}, \mathbf{y}) d\mathbf{f},$$

$$\underbrace{p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y})}_{\text{observed space}} = \int p(y_*|f_*) p(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) df_*.$$

- For a binary classification task with Bernoulli likelihood:

$$p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \overbrace{\sigma(f_*)^{y_*} (1 - \sigma(f_*)^{1-y_*})}^{p(y_*|f_*)} \overbrace{p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X})}^{\text{cond. Gaussian}} \overbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}^{\text{posterior}} df_* d\mathbf{f},$$

$$p(y_* = 1|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \sigma(f_*) p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f}|\mathbf{X}, \mathbf{y}) df_* d\mathbf{f}.$$

Beyond Gaussian likelihood

- GP models for binary classification can be similarly obtained:

$$p(\mathbf{y}|\mathbf{f}) = \prod_{i=1}^N \sigma(f_i)^{y_i} (1 - \sigma(f_i)^{1-y_i}), \quad \mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f).$$

- \mathbf{f} is a **nuisance variable** and we are interested in predictions in the observed space:

$$\underbrace{p(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y})}_{\text{latent space}} = \int p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f}|\mathbf{X}, \mathbf{y}) d\mathbf{f},$$

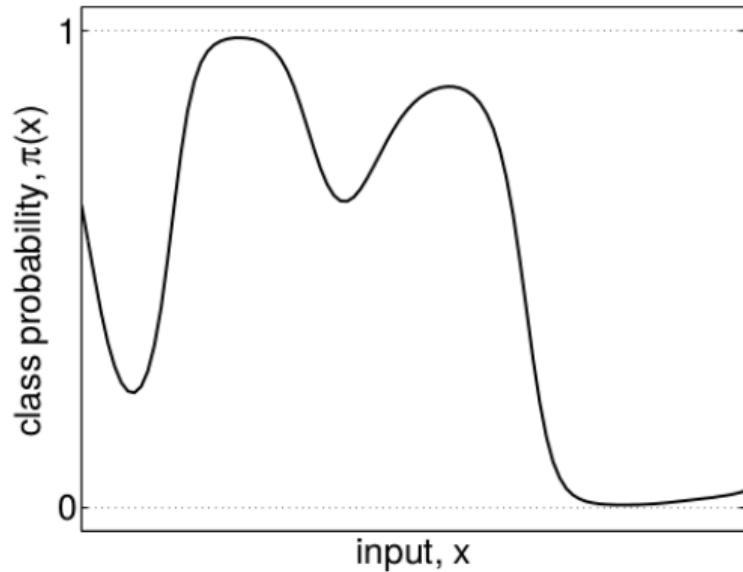
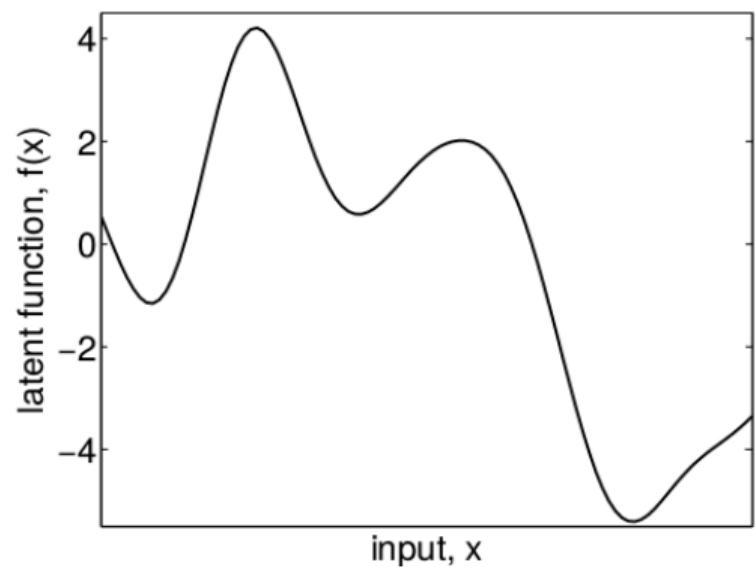
$$\underbrace{p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y})}_{\text{observed space}} = \int p(y_*|f_*) p(f_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) df_*.$$

- For a binary classification task with Bernoulli likelihood:

$$p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \overbrace{\sigma(f_*)^{y_*} (1 - \sigma(f_*)^{1-y_*})}^{p(y_*|f_*)} \overbrace{p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X})}^{\text{cond. Gaussian}} \overbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}^{\text{posterior}} d\mathbf{f} df_*,$$

$$p(y_* = 1|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \sigma(f_*) p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X}) p(\mathbf{f}|\mathbf{X}, \mathbf{y}) d\mathbf{f} df_*.$$

Beyond Gaussian likelihood



(Rasmussen and Williams, 2006)

Beyond Gaussian likelihood

- **Problem:** The posterior $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ is intractable for non-Gaussian likelihood:

$$\underbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}_{\text{posterior}} = \frac{\overbrace{p(\mathbf{y}|\mathbf{f}, \mathbf{X})}^{\text{non-Gaussian}} \overbrace{p(\mathbf{f}|\mathbf{X})}^{\text{Gaussian}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{non-Gaussian}}}.$$

- **Problem:** The marginal likelihood is also intractable for a non-Gaussian likelihood:

$$\begin{aligned}\widehat{p(\mathbf{y}|\mathbf{X})} &= \widehat{\int p(\mathbf{y}|\mathbf{f}, \mathbf{X})p(\mathbf{f}|\mathbf{X})d\mathbf{f}} \\ &= \int \underbrace{p(\mathbf{y}|\mathbf{f}, \mathbf{X})}_{\text{non-Gaussian}} \underbrace{\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f)}_{\text{GP prior}} d\mathbf{f}.\end{aligned}$$

- We need approximations for both the posterior and the marginal likelihood.

Beyond Gaussian likelihood

- **Problem:** The posterior $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ is intractable for non-Gaussian likelihood:

$$\underbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}_{\text{posterior}} = \frac{\overbrace{p(\mathbf{y}|\mathbf{f}, \mathbf{X})}^{\text{non-Gaussian}} \overbrace{p(\mathbf{f}|\mathbf{X})}^{\text{Gaussian}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{non-Gaussian}}}.$$

- **Problem:** The marginal likelihood is also intractable for a non-Gaussian likelihood:

$$\begin{aligned}\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marg. likel.}} &= \int p(\mathbf{y}|\mathbf{f}, \mathbf{X})p(\mathbf{f}|\mathbf{X})d\mathbf{f} \\ &= \int \underbrace{p(\mathbf{y}|\mathbf{f}, \mathbf{X})}_{\text{non-Gaussian}} \underbrace{\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f)}_{\text{GP prior}} d\mathbf{f}.\end{aligned}$$

- We need approximations for both the posterior and the marginal likelihood.

Beyond Gaussian likelihood

- **Problem:** The posterior $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ is intractable for non-Gaussian likelihood:

$$\underbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}_{\text{posterior}} = \frac{\overbrace{p(\mathbf{y}|\mathbf{f}, \mathbf{X})}^{\text{non-Gaussian}} \overbrace{p(\mathbf{f}|\mathbf{X})}^{\text{Gaussian}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{non-Gaussian}}}.$$

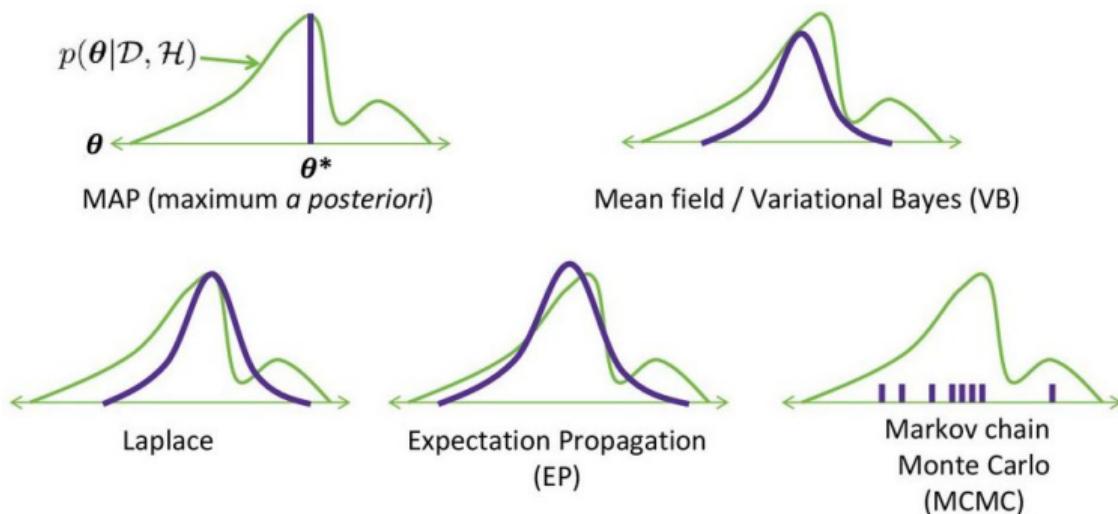
- **Problem:** The marginal likelihood is also intractable for a non-Gaussian likelihood:

$$\begin{aligned}\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marg. likel.}} &= \int p(\mathbf{y}|\mathbf{f}, \mathbf{X})p(\mathbf{f}|\mathbf{X})d\mathbf{f} \\ &= \int \underbrace{p(\mathbf{y}|\mathbf{f}, \mathbf{X})}_{\text{non-Gaussian}} \underbrace{\mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f)}_{\text{GP prior}} d\mathbf{f}.\end{aligned}$$

- We need **approximations** for both the posterior and the marginal likelihood.

Beyond Gaussian likelihood

- The approximations are usually performed by
 - Laplace approximation;
 - Variational inference;
 - Expectation Propagation;
 - Sampling methods (e.g. MCMC).



(Ben-Elazar, *et al.*, 2017)

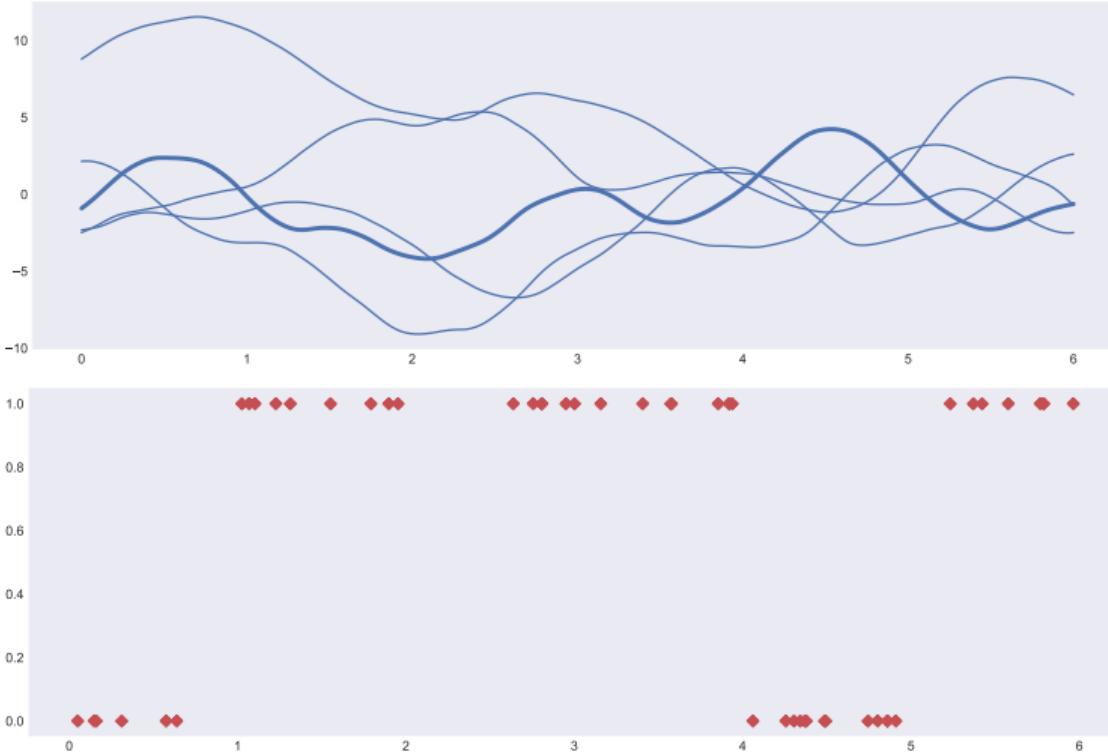
Beyond Gaussian likelihood

- For a **Gaussian approximate posterior** $q(\mathbf{f}|\mathbf{X}, \mathbf{y})$, the prediction becomes:

$$\begin{aligned} p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int \underbrace{p(y_*|f_*)}_{\text{non-Gaussian}} p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X}) \underbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}_{\text{intractable}} d\mathbf{f} df_* \\ &\approx \int p(y_*|f_*) \left[\int \underbrace{p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X})}_{\text{cond. Gaussian}} \underbrace{q(\mathbf{f}|\mathbf{X}, \mathbf{y})}_{\text{Gaussian approx.}} d\mathbf{f} \right] df_*. \end{aligned}$$

- The integral on \mathbf{f} is tractable, since it involves only Gaussians.
- The integral on f_* is usually intractable, but since it is one-dimensional, it can be cheaply approximated by Monte Carlo, Gauss-Hermite quadrature, etc.
- Importantly, the previous frameworks also approximate the marginal likelihood.

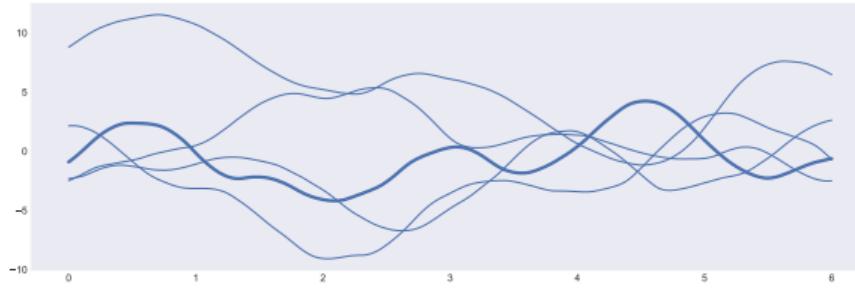
Beyond Gaussian likelihood



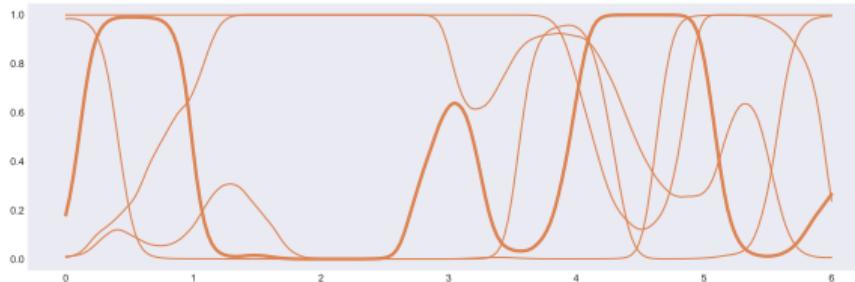
$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f)$$

$$(\mathbf{X}, \mathbf{y}), y_i \in \{0, 1\}, \forall i$$

Beyond Gaussian likelihood

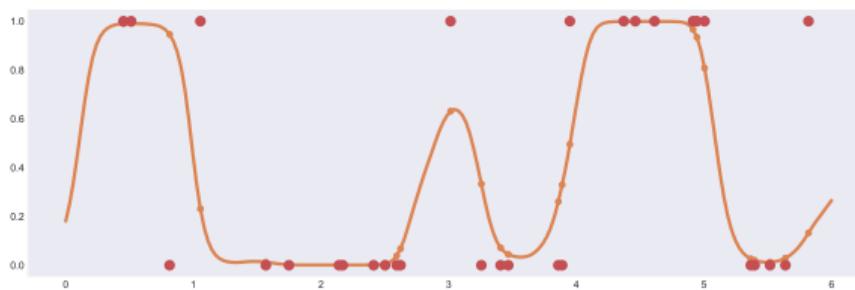


$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f)$$



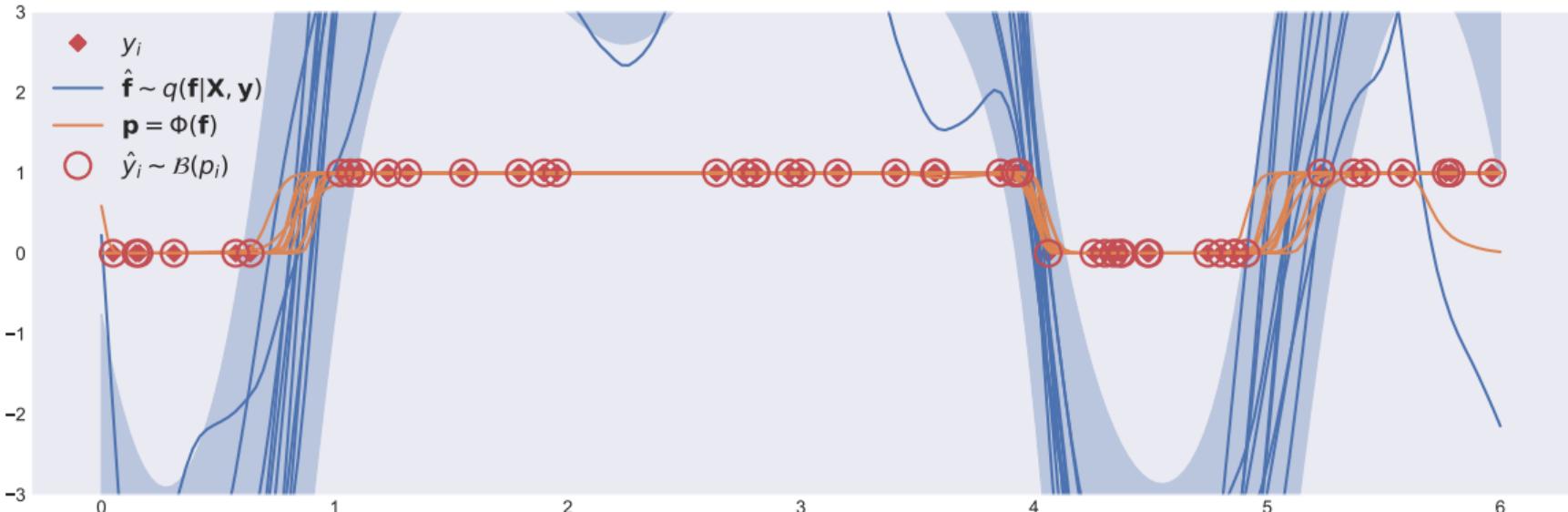
$$\mathbf{p} = \Phi(\mathbf{f}) \triangleq \int_{-\infty}^f \mathcal{N}(f|0, 1) df$$

(Normal CDF)



$$y_i \sim \mathcal{B}(p_i), \forall i$$

Beyond Gaussian likelihood



- Gaussian approximate posterior $q(\mathbf{f}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{f}|\mathbf{m}, \mathbf{\Sigma})$.
- Variational inference with **GPflow** (whose documentation this illustration was based on).
- Latent samples are continuous, but samples in the observed space remain discrete.

Warped GPs

- A **warping function** may be a useful preprocessing step for non-Gaussian observations.
 - For instance, if $y > 0$, $f = \log(y)$.
- We can also consider a parameterised warping function $w(\cdot|\phi)$:

$$f = w(y|\phi).$$

- $w(\cdot|\phi)$ must be monotonic, so that $y = w^{-1}(f|\phi)$ is unambiguous.
- If $p(f|X)$ is a GP, we obtain a **warped GP model** (Snelson *et al.*, 2004).
- A **Bayesian warped GP** is obtained by directly modeling the inverse warping function (Lázaro-Gredilla, 2012):

$$y = g(f(x)) + \epsilon,$$

where both $f(\cdot)$ and $g(\cdot)$ have GP priors and ϵ is an observation noise.

→ Analogous to a 2-layer Deep GP!

Warped GPs

- A **warping function** may be a useful preprocessing step for non-Gaussian observations.
 - For instance, if $y > 0$, $f = \log(y)$.
- We can also consider a parameterised warping function $w(\cdot|\phi)$:

$$f = w(y|\phi).$$

- $w(\cdot|\phi)$ must be monotonic, so that $y = w^{-1}(f|\phi)$ is unambiguous.
- If $p(f|X)$ is a GP, we obtain a **warped GP model** (Snelson *et al.*, 2004).
- A **Bayesian warped GP** is obtained by directly modeling the inverse warping function (Lázaro-Gredilla, 2012):

$$y = g(f(x)) + \epsilon,$$

where both $f(\cdot)$ and $g(\cdot)$ have GP priors and ϵ is an observation noise.

→ Analogous to a 2-layer Deep GP!

Warped GPs

- A **warping function** may be a useful preprocessing step for non-Gaussian observations.
 - For instance, if $y > 0$, $f = \log(y)$.
- We can also consider a parameterised warping function $w(\cdot|\phi)$:

$$f = w(y|\phi).$$

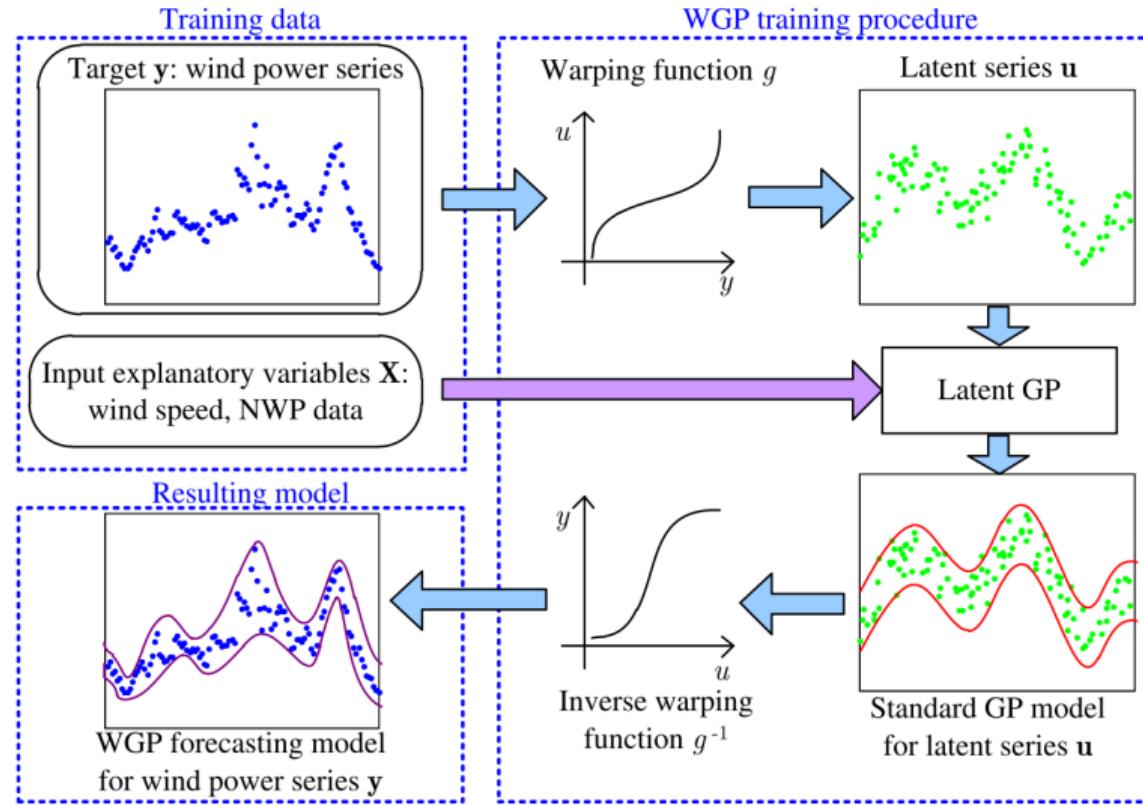
- $w(\cdot|\phi)$ must be monotonic, so that $y = w^{-1}(f|\phi)$ is unambiguous.
- If $p(f|X)$ is a GP, we obtain a **warped GP model** (Snelson *et al.*, 2004).
- A **Bayesian warped GP** is obtained by directly modeling the inverse warping function (Lázaro-Gredilla, 2012):

$$y = g(f(x)) + \epsilon,$$

where both $f(\cdot)$ and $g(\cdot)$ have GP priors and ϵ is an observation noise.

- Analogous to a 2-layer Deep GP!

Warped GPs



Kou, Peng et al. "Sparse online warped Gaussian process for wind power probabilistic forecasting". Applied Energy, 2013.

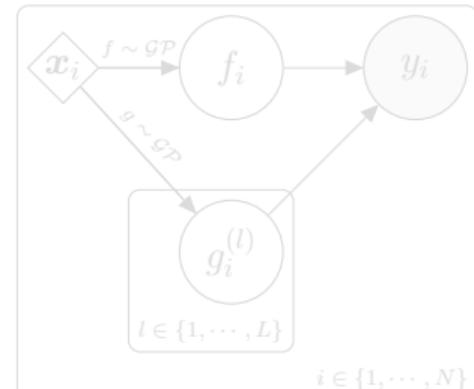
Beyond Gaussian likelihood

- Flexible models can be obtained by considering **multiple GP priors**.
- **Heteroscedastic Gaussian likelihood**, i.e. input-dependent noise variance (Goldberg *et al.*, 1998; Lázaro-Gredilla and Titsias, 2011):

$$y = f(\mathbf{x}) + \epsilon,$$
$$f|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f), \quad \epsilon|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \exp(g(\mathbf{x}))), \quad g|\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_g, \mathbf{K}_g).$$

- Chained GPs with L GP priors for the likelihood parameters ϕ (Saul *et al*, 2016):

$$\phi = [\mathbf{g}^{(1)}, \dots, \mathbf{g}^{(L)}],$$
$$p(\mathbf{y}|\mathbf{f}, \phi) = p(\mathbf{y}|\mathbf{f}, \mathbf{g}^{(1)}, \dots, \mathbf{g}^{(L)}),$$
$$f|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f),$$
$$\mathbf{g}^{(l)}|\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_{g^{(l)}}, \mathbf{K}_{g^{(l)}}), \quad 1 \leq l \leq L.$$



Beyond Gaussian likelihood

- Flexible models can be obtained by considering **multiple GP priors**.
- **Heteroscedastic Gaussian likelihood**, i.e. input-dependent noise variance (Goldberg *et al.*, 1998; Lázaro-Gredilla and Titsias, 2011):

$$y = f(\mathbf{x}) + \epsilon,$$

$$f|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f), \quad \epsilon|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \exp(g(\mathbf{x}))), \quad g|\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_g, \mathbf{K}_g).$$

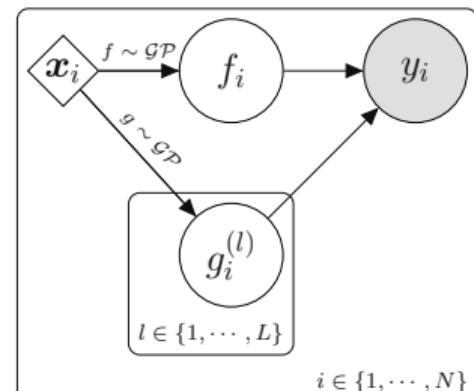
- **Chained GPs** with L GP priors for the likelihood parameters ϕ (Saul *et al*, 2016):

$$\phi = [\mathbf{g}^{(1)}, \dots, \mathbf{g}^{(L)}],$$

$$p(\mathbf{y}|\mathbf{f}, \phi) = p(\mathbf{y}|\mathbf{f}, \mathbf{g}^{(1)}, \dots, \mathbf{g}^{(L)}),$$

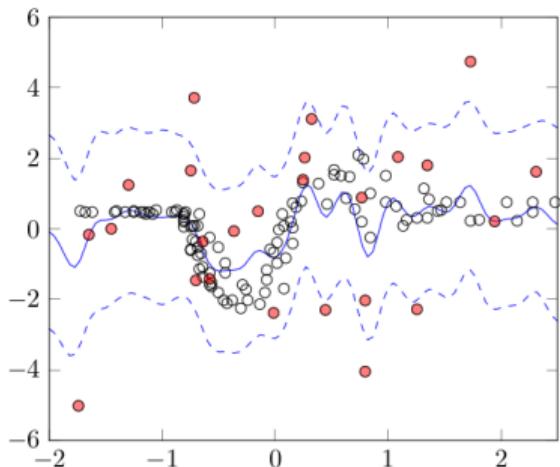
$$\mathbf{f}|\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f),$$

$$\mathbf{g}^{(l)}|\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_{g^{(l)}}, \mathbf{K}_{g^{(l)}}), \quad 1 \leq l \leq L.$$

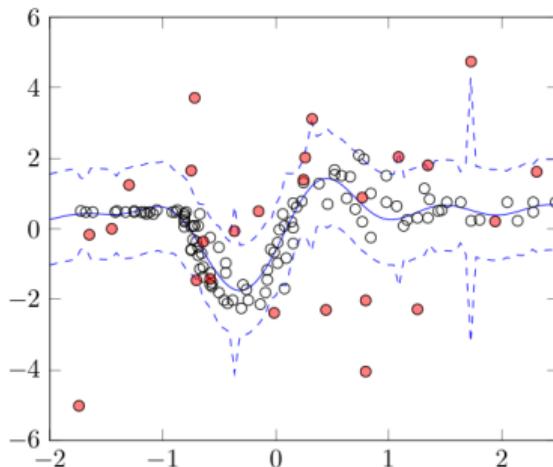


Chained GPs

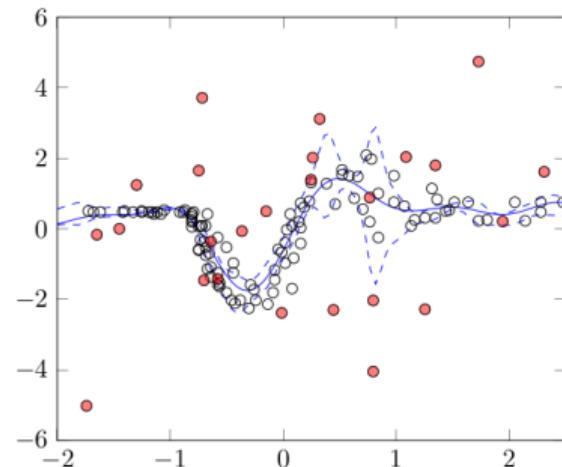
Standard Gaussian Process



Heteroscedastic Gaussian



Heteroscedastic Student-t



Saul, Alan, et al. “**Chained Gaussian processes**”. AISTATS, 2016.

$$y_i \sim \text{St}(\mu = f(\mathbf{x}_i), \sigma^2 = \exp(g(\mathbf{x}_i)), \nu),$$

$$\mathbf{f}|\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_f, \mathbf{K}_f),$$

$$\mathbf{g}|\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_g, \mathbf{K}_g).$$

Beyond Gaussian likelihood

Summary

- GP models are tractable only for Gaussian likelihood.
- Other noise models are equivalent to consider **non-Gaussian likelihoods**.
- We approximate the posterior $\underbrace{p(\mathbf{f}|\mathbf{X}, \mathbf{y})}_{\text{predictions}}$ and the marginal likelihood $\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{model selection}}$.
 - Usual approaches: LA, VI, EP or MCMC.
- The observations can also be **warped**, either by a parametric function or by directly considering a **second GP prior** for the **inverse warping function**.
- By considering **multiple latent functions**, we can obtain flexible **heteroscedastic** likelihoods.

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

Sparse approximations

- We have seen how GP models are flexible and can be applied to both Gaussian and non-Gaussian observations.
- **Problem:** GP expressions involve the terms $|\mathbf{K}_f + \sigma_y^2 \mathbf{I}|$ and $(\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1}$:

$$\underbrace{\log p(\mathbf{y} | \mathbf{X})}_{\text{marg. likelihood}} = -\frac{1}{2} \log |\mathbf{K}_f + \sigma_y^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N}{2} \log(2\pi),$$
$$\underbrace{p(y_* | \mathbf{x}_*, \mathbf{y}, \mathbf{X})}_{\text{predictive dist.}} = \mathcal{N}(y_* | \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*} + \sigma_y^2).$$

- Naive computation scales with $\mathcal{O}(N^3)$.
- Storage scales with $\mathcal{O}(N^2)$.
- **Goal:** Mitigate (or eliminate!) the influence of N in the scaling.
- **Idea:** Maintain a smaller set of points that summarizes the available data.

Sparse approximations

- We have seen how GP models are flexible and can be applied to both Gaussian and non-Gaussian observations.
- **Problem:** GP expressions involve the terms $|\mathbf{K}_f + \sigma_y^2 \mathbf{I}|$ and $(\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1}$:

marg. likelihood

$$\overbrace{\log p(\mathbf{y} | \mathbf{X})}^{\text{marg. likelihood}} = -\frac{1}{2} \log |\mathbf{K}_f + \sigma_y^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N}{2} \log(2\pi),$$

$$\underbrace{p(y_* | \mathbf{x}_*, \mathbf{y}, \mathbf{X})}_{\text{predictive dist.}} = \mathcal{N} \left(y_* | \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*} + \sigma_y^2 \right).$$

→ Naive computation scales with $\mathcal{O}(N^3)$.

→ Storage scales with $\mathcal{O}(N^2)$.

- **Goal:** Mitigate (or eliminate!) the influence of N in the scaling.
- **Idea:** Maintain a smaller set of points that summarizes the available data.

Sparse approximations

- We have seen how GP models are flexible and can be applied to both Gaussian and non-Gaussian observations.
- **Problem:** GP expressions involve the terms $|\mathbf{K}_f + \sigma_y^2 \mathbf{I}|$ and $(\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1}$:

marg. likelihood

$$\overbrace{\log p(\mathbf{y} | \mathbf{X})}^{\text{marg. likelihood}} = -\frac{1}{2} \log |\mathbf{K}_f + \sigma_y^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N}{2} \log(2\pi),$$

$$\underbrace{p(y_* | \mathbf{x}_*, \mathbf{y}, \mathbf{X})}_{\text{predictive dist.}} = \mathcal{N} \left(y_* | \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*} + \sigma_y^2 \right).$$

→ Naive computation scales with $\mathcal{O}(N^3)$.

→ Storage scales with $\mathcal{O}(N^2)$.

- **Goal:** Mitigate (or eliminate!) the influence of N in the scaling.
- **Idea:** Maintain a smaller set of points that summarizes the available data.

Sparse approximations

- We have seen how GP models are flexible and can be applied to both Gaussian and non-Gaussian observations.
- **Problem:** GP expressions involve the terms $|\mathbf{K}_f + \sigma_y^2 \mathbf{I}|$ and $(\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1}$:

marg. likelihood

$$\overbrace{\log p(\mathbf{y} | \mathbf{X})}^{\text{marg. likelihood}} = -\frac{1}{2} \log |\mathbf{K}_f + \sigma_y^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N}{2} \log(2\pi),$$

$$\underbrace{p(y_* | \mathbf{x}_*, \mathbf{y}, \mathbf{X})}_{\text{predictive dist.}} = \mathcal{N} \left(y_* | \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{y}, k_{**} - \mathbf{k}_{f*}^\top (\mathbf{K}_f + \sigma_y^2 \mathbf{I})^{-1} \mathbf{k}_{f*} + \sigma_y^2 \right).$$

→ Naive computation scales with $\mathcal{O}(N^3)$.

→ Storage scales with $\mathcal{O}(N^2)$.

- **Goal:** Mitigate (or eliminate!) the influence of N in the scaling.
- **Idea:** Maintain a smaller set of points that **summarizes the available data**.

Sparse approximations

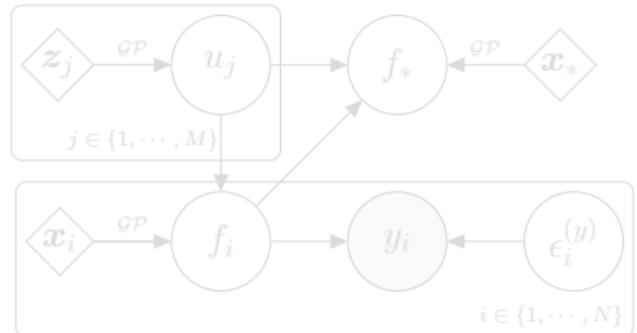
- Consider M inducing variables $\mathbf{u} \in \mathbb{R}^M$ from the same GP prior of \mathbf{f} .
 → \mathbf{u} corresponds to the evaluation of $f(\cdot)$ at M **pseudo-inputs** $\mathbf{z}_j|_{j=1}^M \in \mathbb{R}^D$.

marg. likelihood

$$\overbrace{p(y|X)}^{} = \mathcal{N}(y|\mathbf{0}, K_f + \sigma_y^2 I),$$

$$\underbrace{p(f|X)}_{\text{prior}} = \mathcal{N}(f|\mathbf{0}, K_f), \quad [K_f]_{ii'} = k(x_i, x_{i'}),$$

$$\underbrace{p(u|Z)}_{\text{prior}} = \mathcal{N}(u|\mathbf{0}, K_u), \quad [K_u]_{jj'} = k(z_j, z_{j'}).$$



- If the knowledge from \mathbf{f} is concentrated in \mathbf{u} , the approximate posterior becomes

$$q(f_*) = \int \underbrace{p(f_*|u)}_{\approx p(f_*|u, f)} \underbrace{q(u)}_{\approx p(u|y)} du.$$

→ If $M = N$ and $z_i = x_i, \forall i$, we get $\mathbf{u} = \mathbf{f}$ and recover the original posterior exactly.

Sparse approximations

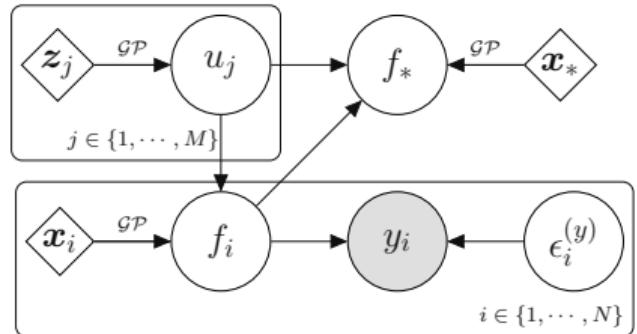
- Consider M inducing variables $\mathbf{u} \in \mathbb{R}^M$ from the same GP prior of \mathbf{f} .
 → \mathbf{u} corresponds to the evaluation of $f(\cdot)$ at M **pseudo-inputs** $\mathbf{z}_j|_{j=1}^M \in \mathbb{R}^D$.

marg. likelihood

$$\overbrace{p(\mathbf{y}|\mathbf{X})} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_f + \sigma_y^2 \mathbf{I}),$$

$$\underbrace{p(\mathbf{f}|\mathbf{X})}_{\text{prior}} = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f), \quad [\mathbf{K}_f]_{ii'} = k(\mathbf{x}_i, \mathbf{x}_{i'}),$$

$$\underbrace{p(\mathbf{u}|\mathbf{Z})}_{\text{prior}} = \mathcal{N}(\mathbf{u}|\mathbf{0}, \mathbf{K}_u), \quad [\mathbf{K}_u]_{jj'} = k(\mathbf{z}_j, \mathbf{z}_{j'}).$$



- If the knowledge from \mathbf{f} is concentrated in \mathbf{u} , the approximate posterior becomes

$$q(f_*) = \int \underbrace{p(f_*|\mathbf{u})}_{\approx p(f_*|\mathbf{u}, \mathbf{f}) \approx p(\mathbf{u}|\mathbf{y})} \underbrace{q(\mathbf{u})}_{\approx p(\mathbf{u}|\mathbf{y})} d\mathbf{u}.$$

→ If $M = N$ and $\mathbf{z}_i = \mathbf{x}_i, \forall i$, we get $\mathbf{u} = \mathbf{f}$ and recover the original posterior exactly.

Sparse approximations

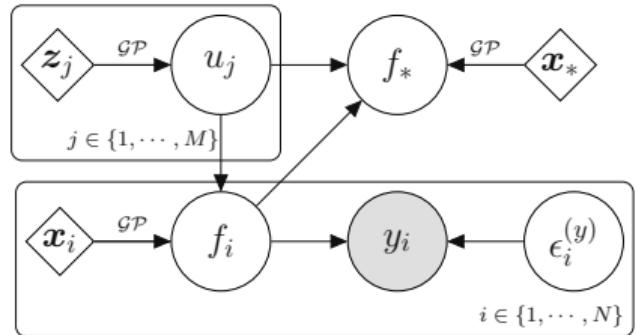
- Consider M inducing variables $\mathbf{u} \in \mathbb{R}^M$ from the same GP prior of \mathbf{f} .
 → \mathbf{u} corresponds to the evaluation of $f(\cdot)$ at M **pseudo-inputs** $\mathbf{z}_j|_{j=1}^M \in \mathbb{R}^D$.

marg. likelihood

$$\overbrace{p(\mathbf{y}|\mathbf{X})} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_f + \sigma_y^2 \mathbf{I}),$$

$$\underbrace{p(\mathbf{f}|\mathbf{X})}_{\text{prior}} = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f), \quad [\mathbf{K}_f]_{ii'} = k(\mathbf{x}_i, \mathbf{x}_{i'}),$$

$$\underbrace{p(\mathbf{u}|\mathbf{Z})}_{\text{prior}} = \mathcal{N}(\mathbf{u}|\mathbf{0}, \mathbf{K}_u), \quad [\mathbf{K}_u]_{jj'} = k(\mathbf{z}_j, \mathbf{z}_{j'}).$$



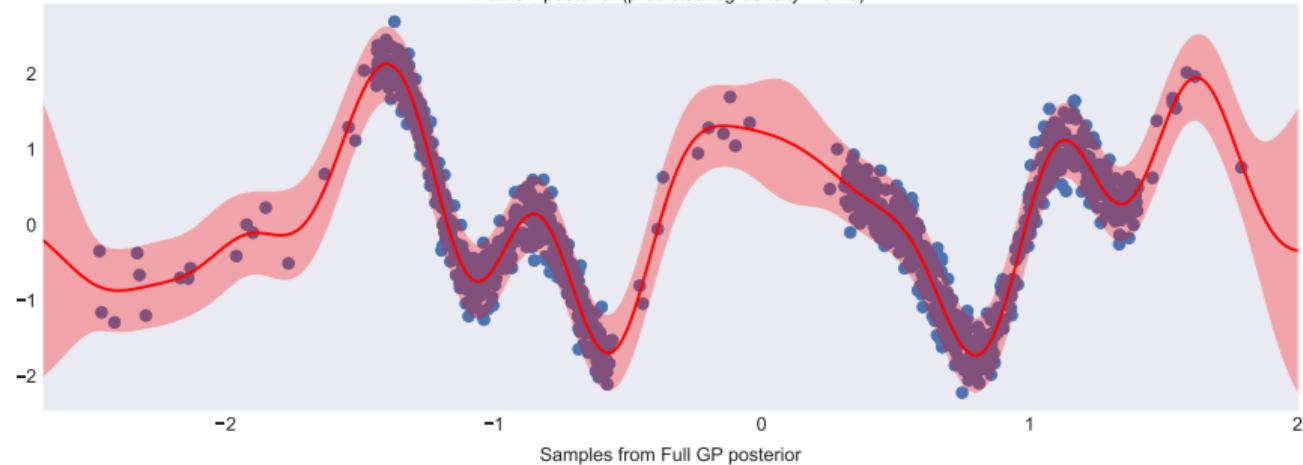
- If the knowledge from \mathbf{f} is concentrated in \mathbf{u} , the **approximate posterior** becomes

$$q(f_*) = \int \underbrace{p(f_*|\mathbf{u})}_{\approx p(f_*|\mathbf{u}, \mathbf{f}) \approx p(\mathbf{u}|\mathbf{y})} \underbrace{q(\mathbf{u})}_{\approx p(\mathbf{u}|\mathbf{y})} \mathrm{d}\mathbf{u}.$$

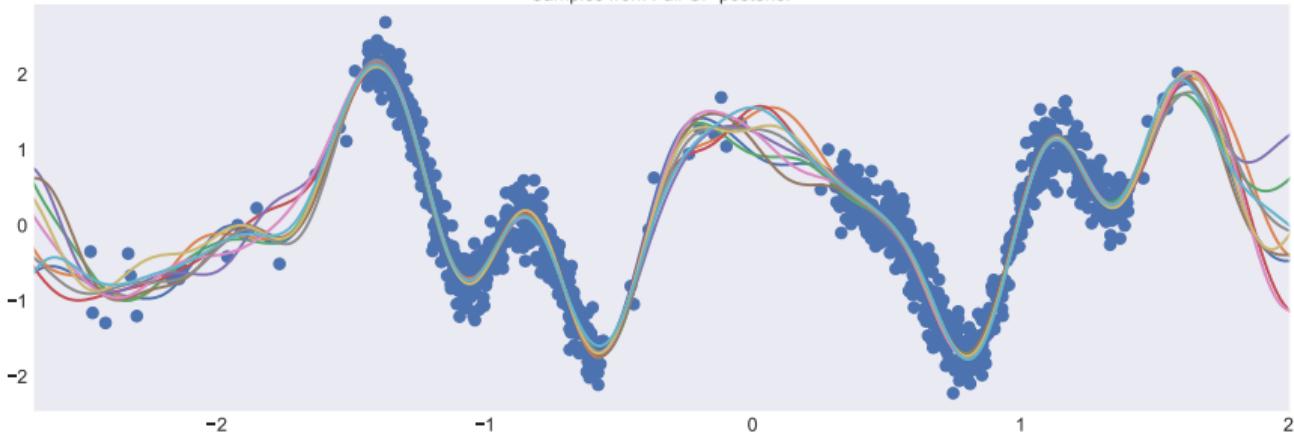
→ If $M = N$ and $\mathbf{z}_i = \mathbf{x}_i, \forall i$, we get $\mathbf{u} = \mathbf{f}$ and recover the original posterior exactly.

Full GP posterior (1083 observations)

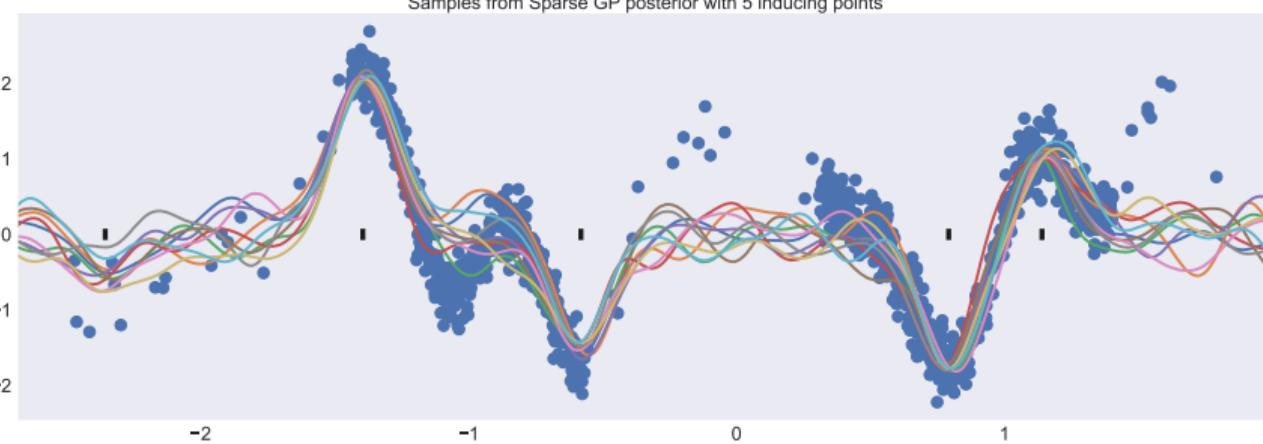
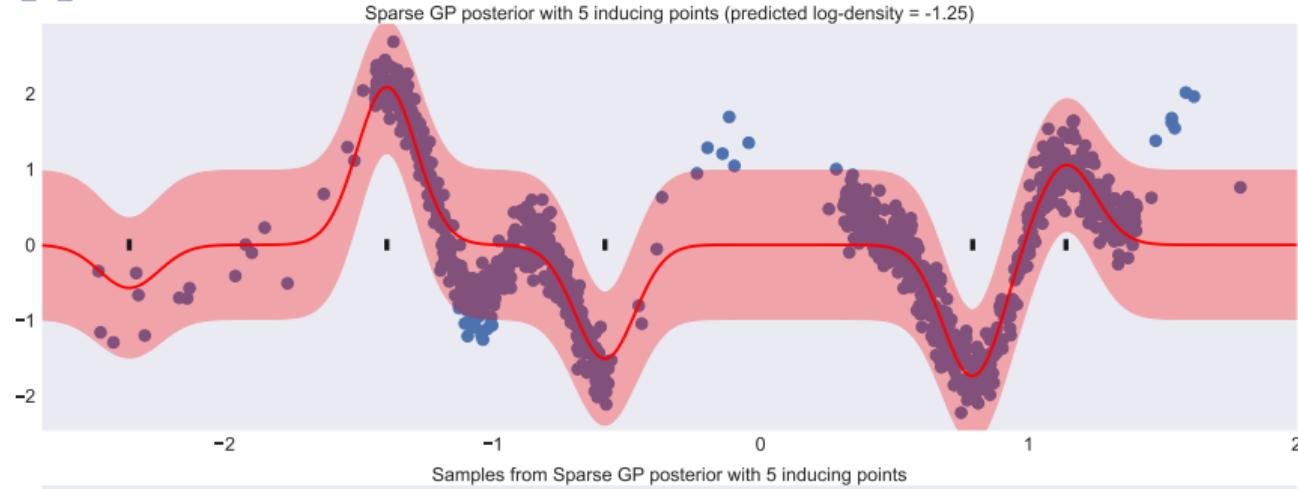
Full GP posterior (predicted log-density = 0.16)



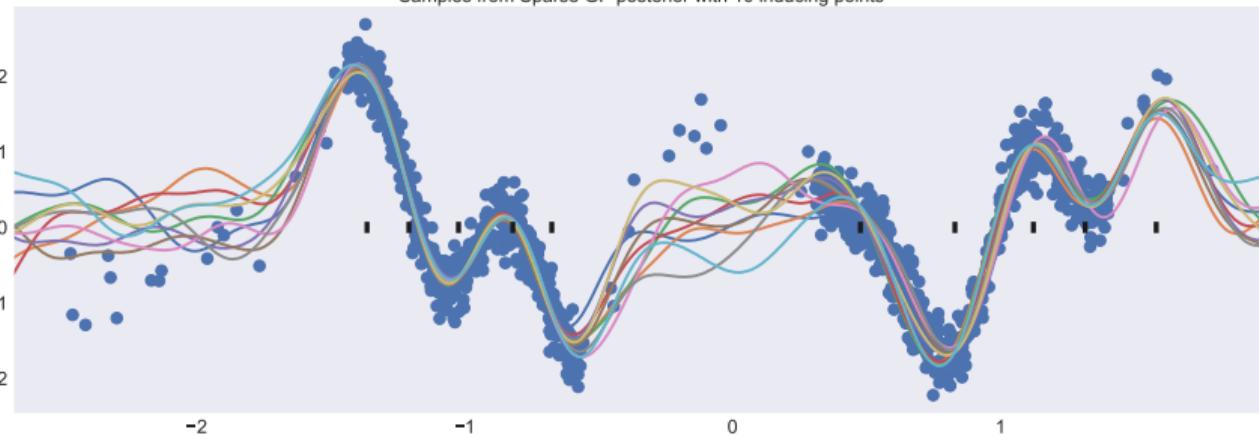
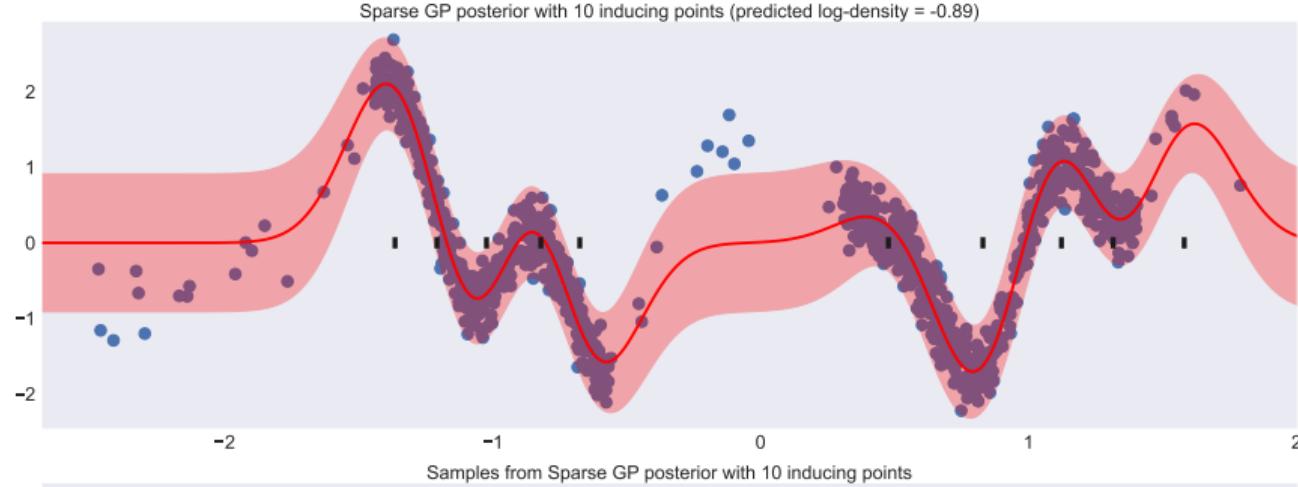
Samples from Full GP posterior



Sparse approximations (1083 observations, 5 inducing points)

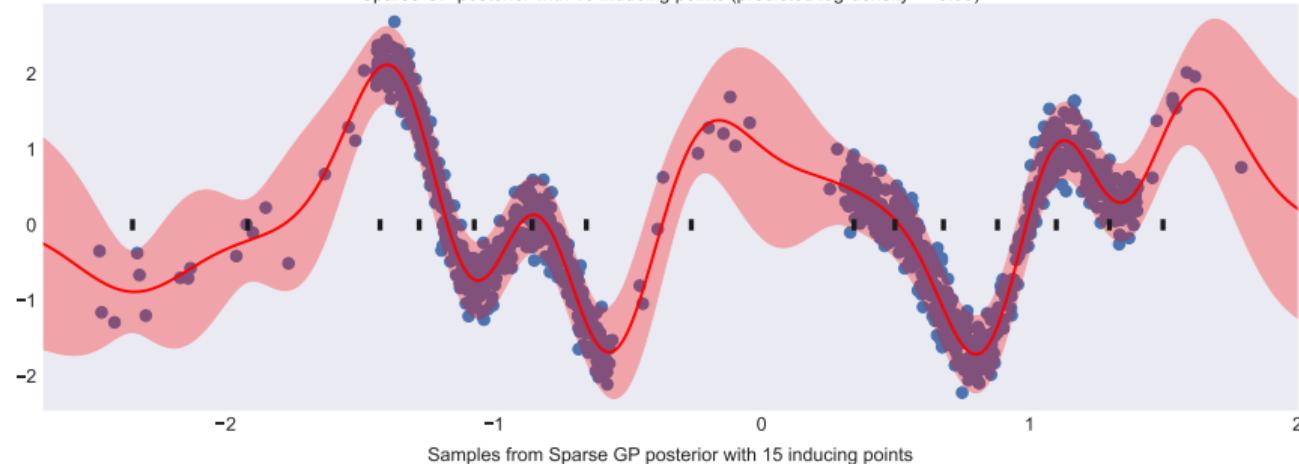


Sparse approximations (1083 observations, 10 inducing points)

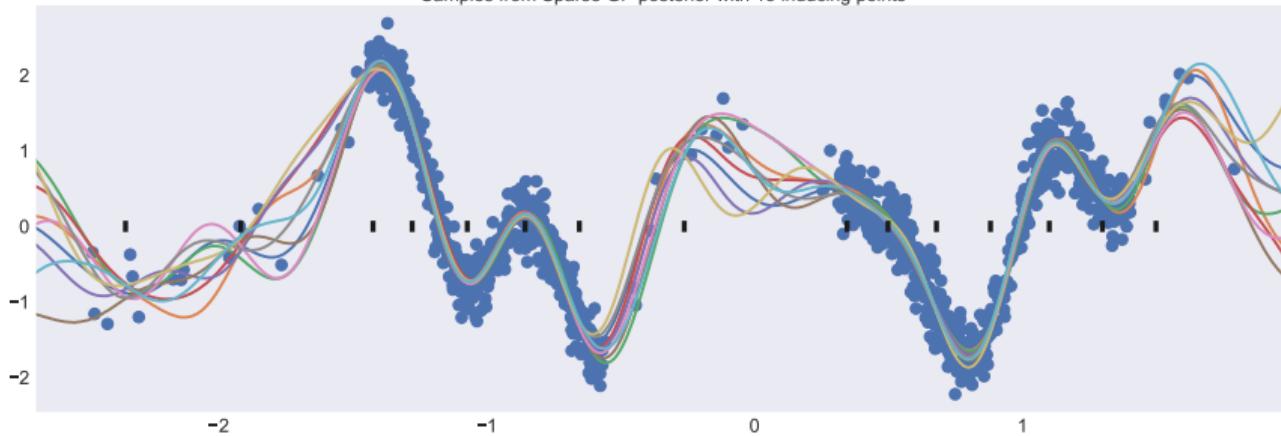


Sparse approximations (1083 observations, 15 inducing points)

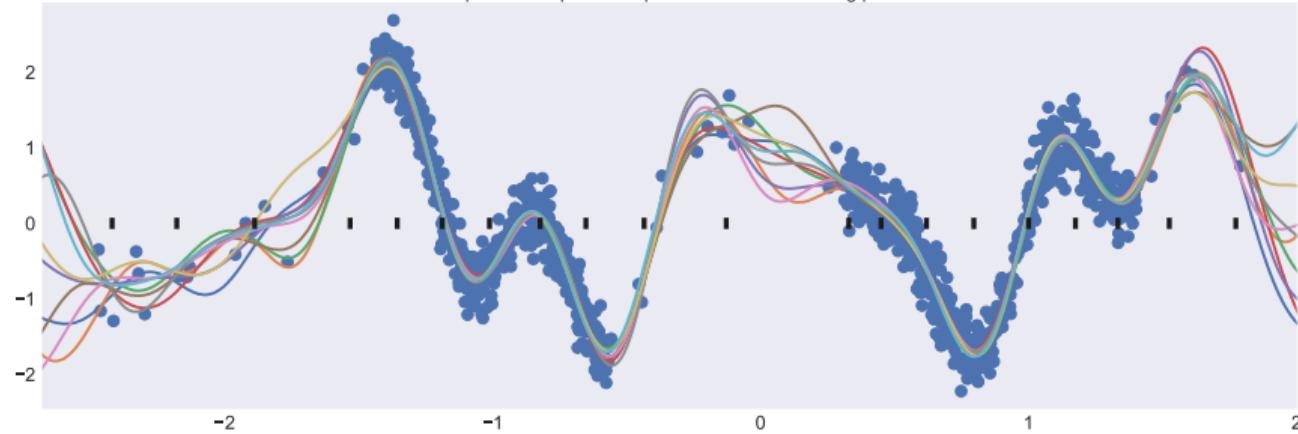
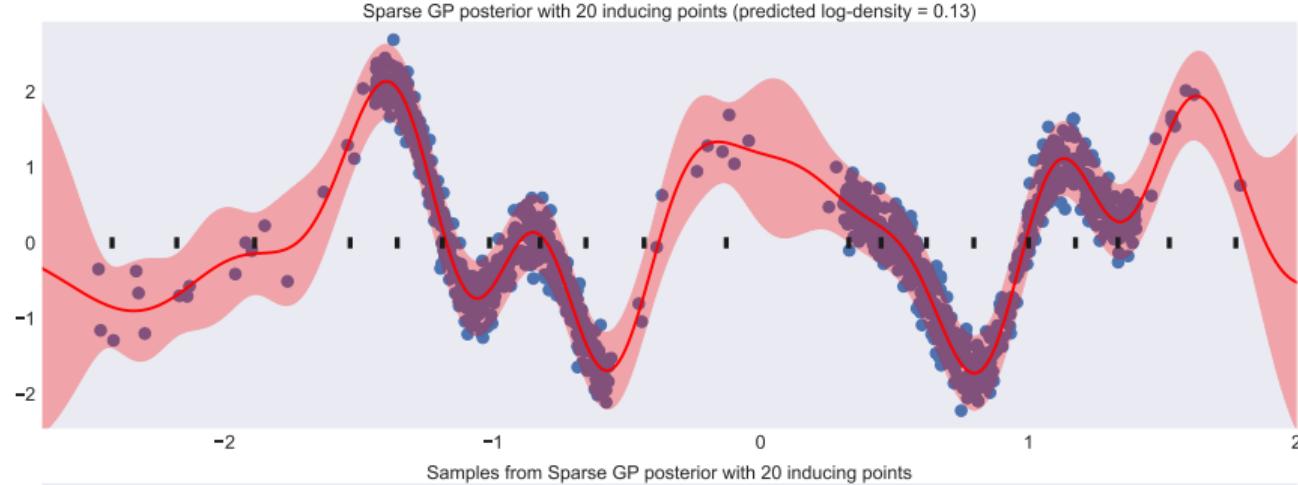
Sparse GP posterior with 15 inducing points (predicted log-density = -0.08)



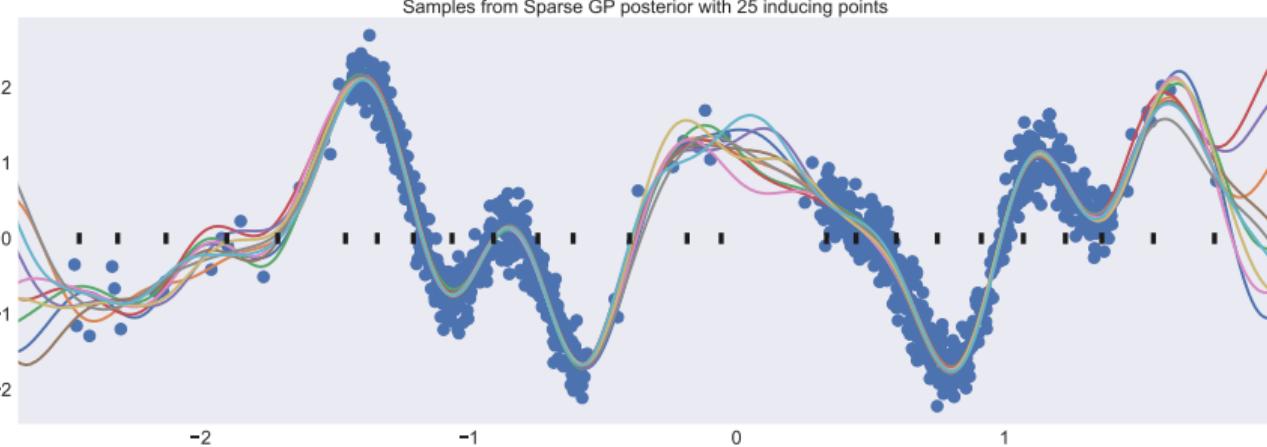
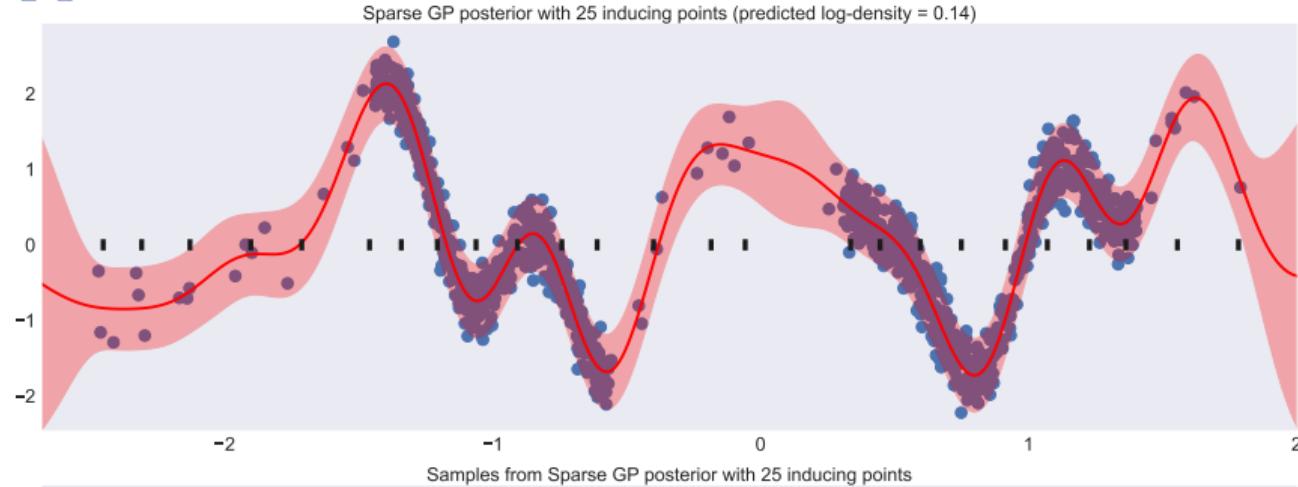
Samples from Sparse GP posterior with 15 inducing points



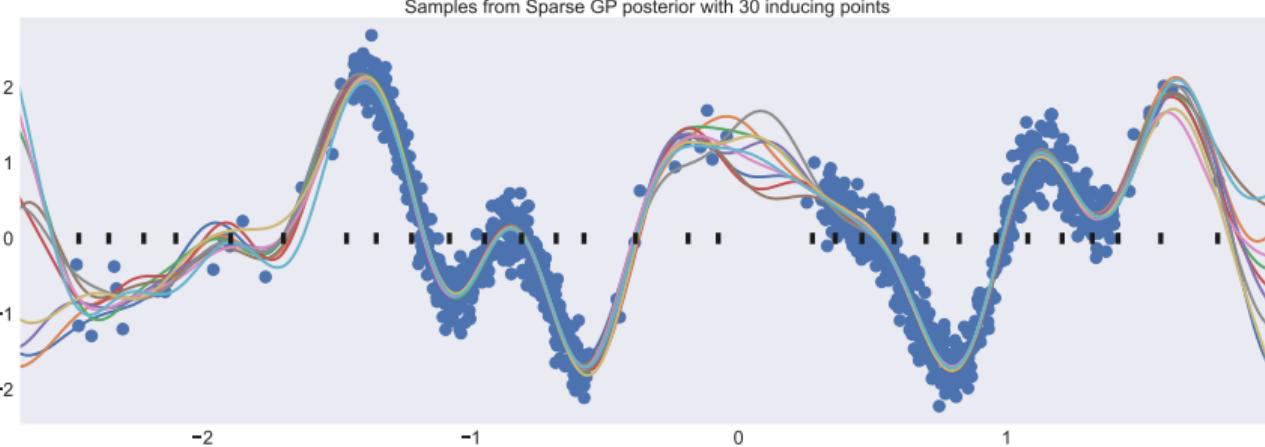
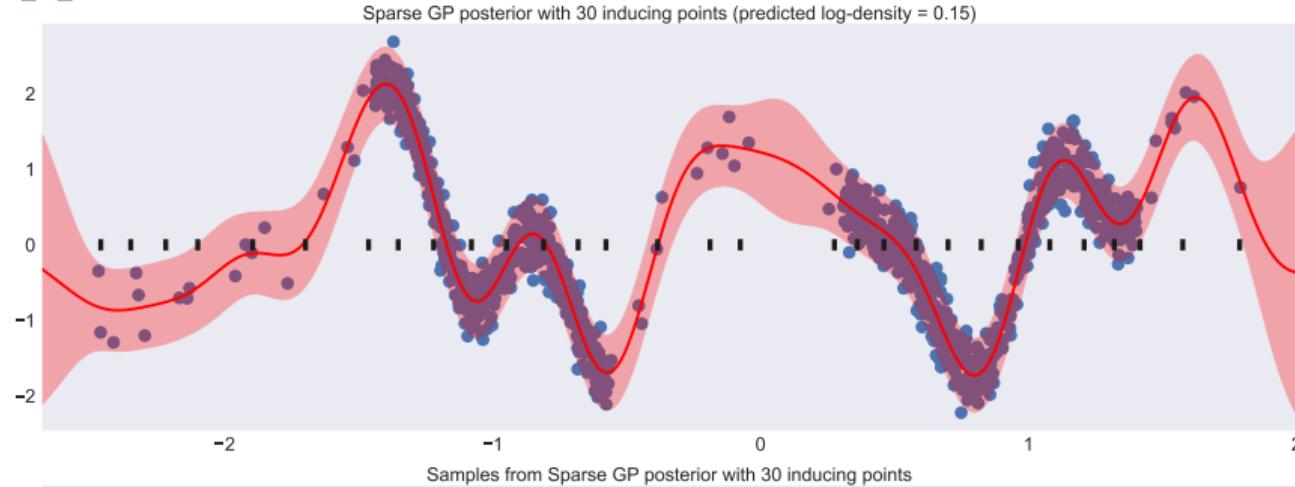
Sparse approximations (1083 observations, 20 inducing points)



Sparse approximations (1083 observations, 25 inducing points)

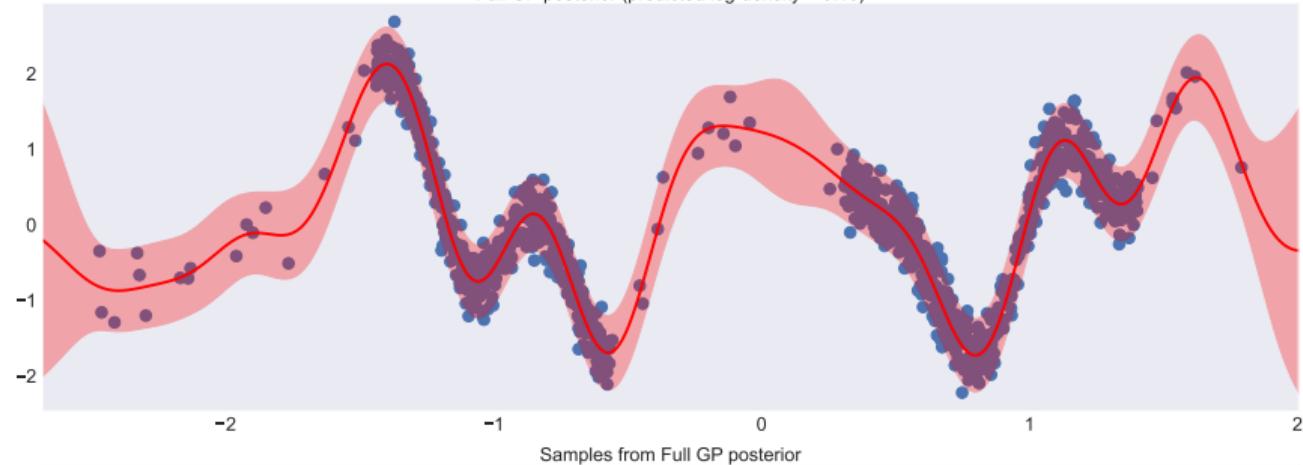


Sparse approximations (1083 observations, 30 inducing points)

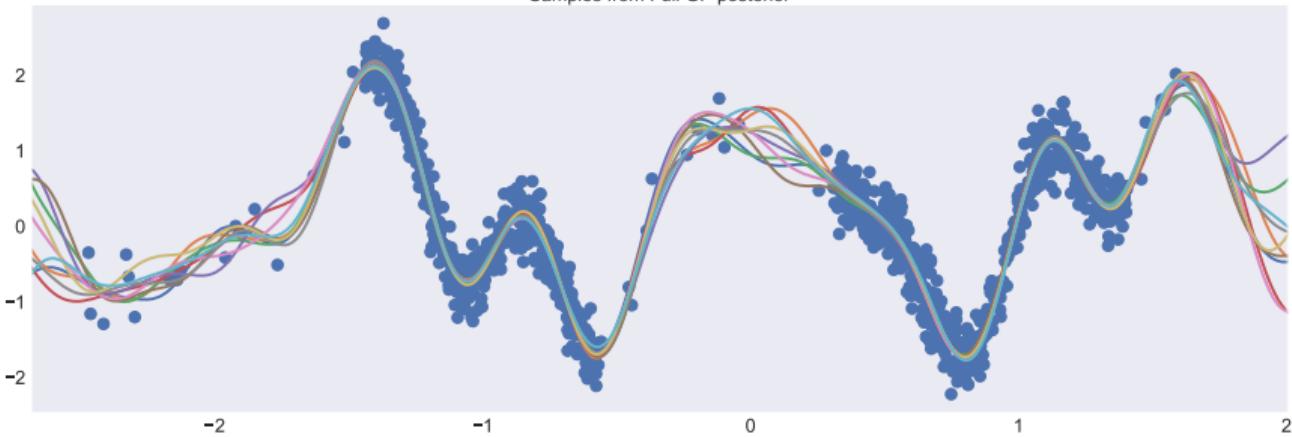


Full GP posterior (1083 observations)

Full GP posterior (predicted log-density = 0.16)



Samples from Full GP posterior



Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the Kullback-Leibler divergence:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the **Kullback-Leibler divergence**:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the **Kullback-Leibler divergence**:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the **Kullback-Leibler divergence**:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the **Kullback-Leibler divergence**:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the **Kullback-Leibler divergence**:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- Given some data, we want an approximation $q(\mathbf{f})$ for the true posterior $p(\mathbf{f}|\mathbf{y})$.
- The approximation quality can be quantified via the **Kullback-Leibler divergence**:

$$\begin{aligned}\text{KL}(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right] \geq 0 \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \frac{p(\mathbf{u}|\mathbf{f})}{p(\mathbf{u}|\mathbf{f})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}|\mathbf{y})} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}} \right] \\ &= \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})} \right] + \overbrace{\log p(\mathbf{y})}^{\text{evidence}} \geq 0, \\ \implies \log p(\mathbf{y}) &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right].\end{aligned}$$

Sparse approximations

- We replace $q(\mathbf{f}) = \int q(\mathbf{f}, \mathbf{u}) d\mathbf{u}$ and (conveniently) choose $q(\mathbf{f}, \mathbf{u}) = q(\mathbf{u}) \underbrace{p(\mathbf{f}|\mathbf{u})}_{\text{cond. prior}}$:

$$\begin{aligned}\underbrace{\log p(\mathbf{y})}_{\text{evidence}} &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{u})}{q(\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{u})} \left[\log \frac{p(\mathbf{u})}{q(\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{y}|\mathbf{f})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})) = \underbrace{\mathcal{L}}_{\text{ELBO}}.\end{aligned}$$

- The kernel hyperparameters and the pseudo-inputs $\mathbf{z}_j|_1^M$ can be optimized by maximizing the **ELBO (evidence lower bound)** \mathcal{L} , improving the approximation.

Sparse approximations

- We replace $q(\mathbf{f}) = \int q(\mathbf{f}, \mathbf{u}) d\mathbf{u}$ and (conveniently) choose $q(\mathbf{f}, \mathbf{u}) = q(\mathbf{u}) \underbrace{p(\mathbf{f}|\mathbf{u})}_{\text{cond. prior}}$:

$$\begin{aligned}\underbrace{\log p(\mathbf{y})}_{\text{evidence}} &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{u})}{q(\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{u})} \left[\log \frac{p(\mathbf{u})}{q(\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{y}|\mathbf{f})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})) = \underbrace{\mathcal{L}}_{\text{ELBO}}.\end{aligned}$$

- The kernel hyperparameters and the pseudo-inputs $\mathbf{z}_j|_1^M$ can be optimized by maximizing the **ELBO (evidence lower bound)** \mathcal{L} , improving the approximation.

Sparse approximations

- We replace $q(\mathbf{f}) = \int q(\mathbf{f}, \mathbf{u}) d\mathbf{u}$ and (conveniently) choose $q(\mathbf{f}, \mathbf{u}) = q(\mathbf{u}) \underbrace{p(\mathbf{f}|\mathbf{u})}_{\text{cond. prior}}$:

$$\begin{aligned}\underbrace{\log p(\mathbf{y})}_{\text{evidence}} &\geq \mathbb{E}_{q(\mathbf{f})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} \left[\log \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{u})}{q(\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{u})} \left[\log \frac{p(\mathbf{u})}{q(\mathbf{u})} \right] \\ &\geq \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})} [\log p(\mathbf{y}|\mathbf{f})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})) = \underbrace{\mathcal{L}}_{\text{ELBO}}.\end{aligned}$$

- The kernel hyperparameters and the pseudo-inputs $\mathbf{z}_j|_1^M$ can be optimized by maximizing the **ELBO (evidence lower bound)** \mathcal{L} , improving the approximation.

Sparse approximations

- What we need to compute the ELBO \mathcal{L} ?

$$\log p(\mathbf{y}) \geq \mathcal{L} = \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})}[\log p(\mathbf{y}|\mathbf{f})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})).$$

- Since the joint $p(\mathbf{f}, \mathbf{u})$ is Gaussian, the conditional $p(\mathbf{f}|\mathbf{u})$ is also Gaussian and given by (with omitted dependence on \mathbf{X} and \mathbf{z}):

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_f & \mathbf{K}_{fu} \\ \mathbf{K}_{fu}^\top & \mathbf{K}_u \end{bmatrix}\right), \quad [\mathbf{K}_{fu}]_{ij} = k(\mathbf{x}_i, \mathbf{z}_j),$$

$$p(\mathbf{f}|\mathbf{u}) = \mathcal{N}(\mathbf{f} | \mathbf{K}_{fu}\mathbf{K}_u^{-1}\mathbf{u}, \mathbf{K}_f - \mathbf{K}_{fu}\mathbf{K}_u^{-1}\mathbf{K}_{fu}^\top).$$

- How to define $q(\mathbf{u})$?

Sparse approximations

- What we need to compute the ELBO \mathcal{L} ?

$$\log p(\mathbf{y}) \geq \mathcal{L} = \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})}[\log p(\mathbf{y}|\mathbf{f})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})).$$

- Since the joint $p(\mathbf{f}, \mathbf{u})$ is Gaussian, the conditional $p(\mathbf{f}|\mathbf{u})$ is also Gaussian and given by (with omitted dependence on \mathbf{X} and \mathbf{z}):

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_f & \mathbf{K}_{fu} \\ \mathbf{K}_{fu}^\top & \mathbf{K}_u \end{bmatrix}\right), \quad [\mathbf{K}_{fu}]_{ij} = k(\mathbf{x}_i, \mathbf{z}_j),$$

$$p(\mathbf{f}|\mathbf{u}) = \mathcal{N}(\mathbf{f} | \mathbf{K}_{fu}\mathbf{K}_u^{-1}\mathbf{u}, \mathbf{K}_f - \mathbf{K}_{fu}\mathbf{K}_u^{-1}\mathbf{K}_{fu}^\top).$$

- How to define $q(\mathbf{u})$?

Sparse approximations

- What we need to compute the ELBO \mathcal{L} ?

$$\log p(\mathbf{y}) \geq \mathcal{L} = \mathbb{E}_{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})}[\log p(\mathbf{y}|\mathbf{f})] - \text{KL}(q(\mathbf{u})||p(\mathbf{u})).$$

- Since the joint $p(\mathbf{f}, \mathbf{u})$ is Gaussian, the conditional $p(\mathbf{f}|\mathbf{u})$ is also Gaussian and given by (with omitted dependence on \mathbf{X} and \mathbf{z}):

$$p(\mathbf{f}, \mathbf{u}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K}_f & \mathbf{K}_{fu} \\ \mathbf{K}_{fu}^\top & \mathbf{K}_u \end{bmatrix}\right), \quad [\mathbf{K}_{fu}]_{ij} = k(\mathbf{x}_i, \mathbf{z}_j),$$

$$p(\mathbf{f}|\mathbf{u}) = \mathcal{N}(\mathbf{f} | \mathbf{K}_{fu}\mathbf{K}_u^{-1}\mathbf{u}, \mathbf{K}_f - \mathbf{K}_{fu}\mathbf{K}_u^{-1}\mathbf{K}_{fu}^\top).$$

- How to define $q(\mathbf{u})$?

Sparse Variational GP

- The original sparse variational approach (Titsias, 2009) analytically optimizes $q(\mathbf{u})$:

$$q(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \hat{\mathbf{m}}, \hat{\mathbf{S}}),$$

$$\hat{\mathbf{m}} = \frac{1}{\sigma_y^2} \mathbf{K}_u \left(\mathbf{K}_u + \frac{1}{\sigma_y^2} \mathbf{K}_{fu}^\top \mathbf{K}_{fu} \right)^{-1} \mathbf{K}_{fu}^\top \mathbf{y},$$

$$\hat{\mathbf{S}} = \mathbf{K}_u \left(\mathbf{K}_u + \frac{1}{\sigma_y^2} \mathbf{K}_{fu}^\top \mathbf{K}_{fu} \right)^{-1} \mathbf{K}_u.$$

- The result is a **collapsed bound** (considering a Gaussian likelihood):

$$\mathcal{L}_{\text{collapsed}} = \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \sigma_y^2 \mathbf{I} + \mathbf{K}_{fz} \mathbf{K}_u^{-1} \mathbf{K}_{fz}^\top) - \frac{1}{2\sigma_y^2} \text{Tr}(\mathbf{K}_f - \mathbf{K}_{fz} \mathbf{K}_u^{-1} \mathbf{K}_{fz}^\top).$$

→ Computation and storage scale with $\mathcal{O}(NM^2)$ and $\mathcal{O}(NM)$.

Sparse Variational GP

- The original sparse variational approach (Titsias, 2009) analytically optimizes $q(\mathbf{u})$:

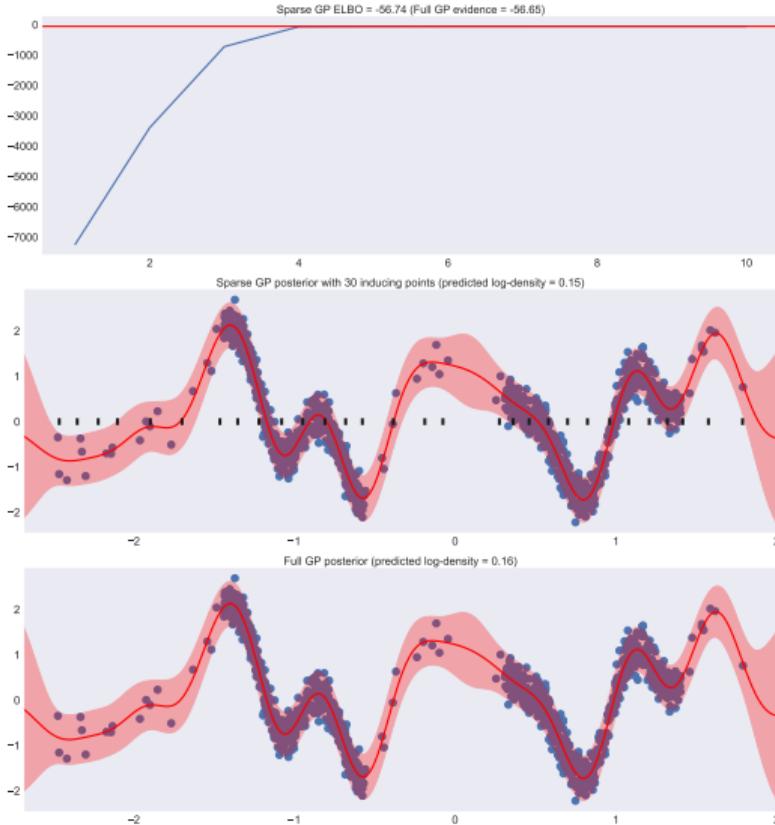
$$\begin{aligned} q(\mathbf{u}) &= \mathcal{N}(\mathbf{u} | \hat{\mathbf{m}}, \hat{\mathbf{S}}), \\ \hat{\mathbf{m}} &= \frac{1}{\sigma_y^2} \mathbf{K}_u \left(\mathbf{K}_u + \frac{1}{\sigma_y^2} \mathbf{K}_{fu}^\top \mathbf{K}_{fu} \right)^{-1} \mathbf{K}_{fu}^\top \mathbf{y}, \\ \hat{\mathbf{S}} &= \mathbf{K}_u \left(\mathbf{K}_u + \frac{1}{\sigma_y^2} \mathbf{K}_{fu}^\top \mathbf{K}_{fu} \right)^{-1} \mathbf{K}_u. \end{aligned}$$

- The result is a **collapsed bound** (considering a Gaussian likelihood):

$$\mathcal{L}_{\text{collapsed}} = \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \sigma_y^2 \mathbf{I} + \mathbf{K}_{fz} \mathbf{K}_u^{-1} \mathbf{K}_{fz}^\top) - \frac{1}{2\sigma_y^2} \text{Tr}(\mathbf{K}_f - \mathbf{K}_{fz} \mathbf{K}_u^{-1} \mathbf{K}_{fz}^\top).$$

→ Computation and storage scale with $\mathcal{O}(NM^2)$ and $\mathcal{O}(NM)$.

Sparse Variational GP



Observations

- The ELBO **monotonically increases**, which always **improves** the approximation.
- The gradients of the ELBO can be computed **analytically**.
- The exact gradients enable the use of **efficient optimization methods**, e.g. L-BFGS.

Stochastic Variational Inference for GPs

- SVI requires a set of **global parameters** and a **factorized ELBO**.
- SVGP explicitly parameterizes $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$ and obtains an **uncollapsed bound** that enables minibatch updates (Hensman *et al.*, 2013):

$$\mathcal{L}_{\text{uncollapsed}} = \sum_{i=1} \mathcal{L}_i - \text{KL}(q(\mathbf{u})||p(\mathbf{u})),$$

$$\text{where } \mathcal{L}_i = \log \mathcal{N}(y_i | \mathbf{k}_i^\top \mathbf{K}_u^{-1} \mathbf{m}, \sigma_y^2) - \frac{1}{2\sigma_y^2} ([\mathbf{K}_f]_{ii} - \mathbf{k}_i^\top \mathbf{K}_u^{-1} \mathbf{k}_i) - \frac{1}{2} \text{Tr} \left(\frac{1}{\sigma_y^2} \mathbf{S} \mathbf{K}_u^{-1} \mathbf{k}_i \mathbf{k}_i^\top \right)$$

$$\text{and } \mathbf{k}_i^\top = [\mathbf{K}_{fz}]_{i:}.$$

→ Given a minibatch size B , computation and storage scale with $\mathcal{O}(BM^2 + M^3)$ and $\mathcal{O}(BM + M^2)$.

Stochastic Variational Inference for GPs

- SVI requires a set of **global parameters** and a **factorized ELBO**.
- SVGP explicitly parameterizes $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$ and obtains an **uncollapsed bound** that enables minibatch updates (Hensman *et al.*, 2013):

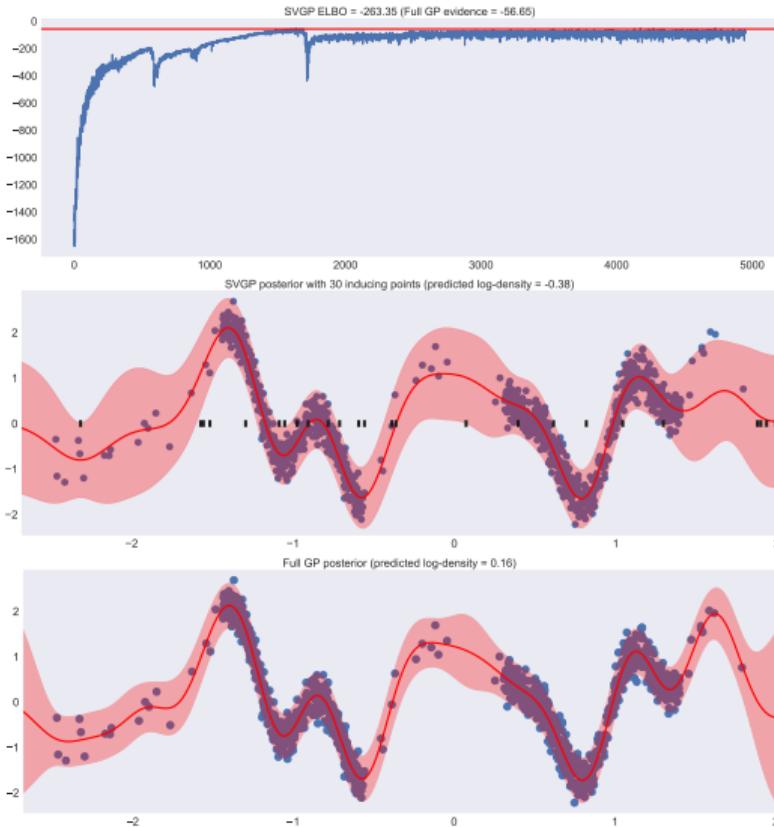
$$\mathcal{L}_{\text{uncollapsed}} = \sum_{i=1} \mathcal{L}_i - \text{KL}(q(\mathbf{u})||p(\mathbf{u})),$$

$$\text{where } \mathcal{L}_i = \log \mathcal{N}(y_i | \mathbf{k}_i^\top \mathbf{K}_u^{-1} \mathbf{m}, \sigma_y^2) - \frac{1}{2\sigma_y^2} ([\mathbf{K}_f]_{ii} - \mathbf{k}_i^\top \mathbf{K}_u^{-1} \mathbf{k}_i) - \frac{1}{2} \text{Tr} \left(\frac{1}{\sigma_y^2} \mathbf{S} \mathbf{K}_u^{-1} \mathbf{k}_i \mathbf{k}_i^\top \right)$$

$$\text{and } \mathbf{k}_i^\top = [\mathbf{K}_{fz}]_{i,:}.$$

→ Given a minibatch size B , computation and storage scale with $\mathcal{O}(BM^2 + M^3)$ and $\mathcal{O}(BM + M^2)$.

Stochastic Variation Inference for GPs



Observations

- The **ELBO does not increase monotonically** due to the use of minibatches.
- The **noisy** gradients of the ELBO can be computed **analytically**.
- More complex **stochastic optimization methods** may be used, e.g. Adam.

Sparse approximations

- In both cases, prediction is performed by integrating out the inducing variables \mathbf{u} using the **variational posterior** $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$:

$$\begin{aligned} q(f_*) &= \int p(f_* | \mathbf{u}) q(\mathbf{u}) d\mathbf{u} \\ &= \int \mathcal{N}(f_* | \mathbf{k}_{*z} \mathbf{K}_u^{-1} \mathbf{u}, K_* - \mathbf{k}_{*z} \mathbf{K}_u^{-1} \mathbf{k}_{*z}^\top) \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{S}) d\mathbf{u} \\ &= \mathcal{N}\left(f_* \middle| \mathbf{k}_{*z} \mathbf{K}_u^{-1} \mathbf{m}, K_* - \mathbf{k}_{*z} \mathbf{K}_u^{-1} (\mathbf{K}_u - \mathbf{S}) \mathbf{K}_u^{-1} \mathbf{k}_{*z}^\top\right). \end{aligned}$$

- Recall that:
 - In the collapsed bound, the moments $\hat{\mathbf{m}}$ and $\hat{\mathbf{S}}$ are optimally obtained.
 - In the uncollapsed bound, the moments \mathbf{m} and \mathbf{S} are variational parameters updated iteratively using the ELBO.

Sparse approximations

- In both cases, prediction is performed by integrating out the inducing variables \mathbf{u} using the **variational posterior** $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}|\mathbf{m}, \mathbf{S})$:

$$\begin{aligned} q(f_*) &= \int p(f_* | \mathbf{u}) q(\mathbf{u}) d\mathbf{u} \\ &= \int \mathcal{N}(f_* | \mathbf{k}_{*z} \mathbf{K}_u^{-1} \mathbf{u}, \mathbf{K}_* - \mathbf{k}_{*z} \mathbf{K}_u^{-1} \mathbf{k}_{*z}^\top) \mathcal{N}(\mathbf{u} | \mathbf{m}, \mathbf{S}) d\mathbf{u} \\ &= \mathcal{N}\left(f_* \mid \mathbf{k}_{*z} \mathbf{K}_u^{-1} \mathbf{m}, \mathbf{K}_* - \mathbf{k}_{*z} \mathbf{K}_u^{-1} (\mathbf{K}_u - \mathbf{S}) \mathbf{K}_u^{-1} \mathbf{k}_{*z}^\top\right). \end{aligned}$$

- Recall that:
 - In the **collapsed bound**, the moments $\hat{\mathbf{m}}$ and $\hat{\mathbf{S}}$ are **optimally obtained**.
 - In the **uncollapsed bound**, the moments \mathbf{m} and \mathbf{S} are variational parameters **updated iteratively** using the ELBO.

Sparse approximations

Summary

- Sparse models enable scaling GP models to **large datasets**.
 - We considered **inducing variables** and corresponding **pseudo-inputs**.
 - integrated out
 - optimized
 - A **variational approach** approximates the posterior of the inducing variables and obtains a lower bound to the model evidence, the **ELBO**.
 - Kernel hyperparameters and variational parameters are **optimized** using the ELBO.
 - Two ELBO's: **collapsed** (Titsias, 2009) and **uncollapsed** (Hensman *et al.*, 2013).
 - Predictions are performed using the **variational posterior** of the inducing points.

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

Current trends on kernel design

- The properties of a zero mean GP model come from its **kernel function**:

$$y_i = f_i + \epsilon_i, \quad f_i = f(\mathbf{x}_i), \quad i \in \{1, \dots, N\},$$

$$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f),$$

$$[\mathbf{K}_f]_{ij} = k(\mathbf{x}_i, \mathbf{x}_j).$$

- Functions that generate a **positive semidefinite** matrix \mathbf{K}_f are valid kernels.
- We have seen some of the most used kernels: Linear, Squared Exponential, Matérn, Rational Quadratic, Periodic, etc.
- How to find suitable kernel functions for a given problem?

Current trends on kernel design

- The properties of a zero mean GP model come from its **kernel function**:

$$\begin{aligned}y_i &= f_i + \epsilon_i, \quad f_i = f(\mathbf{x}_i), \quad i \in \{1, \dots, N\}, \\ \mathbf{f} &\sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f), \\ [\mathbf{K}_f]_{ij} &= k(\mathbf{x}_i, \mathbf{x}_j).\end{aligned}$$

- Functions that generate a **positive semidefinite** matrix \mathbf{K}_f are valid kernels.
- We have seen some of the most used kernels: Linear, Squared Exponential, Matérn, Rational Quadratic, Periodic, etc.
- How to find suitable kernel functions for a given problem?

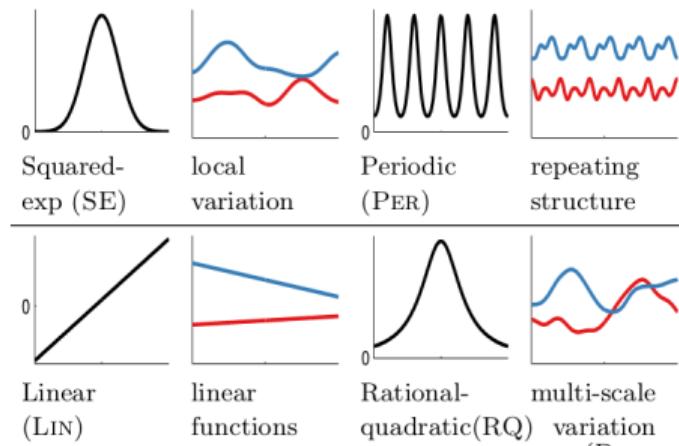
Current trends on kernel design

- The properties of a zero mean GP model come from its **kernel function**:

$$\begin{aligned}y_i &= f_i + \epsilon_i, \quad f_i = f(\mathbf{x}_i), \quad i \in \{1, \dots, N\}, \\ \mathbf{f} &\sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f), \\ [\mathbf{K}_f]_{ij} &= k(\mathbf{x}_i, \mathbf{x}_j).\end{aligned}$$

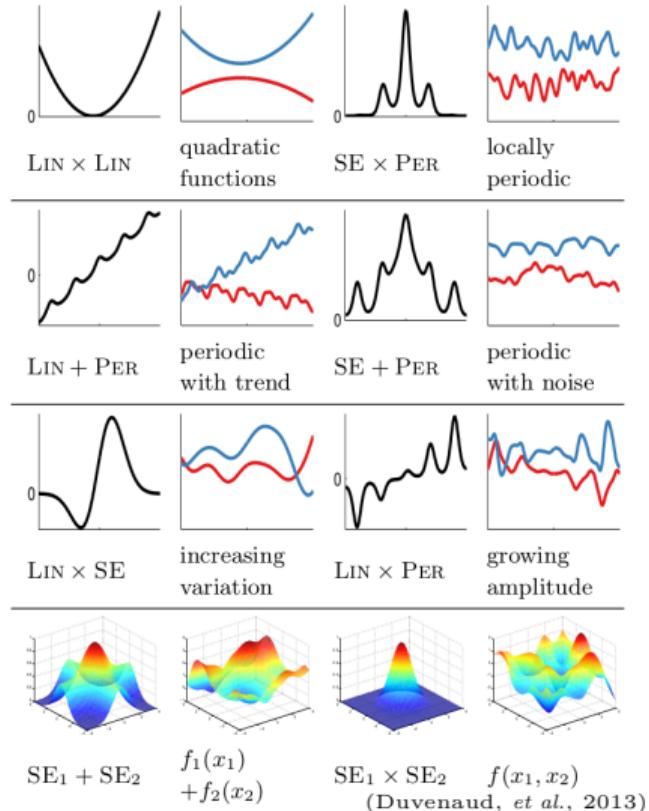
- Functions that generate a **positive semidefinite** matrix \mathbf{K}_f are valid kernels.
- We have seen some of the most used kernels: Linear, Squared Exponential, Matérn, Rational Quadratic, Periodic, etc.
- How to find suitable kernel functions for a given problem?

Composition of base kernels



(Duvenaud, et al., 2013)

- Duvenaud, et al. “**Structure discovery in nonparametric regression through compositional kernel search**”, ICML, 2013.
- Minka, et al. “**Bayesian optimization for automated model selection**”. ICML AutoML Workshop, 2016.
- Kim and Teh. “**Scaling up the Automatic Statistician: Scalable structure discovery using Gaussian processes**”. AISTATS, 2018.
- Teng, et al. “**Scalable variational Bayesian kernel selection for sparse Gaussian process regression**”. AAAI, 2020.

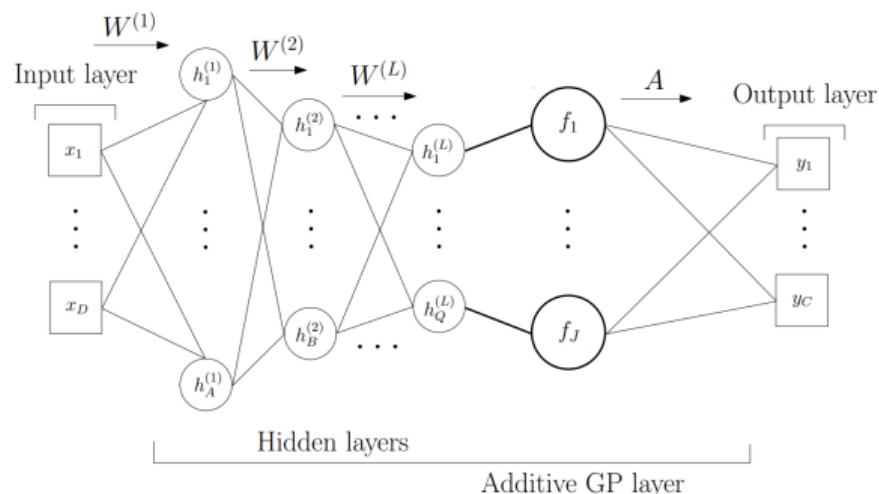


Deep Kernel Learning

- Mapping the input space by an arbitrary function $g(\cdot)$ results in a valid kernel:

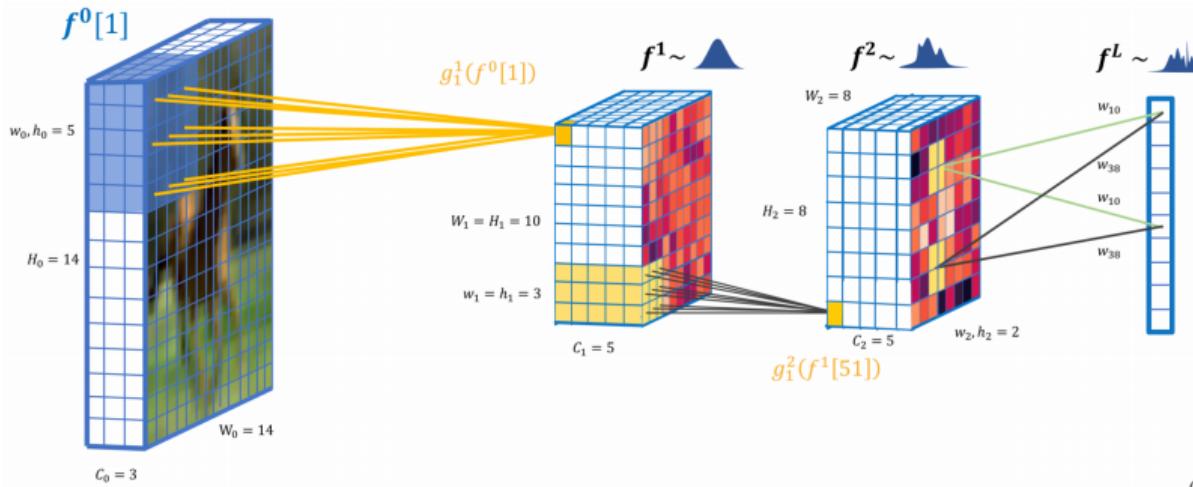
$$\tilde{k}(\mathbf{x}_i, \mathbf{x}_j) = k(g(\mathbf{x}_i), g(\mathbf{x}_j)).$$

- Calandra *et al.* (2016) model $g(\cdot)$ with a **neural network** parameterized by θ that is **jointly optimized** following the gradients of the marginal likelihood $\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial \mathcal{L}}{\partial \tilde{K}} \frac{\partial \tilde{K}}{\partial g_\theta} \frac{\partial g_\theta}{\partial \theta}$.
- Wilson *et al.* (2016) scale the approach to large datasets.



(Wilson, *et al.*, 2016)

Convolutional Gaussian Processes



(Blomqvist, et al., 2018)

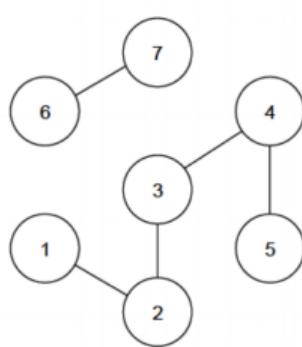
- Sum of patch responses $g(\mathbf{x}^{[p]})$ of a given image \mathbf{x} :

$$f(\mathbf{x}) = \sum_{p=1}^P g(\mathbf{x}^{[p]}), \quad \mathbf{g} \sim \mathcal{N}(\mathbf{0}, k_g(\mathbf{x}, \mathbf{x}')),$$

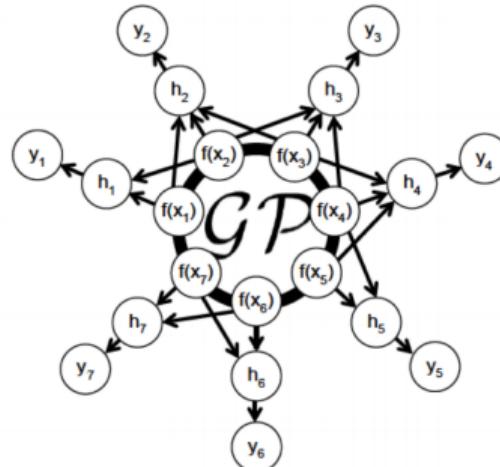
$$\mathbf{f} \sim \mathcal{N} \left(\mathbf{0}, \sum_{p=1}^P \sum_{p'=1}^P k_g(\mathbf{x}^{[p]}, \mathbf{x}^{[p']}) \right).$$

- Van der Wilk, et al. “Convolutional Gaussian processes”. NeurIPS, 2017.
- Blomqvist, et al. “Deep convolutional Gaussian processes”. ECMLPKDD, 2018.
- Dutordoir, et al. “Bayesian image classification with deep convolutional Gaussian processes”. AISTATS, 2020.

Gaussian Processes over graphs



Observed
Graph



- Let $\mathbf{A} \in \{0, 1\}^{N \times N}$ be an **adjacency matrix**, we have:

$$p(\mathbf{y}, \mathbf{h} | \mathbf{X}, \mathbf{A}) = p(\mathbf{h} | \mathbf{X}, \mathbf{A}) \prod_{i=1}^N p(y_i, h_i),$$

$$p(\mathbf{h} | \mathbf{X}, \mathbf{A}) = \mathcal{N}(\mathbf{0}, \mathbf{P} \mathbf{K}_f \mathbf{P}^\top), \quad p(\mathbf{f} | \mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{K}_f),$$

$$\mathbf{P} = (\mathbf{I} + \mathbf{D})^{-1}(\mathbf{I} + \mathbf{A}),$$

where \mathbf{D} is the **degree matrix**.

(Ng, et al., 2018)

- Ng et al. “**Bayesian Semi-supervised Learning with Graph Gaussian Processes**”. NeurIPS, 2018.
- Walker and Glocker. “**Graph convolutional Gaussian processes**”. ICML, 2019.
- Borovitskiy, et al. “**Matérn Gaussian processes on graphs**”. AISTATS, 2021.

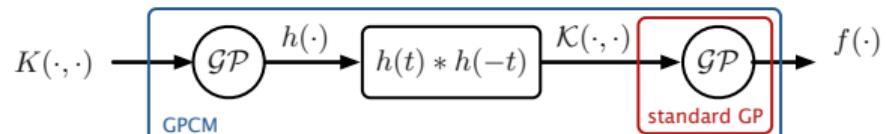
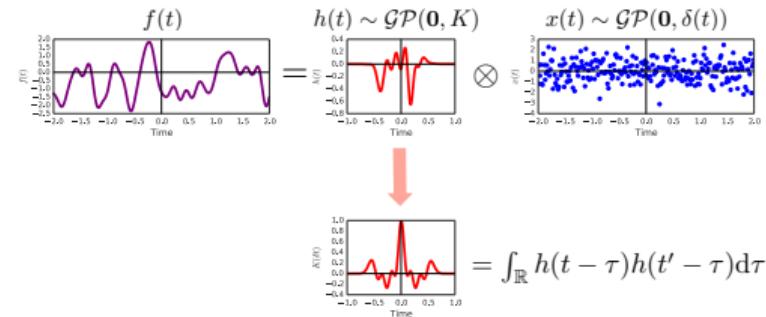
Nonparametric kernels

- Gaussian Process Convolution Model (GPCM)

- i) draw from the filter function $h \sim \mathcal{GP}(\mathbf{0}, K)$;
- ii) convolve it with a white noise process $x(t)$:

$$f(t) = \int_{\mathbb{R}} h(t - \tau)x(\tau)d\tau.$$

- Conditionally, $f|h \sim \mathcal{GP}(\mathbf{0}, \mathcal{K})$ is a GP with a nonparametric covariance function.
- Closely related to linear time-invariant (LTI) systems, with $h(t)$ acting as an impulse response.



Tobar, et al. “Learning stationary time series using Gaussian processes with nonparametric kernels”, Neurips, 2015.

Current trends on kernel design

Summary

- The choice of the **kernel function** defines the properties of a GP model.
- New kernels can be designed by **manually or automatically combining** base kernels.
- **Deep kernel learning** enables flexible kernels jointly optimized with the other GP hyperparameters.
- Specific data structures can be handled with **convolutional or graph kernels**.
- It is possible to use convolutions of processes to obtain **nonparametric kernels** and work with GPs in the spectral domain.

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

From GPLVM to Deep GPs

- So far we have considered supervised learning tasks of the form:

$$y_i = f_i + \epsilon_i, \quad f_i = f(\mathbf{x}_i),$$
$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f).$$

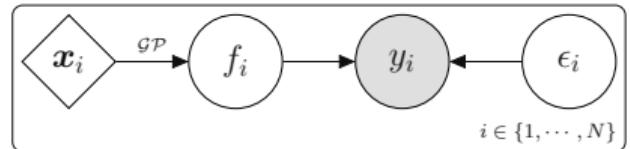
- $y_i|_{i=1}^N$ are **noisy observations**;
- $f_i|_{i=1}^N$ are **latent random variables**;
- $\mathbf{x}_i|_{i=1}^N$ are **deterministic observed inputs**.

- For multidimensional observations $\mathbf{Y} \in \mathbb{R}^{N \times D_y}$ and random inputs \mathbf{X} :



Unobserved inputs X

nonlinear unsupervised projection of Y



Observed inputs X

supervised learning with noisy inputs

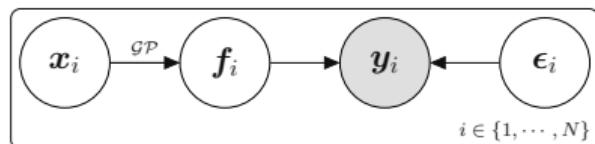
From GPLVM to Deep GPs

- So far we have considered supervised learning tasks of the form:

$$y_i = f_i + \epsilon_i, \quad f_i = f(\mathbf{x}_i), \\ p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{f}|\mathbf{0}, \mathbf{K}_f).$$

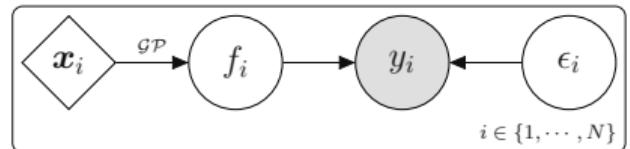
- $y_i|_{i=1}^N$ are **noisy observations**;
- $f_i|_{i=1}^N$ are **latent random variables**;
- $\mathbf{x}_i|_{i=1}^N$ are **deterministic observed inputs**.

- For multidimensional observations $\mathbf{Y} \in \mathbb{R}^{N \times D_y}$ and random inputs \mathbf{X} :



Unobserved inputs X

nonlinear unsupervised projection of \mathbf{Y}



Observed inputs X

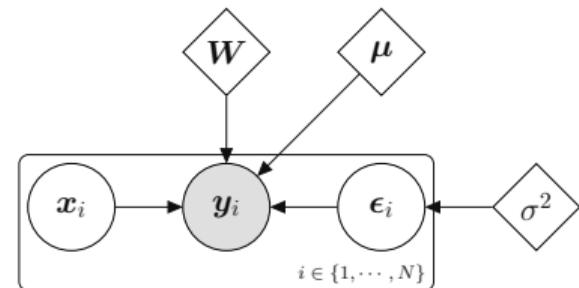
supervised learning with noisy inputs

PPCA - Probabilistic PCA

- The PPCA is one of the simplest models with **continuous latent variables**:

$$\underbrace{p(\mathbf{X})}_{\text{prior}} = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I}), \quad p(\boldsymbol{\epsilon}_i) = \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}),$$

$$\underbrace{p(\mathbf{Y}|\mathbf{X}, \mathbf{W}, \boldsymbol{\mu}, \sigma^2)}_{\text{likelihood}} = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | \mathbf{W}\mathbf{x}_i + \boldsymbol{\mu}, \sigma^2 \mathbf{I}).$$



- Since the PPCA is linear, we have analytical marginal likelihood and posterior:

$$\underbrace{p(\mathbf{Y}|\mathbf{W}, \boldsymbol{\mu}, \sigma^2)}_{\text{marg. likelihood}} = \int \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | \mathbf{W}\mathbf{x}_i + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I}) d\mathbf{x}_i = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}),$$

$$\underbrace{p(\mathbf{X}|\mathbf{Y}, \hat{\mathbf{W}}, \hat{\boldsymbol{\mu}}, \hat{\sigma}^2)}_{\text{posterior}} = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \hat{\mathbf{M}}^{-1} \hat{\mathbf{W}}^\top (\mathbf{y}_i - \hat{\boldsymbol{\mu}}), \hat{\sigma}^2 \hat{\mathbf{M}}^{-1}), \quad \text{where } \hat{\mathbf{M}} = (\hat{\mathbf{W}}^\top \hat{\mathbf{W}} + \hat{\sigma}^2 \mathbf{I})^{-1}.$$

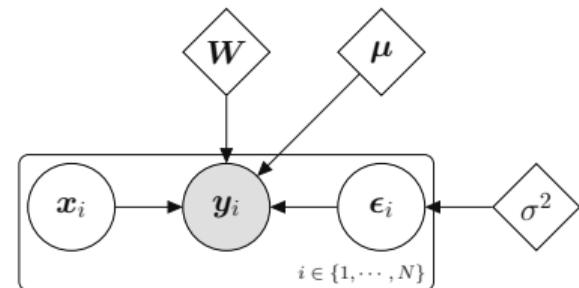
- ML estimates $\hat{\boldsymbol{\mu}}$, $\hat{\mathbf{W}}$ and $\hat{\sigma}^2$ can be obtained in closed form or via EM algorithm.

PPCA - Probabilistic PCA

- The PPCA is one of the simplest models with **continuous latent variables**:

$$\underbrace{p(\mathbf{X})}_{\text{prior}} = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I}), \quad p(\boldsymbol{\epsilon}_i) = \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}),$$

$$\underbrace{p(\mathbf{Y}|\mathbf{X}, \mathbf{W}, \boldsymbol{\mu}, \sigma^2)}_{\text{likelihood}} = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | \mathbf{W}\mathbf{x}_i + \boldsymbol{\mu}, \sigma^2 \mathbf{I}).$$



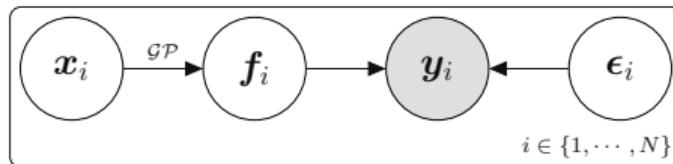
- Since the PPCA is **linear**, we have analytical marginal likelihood and posterior:

$$\underbrace{p(\mathbf{Y}|\mathbf{W}, \boldsymbol{\mu}, \sigma^2)}_{\text{marg. likelihood}} = \int \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | \mathbf{W}\mathbf{x}_i + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I}) d\mathbf{x}_i = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}),$$

$$\underbrace{p(\mathbf{X}|\mathbf{Y}, \hat{\mathbf{W}}, \hat{\boldsymbol{\mu}}, \hat{\sigma}^2)}_{\text{posterior}} = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \hat{\mathbf{M}}^{-1} \hat{\mathbf{W}}^\top (\mathbf{x}_i - \hat{\boldsymbol{\mu}}), \sigma^2 \hat{\mathbf{M}}^{-1}), \quad \text{where } \hat{\mathbf{M}} = (\hat{\mathbf{W}}^\top \hat{\mathbf{W}} + \hat{\sigma}^2 \mathbf{I})^{-1}.$$

- ML estimates $\hat{\boldsymbol{\mu}}$, $\hat{\mathbf{W}}$ and $\hat{\sigma}^2$ can be obtained in closed form or via EM algorithm.

GPLVM - Gaussian Process Latent Variable Model

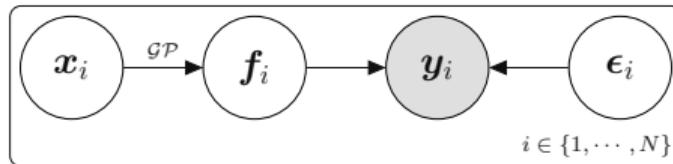


- Inspired by the PPCA formulation, Lawrence (2005) proposed the GPLVM:

$$\begin{aligned} p(\mathbf{X}) &= \prod_{i=1}^N \underbrace{\mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I})}_{\text{input prior}}, & p(\mathbf{F} | \mathbf{X}) &= \prod_{d=1}^{D_y} \underbrace{\mathcal{N}(\mathbf{f}_{:d} | \mathbf{0}, \mathbf{K}_f)}_{\text{GP prior}}, \\ p(\mathbf{Y} | \mathbf{F}, \mathbf{X}) &= \prod_{d=1}^{D_y} \underbrace{\mathcal{N}(\mathbf{y}_{:d} | \mathbf{f}_{:d}, \sigma^2 \mathbf{I})}_{\text{likelihood}} \mathcal{N}(\mathbf{f}_{:d} | \mathbf{0}, \mathbf{K}_f), \\ p(\mathbf{Y} | \mathbf{X}) &= \prod_{d=1}^{D_y} \underbrace{\mathcal{N}(\mathbf{y}_{:d} | \mathbf{0}, \mathbf{K}_f + \sigma^2 \mathbf{I})}_{\text{marginal likelihood}}. \end{aligned}$$

- Marginalization of \mathbf{F} is analytic, but of the random input \mathbf{X} is not.
→ \mathbf{X} appears in a nonlinear way inside the covariance matrix \mathbf{K}_f .

GPLVM - Gaussian Process Latent Variable Model

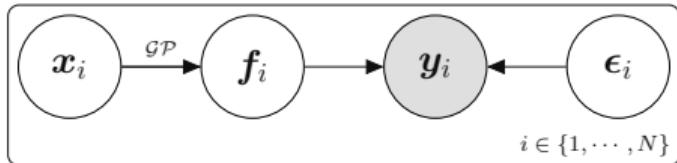


- Inspired by the PPCA formulation, Lawrence (2005) proposed the GPLVM:

$$\begin{aligned} p(\mathbf{X}) &= \prod_{i=1}^N \underbrace{\mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I})}_{\text{input prior}}, & p(\mathbf{F} | \mathbf{X}) &= \prod_{d=1}^{D_y} \underbrace{\mathcal{N}(\mathbf{f}_{:d} | \mathbf{0}, \mathbf{K}_f)}_{\text{GP prior}}, \\ p(\mathbf{Y} | \mathbf{F}, \mathbf{X}) &= \prod_{d=1}^{D_y} \underbrace{\mathcal{N}(\mathbf{y}_{:d} | \mathbf{f}_{:d}, \sigma^2 \mathbf{I})}_{\text{likelihood}} \mathcal{N}(\mathbf{f}_{:d} | \mathbf{0}, \mathbf{K}_f), \\ p(\mathbf{Y} | \mathbf{X}) &= \prod_{d=1}^{D_y} \underbrace{\mathcal{N}(\mathbf{y}_{:d} | \mathbf{0}, \mathbf{K}_f + \sigma^2 \mathbf{I})}_{\text{marginal likelihood}}. \end{aligned}$$

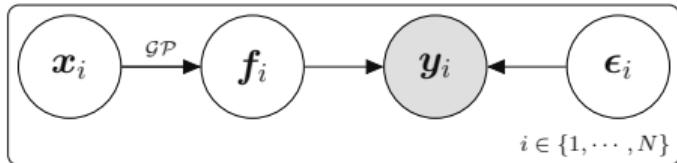
- Marginalization of \mathbf{F} is analytic, but of the random input \mathbf{X} is not.
→ \mathbf{X} appears in a nonlinear way inside the covariance matrix \mathbf{K}_f .

GPLVM - Gaussian Process Latent Variable Model



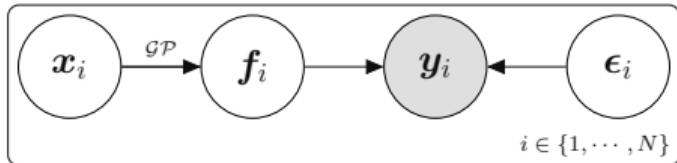
- Lawrence (2005) avoids the intractability by optimizing \mathbf{X} in a MAP fashion.
- Titsias and Lawrence (2010) follow a variational approximation to marginalize \mathbf{X} .
 - The posterior of \mathbf{X} is approximated by $q(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(x_i | m_i, S_i)$.
 - Enables tuning of the latent space dimension via ARD.
 - More robust to overfitting.
- Contrary to standard GPs, the GPLVM is a **generative model**.
- Generation of new observations is easy, since the GP prior maps from x_* to y_* .
- Latent projection given a new observation y_* is nontrivial.
 - E.g. neural networks can be used in an autoencoder fashion (Dai *et al.*, 2016).

GPLVM - Gaussian Process Latent Variable Model



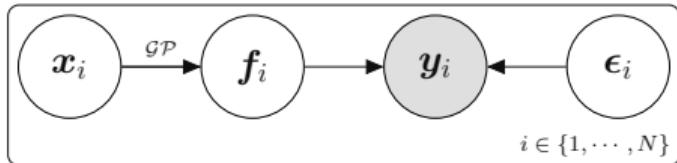
- Lawrence (2005) avoids the intractability by optimizing \mathbf{X} in a MAP fashion.
- Titsias and Lawrence (2010) follow a variational approximation to marginalize \mathbf{X} .
 - The posterior of \mathbf{X} is approximated by $q(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{m}_i, \mathbf{S}_i)$.
 - Enables tuning of the latent space dimension via ARD.
 - More robust to overfitting.
- Contrary to standard GPs, the GPLVM is a **generative model**.
- Generation of new observations is easy, since the GP prior maps from \mathbf{x}_* to \mathbf{y}_* .
- Latent projection given a new observation \mathbf{y}_* is nontrivial.
 - E.g. neural networks can be used in an autoencoder fashion (Dai *et al.*, 2016).

GPLVM - Gaussian Process Latent Variable Model



- Lawrence (2005) avoids the intractability by optimizing \mathbf{X} in a MAP fashion.
- Titsias and Lawrence (2010) follow a variational approximation to marginalize \mathbf{X} .
 - The posterior of \mathbf{X} is approximated by $q(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{m}_i, \mathbf{S}_i)$.
 - Enables tuning of the latent space dimension via ARD.
 - More robust to overfitting.
- Contrary to standard GPs, the GPLVM is a **generative model**.
- Generation of new observations is easy, since the GP prior maps from \mathbf{x}_* to \mathbf{y}_* .
- Latent projection given a new observation \mathbf{y}_* is nontrivial.
 - E.g. neural networks can be used in an autoencoder fashion (Dai *et al.*, 2016).

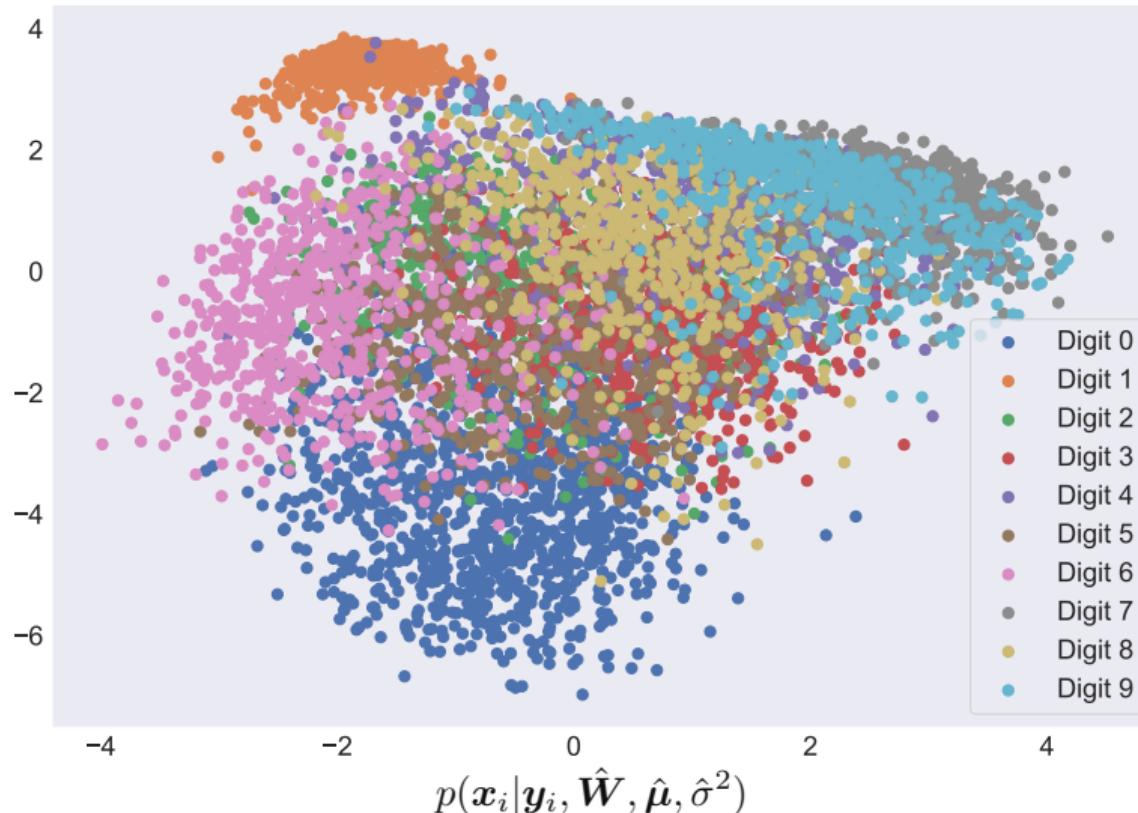
GPLVM - Gaussian Process Latent Variable Model



- Lawrence (2005) avoids the intractability by optimizing \mathbf{X} in a MAP fashion.
- Titsias and Lawrence (2010) follow a variational approximation to marginalize \mathbf{X} .
 - The posterior of \mathbf{X} is approximated by $q(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{m}_i, \mathbf{S}_i)$.
 - Enables tuning of the latent space dimension via ARD.
 - More robust to overfitting.
- Contrary to standard GPs, the GPLVM is a **generative model**.
- Generation of new observations is easy, since the GP prior maps from \mathbf{x}_* to \mathbf{y}_* .
- Latent projection given a new observation \mathbf{y}_* is nontrivial.
 - E.g. neural networks can be used in an autoencoder fashion (Dai *et al.*, 2016).

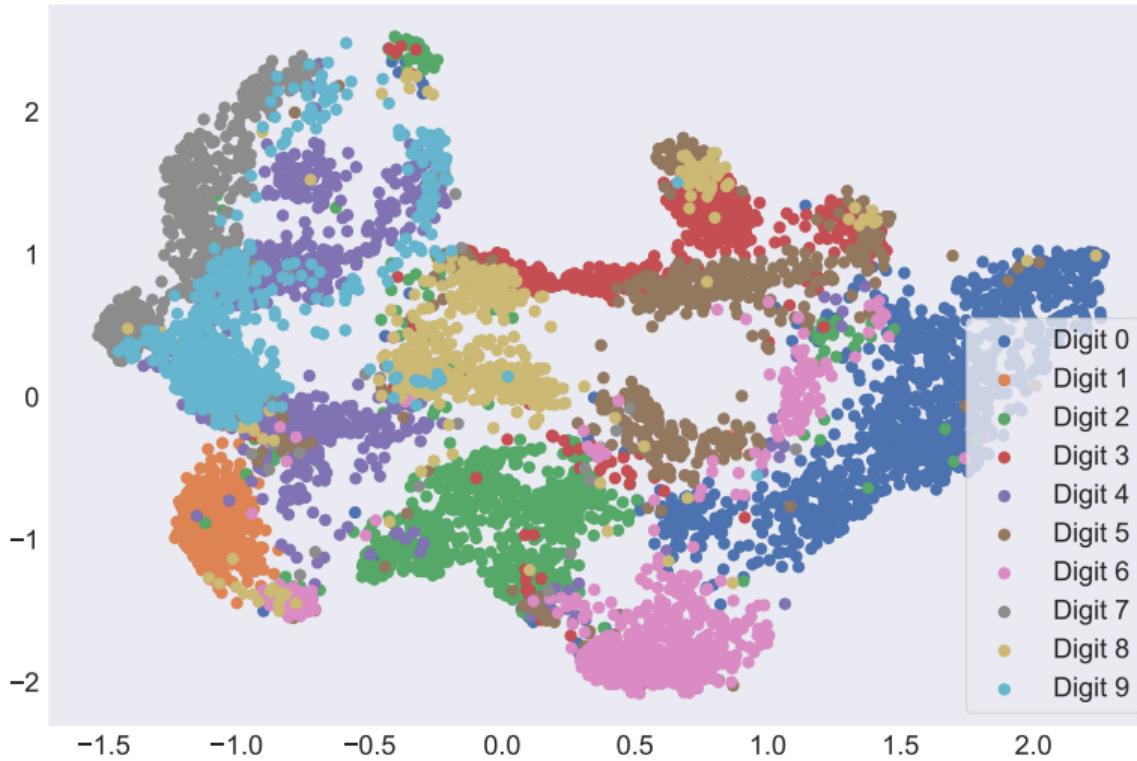
USPS digits - 7291 training samples

256 → 2 dimension projection with PCA



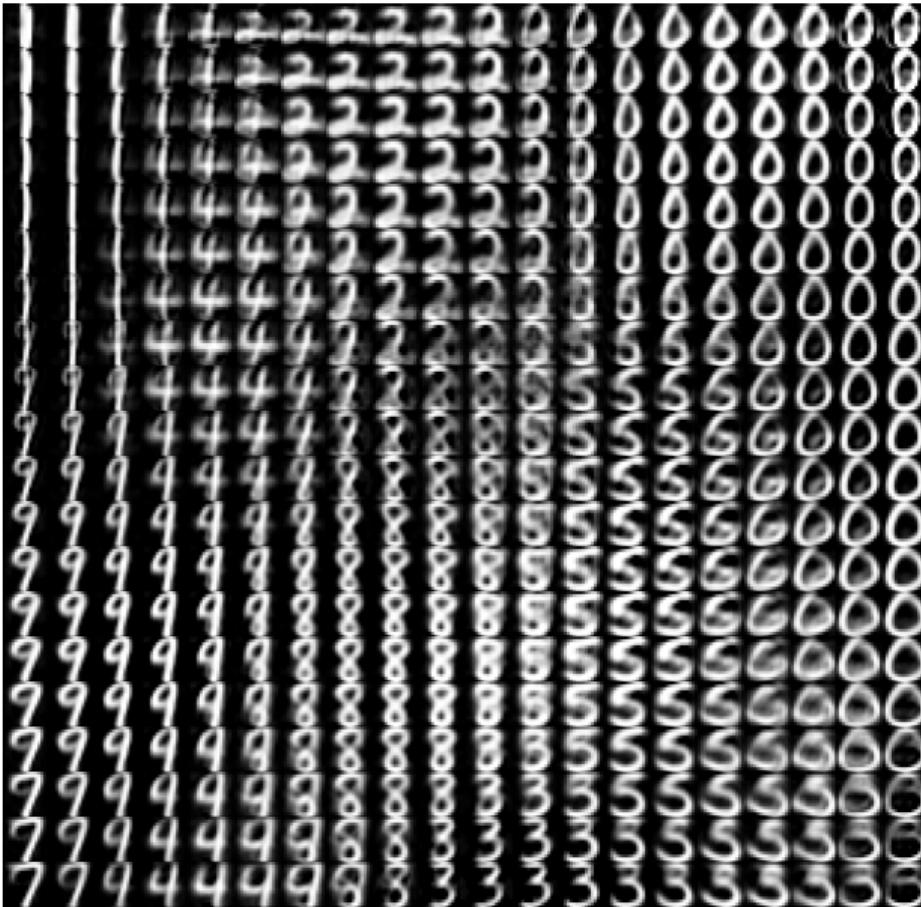
USPS digits - 7291 training samples

256 → 2 dimension projection with GPLVM (50 inducing points)

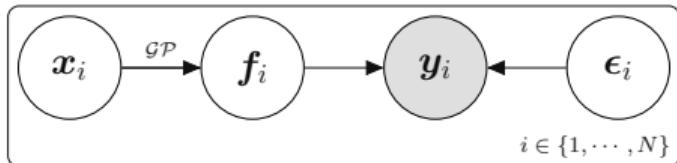


$$p(\mathbf{x}_i|\mathbf{y}_i) \approx q(\mathbf{x}_i|\mathbf{y}_i) = \mathcal{N}(\mathbf{x}_i|\mathbf{m}_i, \mathbf{S}_i).$$

USPS digits - GPLVM samples

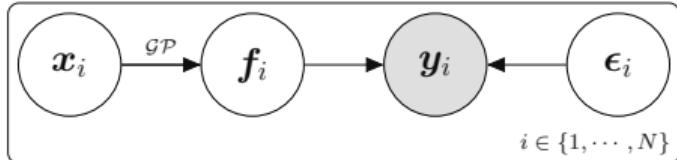


GPLVM - Gaussian Process Latent Variable Model



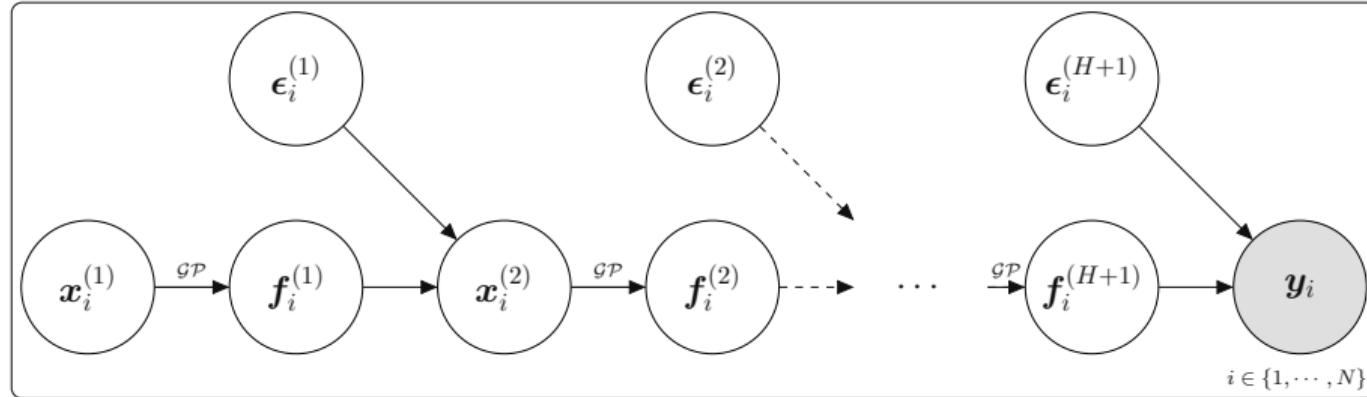
- Instead of a simple prior $p(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I})$, we could have different choices:
 - **Discriminative priors** for classification (Urtasun and Darrell, 2007);
 - **Dynamical priors** with temporal dependence (Wang *et al.*, 2007; Damianou *et al.*, 2011; Frigola *et al.*, 2014; Mattos *et al.*, 2016);
 - Mixture of **visible and missing** data (Damianou and Lawrence, 2015).
 - **Conditional** variables (Dutordoir *et al.*, 2016);
 - **Other(s) GP prior(s)** (Lawrence and Moore, 2007; Damianou and Lawrence, 2013; Mattos *et al.*, 2016; Bui *et al.*, 2016; Salimbeni and Deisenroth, 2017).

GPLVM - Gaussian Process Latent Variable Model



- Instead of a simple prior $p(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I})$, we could have different choices:
 - **Discriminative priors** for classification (Urtasun and Darrell, 2007);
 - **Dynamical priors** with temporal dependence (Wang *et al.*, 2007; Damianou *et al.*, 2011; Frigola *et al.*, 2014; Mattos *et al.*, 2016);
 - Mixture of **visible and missing** data (Damianou and Lawrence, 2015).
 - **Conditional** variables (Dutordoir *et al.*, 2016);
 - **Other(s) GP prior(s)** (Lawrence and Moore, 2007; Damianou and Lawrence, 2013; Mattos *et al.*, 2016; Bui *et al.*, 2016; Salimbeni and Deisenroth, 2017).

Deep Gaussian Processes

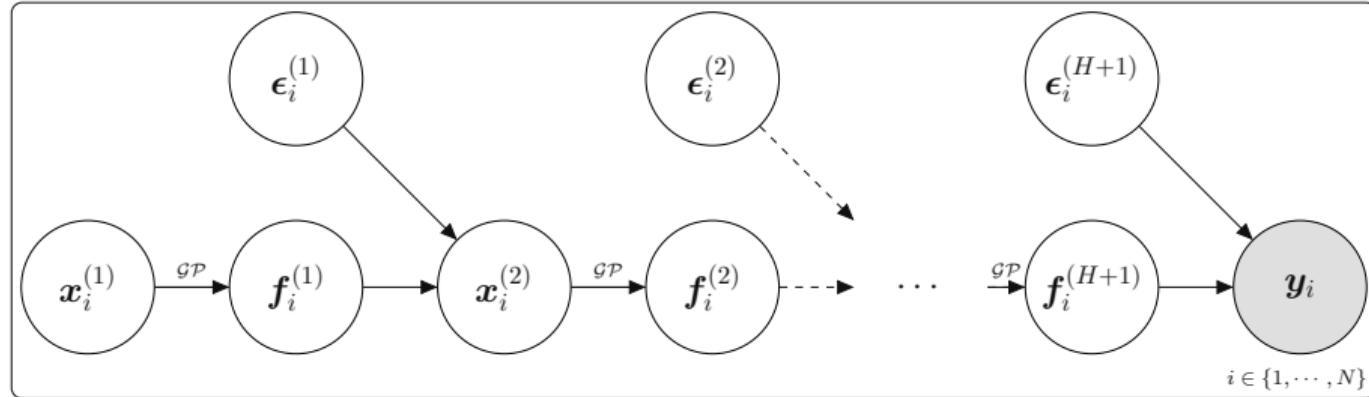


- A **hierarchy** of H GP priors gives rise to a **deep** GP model with H hidden layers:

$$\begin{aligned} \mathbf{y}_i &= f^{(H+1)}(\mathbf{x}_i^{(H+1)}) + \epsilon_i^{(H+1)}, & \mathbf{f}_{:d}^{(H+1)} &\sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f^{(H+1)}) \\ \mathbf{x}_i^{(h+1)} &= f^{(h)}(\mathbf{x}_i^{(h)}) + \epsilon_i^{(h)}, & \mathbf{f}_{:d}^{(h)} &\sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f^{(h)}), \quad 1 \leq h \leq H. \end{aligned}$$

- The intractabilities require approximate inference and inducing point approaches.

Deep Gaussian Processes

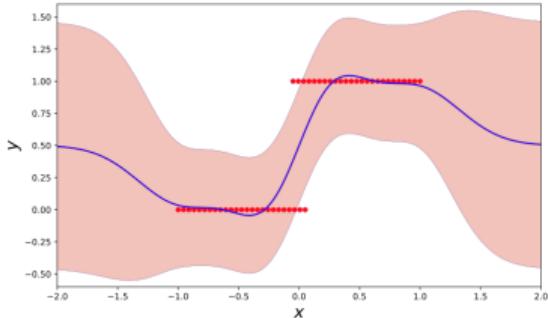


- A **hierarchy** of H GP priors gives rise to a **deep** GP model with H hidden layers:

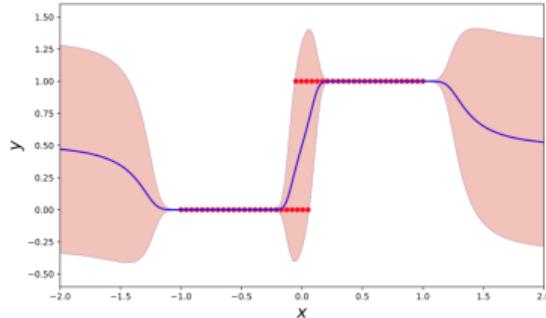
$$\mathbf{y}_i = f^{(H+1)}(\mathbf{x}_i^{(H+1)}) + \epsilon_i^{(H+1)}, \quad \mathbf{f}_{:d}^{(H+1)} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f^{(H+1)})$$
$$\mathbf{x}_i^{(h+1)} = f^{(h)}(\mathbf{x}_i^{(h)}) + \epsilon_i^{(h)}, \quad \mathbf{f}_{:d}^{(h)} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_f^{(h)}), \quad 1 \leq h \leq H.$$

- The intractabilities require approximate inference and inducing point approaches.

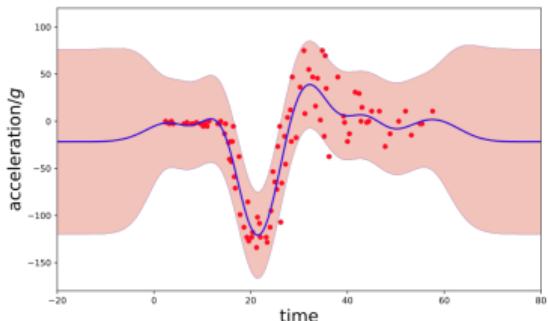
Deep Gaussian Processes



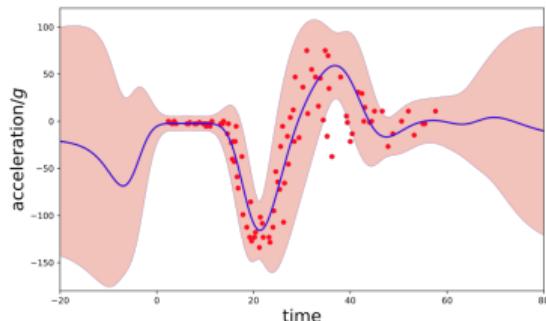
(a) Shallow GP.



(b) Deep GP.



(c) Shallow GP.



(d) Deep GP.

(Lawrence, 2020)

Deep Gaussian Processes

- But what are the advantages of deep GPs?
 - The resulting model **cannot be explained** by a single GP.
 - Multiple layers **alleviate the choice** of the kernel function.
 - Enables **hierarchical** feature learning.
 - Deep GPs are **related to deep neural networks** with infinite hidden units.
- Supported by modern frameworks, such as **GPyTorch** and **GPflux**.
- Research on deep GPs has a wide range of topics:
 - Variational Auto-encoded Deep Gaussian Processes (Dai *et al.*, 2016);
 - Recurrent Gaussian Processes (Mattos *et al.*, 2016);
 - Doubly Stochastic Variational Inference for Deep GPs (Salimbeni and Deisenroth, 2017);
 - Deep Convolutional GPs (Blomqvist *et al.*, 2018);
 - Deep GPs with Importance-Weighted Variational Inference (Salimbeni *et al.*, 2019);
 - Stochastic Deep Gaussian Processes over Graphs (Li *et al.*, 2020);
 - Global Inducing Point Variational Posteriors for BNNs and DGPs (Ober and Aitchison, 2021).

Deep Gaussian Processes

- But what are the advantages of deep GPs?
 - The resulting model **cannot be explained** by a single GP.
 - Multiple layers **alleviate the choice** of the kernel function.
 - Enables **hierarchical** feature learning.
 - Deep GPs are **related to deep neural networks** with infinite hidden units.
- Supported by modern frameworks, such as **GPyTorch** and **GPflux**.
- Research on deep GPs has a wide range of topics:
 - Variational Auto-encoded Deep Gaussian Processes (Dai *et al.*, 2016);
 - Recurrent Gaussian Processes (Mattos *et al.*, 2016);
 - Doubly Stochastic Variational Inference for Deep GPs (Salimbeni and Deisenroth, 2017);
 - Deep Convolutional GPs (Blomqvist *et al.*, 2018);
 - Deep GPs with Importance-Weighted Variational Inference (Salimbeni *et al.*, 2019);
 - Stochastic Deep Gaussian Processes over Graphs (Li *et al.*, 2020);
 - Global Inducing Point Variational Posteriors for BNNs and DGPs (Ober and Aitchison, 2021).

Deep Gaussian Processes

- But what are the advantages of deep GPs?
 - The resulting model **cannot be explained** by a single GP.
 - Multiple layers **alleviate the choice** of the kernel function.
 - Enables **hierarchical** feature learning.
 - Deep GPs are **related to deep neural networks** with infinite hidden units.
- Supported by modern frameworks, such as **GPyTorch** and **GPflux**.
- Research on deep GPs has a wide range of topics:
 - Variational Auto-encoded Deep Gaussian Processes (Dai *et al.*, 2016);
 - Recurrent Gaussian Processes (Mattos *et al.*, 2016);
 - Doubly Stochastic Variational Inference for Deep GPs (Salimbeni and Deisenroth, 2017);
 - Deep Convolutional GPs (Blomqvist *et al.*, 2018);
 - Deep GPs with Importance-Weighted Variational Inference (Salimbeni *et al.*, 2019);
 - Stochastic Deep Gaussian Processes over Graphs (Li *et al.*, 2020);
 - Global Inducing Point Variational Posteriors for BNNs and DGPs (Ober and Aitchison, 2021).

From GPLVM to Deep GPs

Summary

- GPLVM enables **unsupervised learning** with GP priors in a **generative setting**.
- GP learning with **missing input data** uses the same GPLVM inference tools.
- A **hierarchy** of GPLVM blocks results in a **deep GP model**.
- The multilayer **composition of processes** cannot be represented by a single GP.
- Bayesian inference with GPLVMs and deep GPs is hard, but some **open source modern frameworks** are available.

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

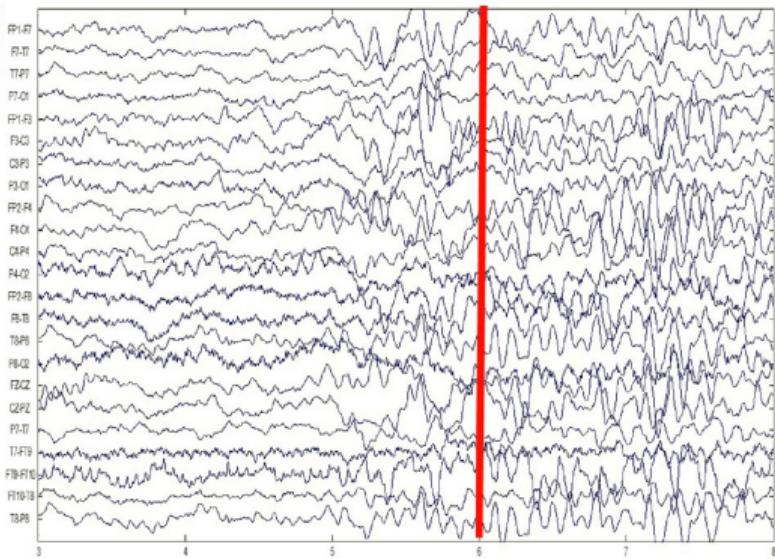
From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

The need for multivariate processing

shared sources of uncertainty, relationship across measurements, only some channels are observed at some locations (incomplete measurements)



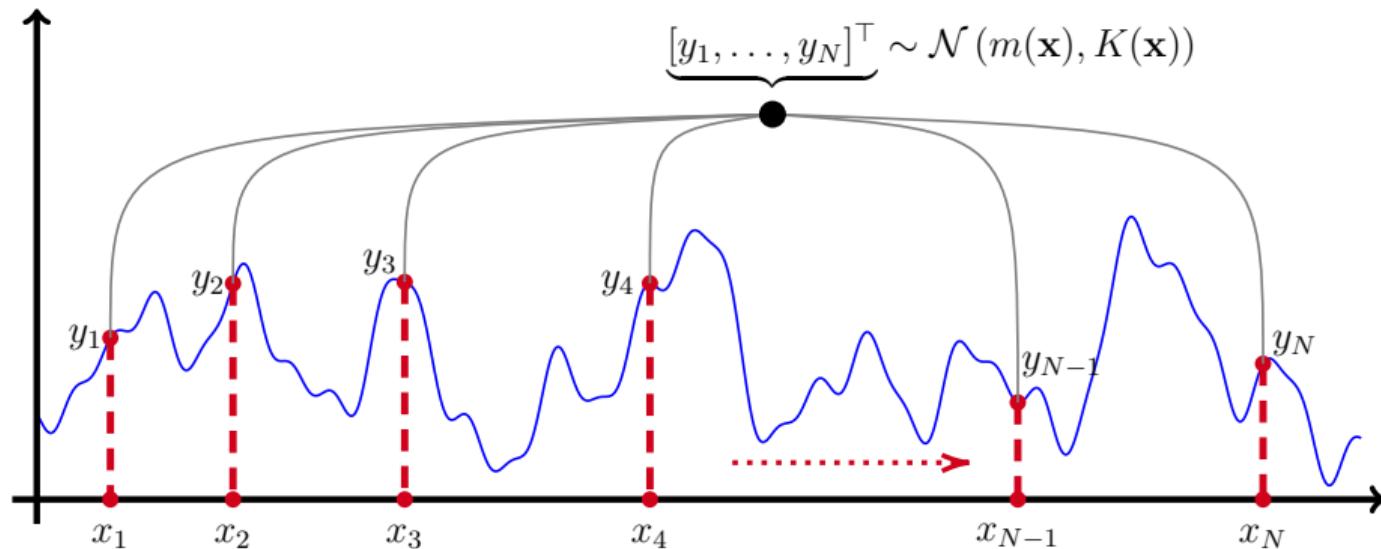
Left: EMOTIV EPOC+, <https://lucid.me/blog/wed-love-try-emotiv-epoch>

Right: Xun, G., Jia, X. & Zhang, A. Detecting epileptic seizures with electroencephalogram via a context-learning model. BMC Med Inform Decis Mak 16, 70 (2016)

Gaussian processes

a 1-slide reminder

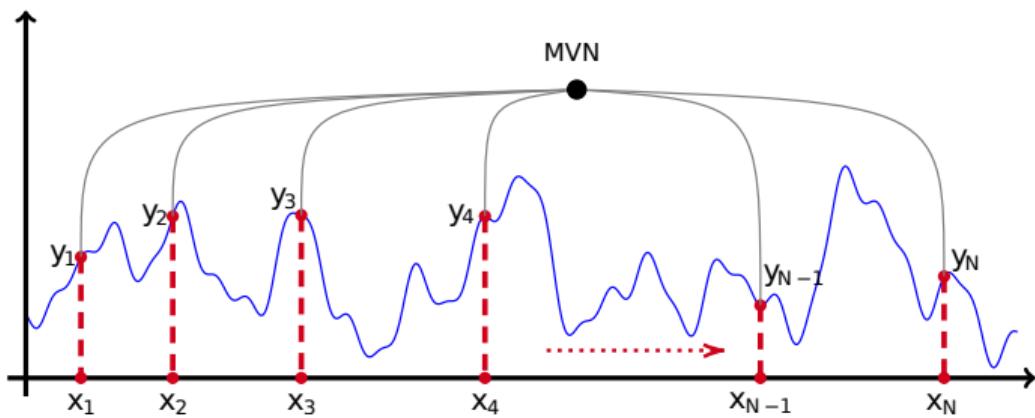
Definition: A GP is a stochastic process such that any finite collection of values follows a multivariate normal distribution.



Notation: $f \sim \mathcal{GP}(\mu, K) \iff f(\mathbf{x}) \sim \mathcal{N}(m(\mathbf{x}), K(\mathbf{x}, \mathbf{x})), \forall \mathbf{x} \in \mathcal{X}^n, n \in \mathbb{N}$

Def: Multioutput Gaussian processes

Definition: An MOGP is a **vector-valued** stochastic process such that any finite collection of values **vectors** follow a multivariate normal distribution.

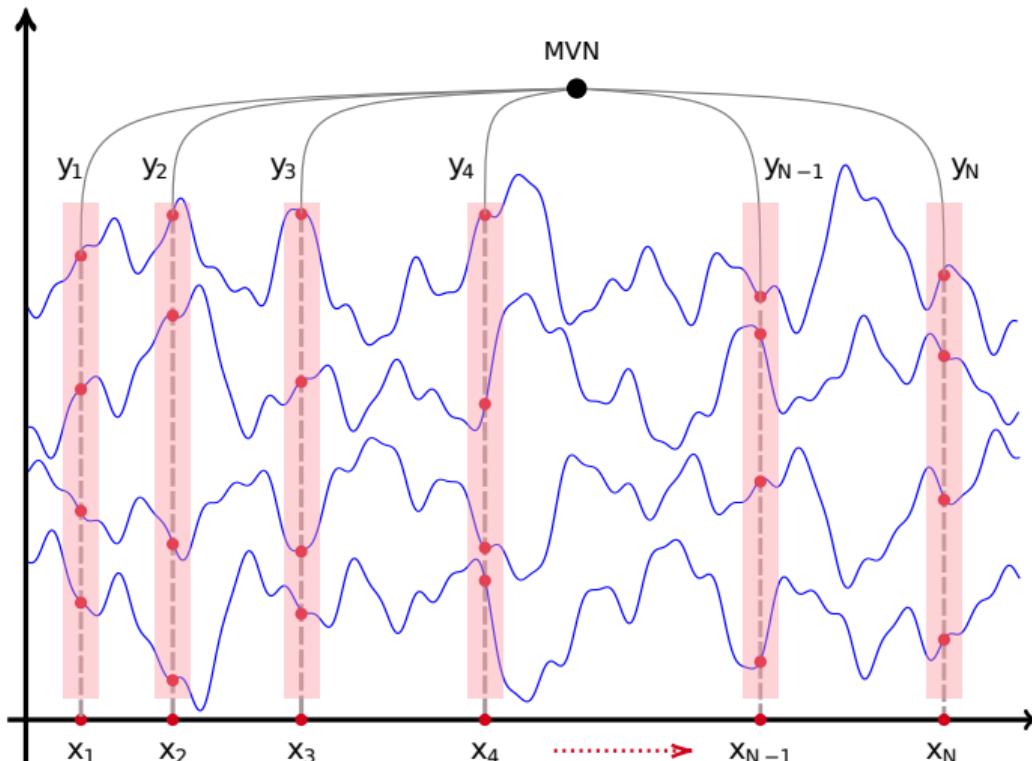


However, due to the marginalisation property of the MVN, vectors need not be chosen in full. Therefore:

Definition: An MOGP is a vector-valued stochastic process such that any finite collection of values, **chosen from any channel at any time**, follows a multivariate normal distribution.

Def: Multioutput Gaussian processes

Definition: An MOGP is a **vector-valued** stochastic process such that any finite collection of values **vectors** follow a multivariate normal distribution.

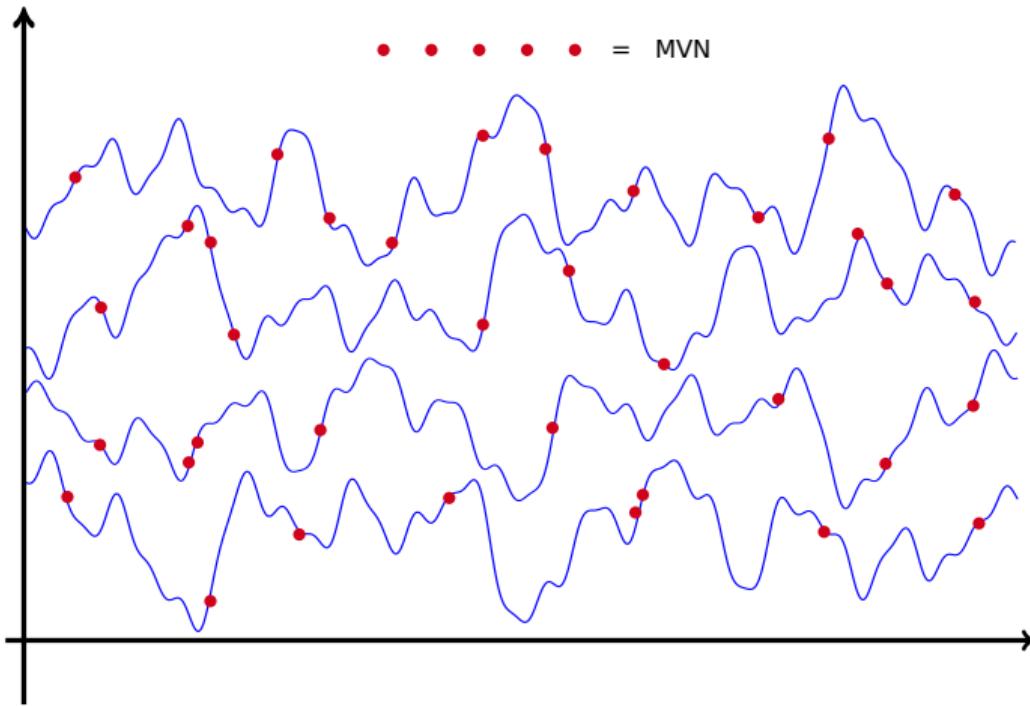


However, due to the marginalisation property of the MVN, vectors need not be chosen in full. Therefore:

Definition: An MOGP is a vector-valued stochastic process such that any finite collection of values, **chosen from any channel at any time**, follows a multivariate normal distribution.

Def: Multioutput Gaussian processes

Definition: An MOGP is a **vector-valued** stochastic process such that any finite collection of values **vectors** follow a multivariate normal distribution.



However, due to the marginalisation property of the MVN, vectors need not be chosen in full. Therefore:

Definition: An MOGP is a vector-valued stochastic process such that any finite collection of values, **chosen from any channel at any time**, follows a multivariate normal distribution.

How do we specify the covariance of an MOGP?

Let us introduce some notation:

- $f = [f_1, f_2, \dots, f_m]^\top \sim \mathcal{MOP}(\mu, K)$
- Number of channels = m

An MOGP is specified by its covariance function which is a four-way array:

$$K : \{1, 2, \dots, m\}^2 \times \mathcal{X}^2 \rightarrow \mathbb{R} \quad (25)$$

$$(i, j, x, x') \mapsto K_{ij}(x, x'), \quad (26)$$

meaning that for a pair of values $f_i(x)$ and $f_j(x')$, we have $\mathbb{V}[f_i(x), f_j(x')] = K_{ij}(x, x')$.

Properties. $\mathcal{K}(x, x') : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^{m \times m}$ defined element-wise by $[\mathcal{K}(x, x')]_{ij} = k_{ij}(x, x')$ is:

- Symmetric, i.e., $\mathcal{K}(x, x') = \mathcal{K}(x', x)^\top, \forall x, x' \in \mathcal{X}$, and
- Positive definite, i.e., $\forall N \in \mathbb{N}, \{c_p\}_{p=1}^N \subseteq \mathbb{R}^m, \{x_p\}_{p=1}^N \subseteq \mathcal{X}$, we have:

$$\sum_{i,j=1}^m \sum_{p,q=1}^N c_{pi} c_{qj} k_{ij}(x_p, x_q) \geq 0. \quad (27)$$

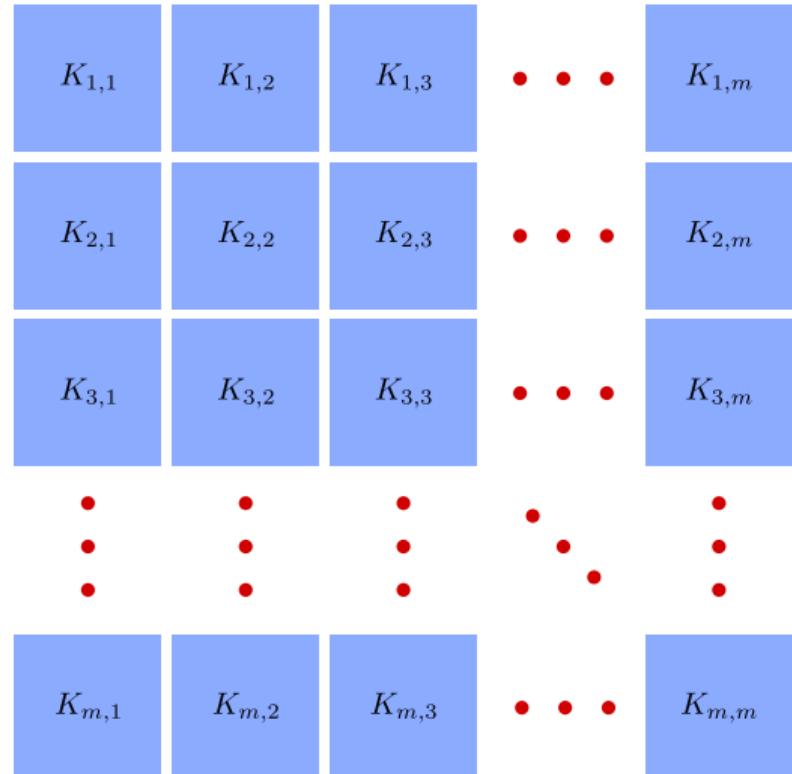
Multioutput covariance

A few observations

Visualisation: Since the kernel has 2 discrete and 2 continuous dimensions, it can be shown by *slicing* the covariance over its discrete dimensions (channels) and show each slice.

Parametrisation: This is the tricky part. In the sliced representation the covariance is not continuous, since it has channel jumps. That makes parametrisation quite difficult.

Stacked representation: This covariance corresponds to the concatenation of all datapoints stacked, where one loses channel information. Useful for coding and understanding positive semidefiniteness via index reshaping.



The assumption of stationarity

From 4D to 3D

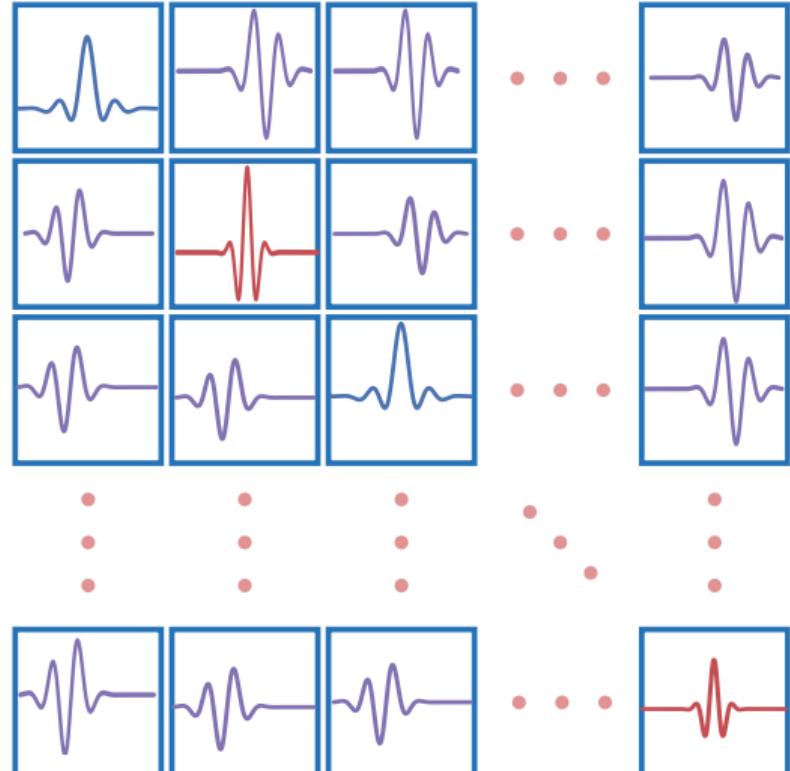
Recall that a covariance is stationary iff

$$K(x, x') = K(x - x'). \quad (28)$$

This condition, for MOGPs, turns into

$$K_{i,j}(x, x') = K_{i,j}(x - x'), \forall i, j = 1, \dots, m. \quad (29)$$

Remark: We've got one fewer dimension! and though it might sound like it isn't much, it'll allow for the design of vector-valued covariance kernels via direct parametrisation (just like we did for the scalar case)



How to design the kernel?

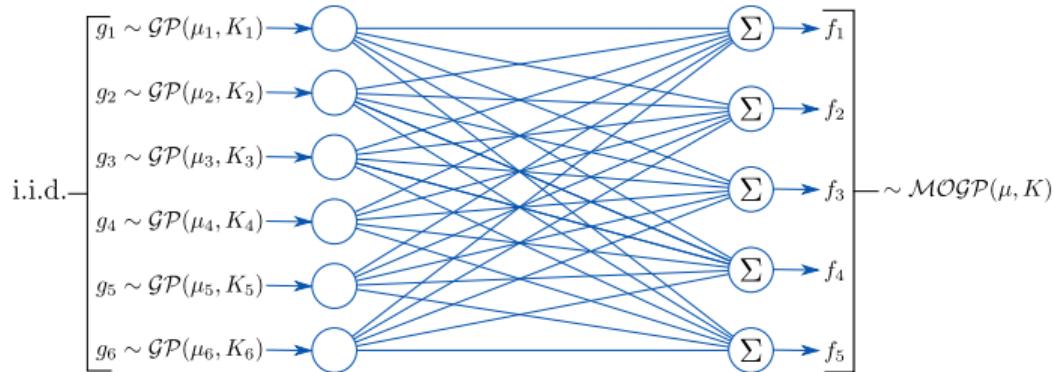
Challenge: Due to the positive semidefiniteness condition, designing MOGP kernels via direct parametrisation is tough if one aims to move from the trivial choices.

- We want kernels to be as expressive as possible
 - Hopefully all structure is discovered (we want to be agnostic)
 - Kernel parametrisation are constrained in a very (parameter-wise) unintuitive way

Let us start with something simple

		Time τ			
		$k_{11}(\tau_1)$	$k_{12}(\tau_1)$	\cdots	$k_{1m}(\tau_1)$
		$k_{11}(\tau_0)$	$k_{12}(\tau_0)$	\cdots	$k_{1m}(\tau_0)$
		$k_{21}(\tau_1)$	$k_{22}(\tau_1)$	\cdots	\vdots
		$k_{11}(\tau)$	$k_{12}(\tau)$	\cdots	\vdots
		\vdots	\vdots	\vdots	\vdots
		$k_{21}(\tau)$	$k_{22}(\tau)$	\cdots	$k_{mm}(\tau_1)$
		\vdots	\vdots	\cdots	\vdots
		$k_{m1}(\tau_0)$	\cdots	\cdots	$k_{mm}(\tau_0)$
		$k_{m1}(\tau)$	\cdots	\cdots	$k_{mm}(\tau)$
Channel		Channel			

MOGP design, first chapter: Induced kernels



Denote:

We have that $f = \theta^\top g$ (or $f_j = \theta_j^\top g$) and thus

- $\theta \in \mathbb{R}^{n \times m}$ the weights
- $g = [g_1, g_2, \dots, g_n]^\top$
- $f = [f_1, f_2, \dots, f_m]^\top$

$$\langle f \rangle = \theta^\top \langle g \rangle = \theta^\top \mu \quad (30)$$

$$\mathbb{V}[f(x), f(x')] = \mathbb{V}[\theta^\top g(x), \theta^\top g(x')] \quad (31)$$

$$= \theta^\top \mathbb{V}[g(x), g(x')] \theta \quad (32)$$

$$= \sum_{q=1}^Q \theta_q \theta_q^\top K_q(x - x') \quad (33)$$

Particular cases of MOGP kernels induced by linear mixing⁴

- **The Linear Model of Coregionalisation**¹: Some of the Q latent GPs, though independent, share the same covariance. Grouping GPs with common covariance, the kernel is

$$\mathbb{V}[f(x), f(x')] = \sum_{q=1}^Q \underbrace{\sum_{r=1}^{R_q} \theta_q^{(r)} \theta_q^{(r)\top}}_{C_q} K_q(x - x'), \quad (35)$$

where $\theta_q^{(r)}$ is the vector of weights associated to the r -th member of the q -th group of inputs. The matrices $C_q, q \in \{1, 2, \dots, m\}$, are called the *coregionalisation matrices*

- **The Intrinsic Model of Coregionalisation IMC**²: All kernels are equal ($Q = 1$ above)
- **The Semiparametric Latent Factor Model**³: The weights θ are free

¹A. G. Journel and C. J. Huijbregts. Mining Geostatistics. Academic Press, London, 1978.

²P. Goovaerts. Geostatistics For Natural Resources Evaluation. Oxford University Press, USA, 1997.

³Y. W. Teh, M. Seeger, and M. I. Jordan. Semiparametric latent factor models. AISTATS, pp. 333–340, 2005.

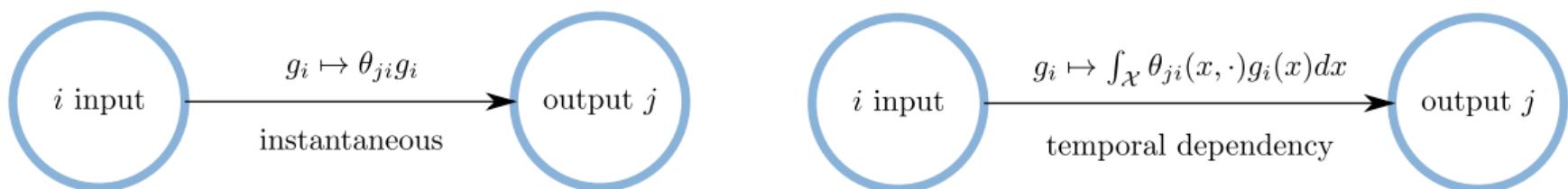
⁴M. Álvarez, L Rossaco, N. Lawrence, Kernels for Vector-Valued Functions: A Review, FTML, 2012.

A drawback of the above construction

The mixing model above is *instantaneous* meaning that the output at x :

- $K(x, x')$ only depends on $K_q(x, x')$ ($\forall q$)
- temporal structure in f is only given by that of g : mixture only introduces channel correlations.
- kernels are separable, i.e., $K_{i,j}(x, x') = B_{i,j}K(x, x')$

$f(x) \sim \mathcal{GP}(m(x), K(x, x'))$ **Fix:** Let us consider a more general linear operator



$$K(f_i(x), f_j(x')) = \sum_q \int \int \theta_{iq}(x - z)\theta_{jq}(x' - z')K_q(z, z')dz dz'$$

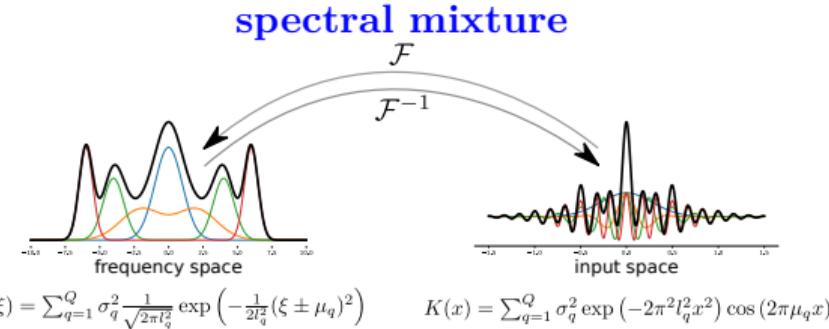
Warning: For the convolution to converge, the kernel has to be square integrable (one particular case are kernels with compact support).

MOGP design, second chapter: spectral parametrisation

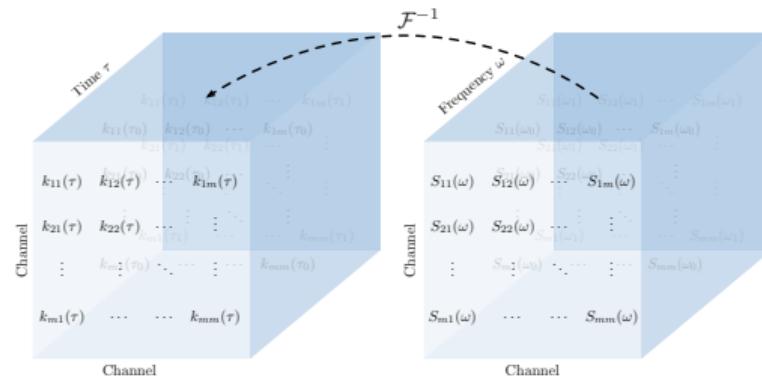
Idea: Unconstrained parametrisation of MO kernel and then transform it - a sort of *reparametrisation trick*.

How: Use spectral mixture but for multioutput, i.e., parametrise the crosspectral power spectrum by a Gaussian mixture and anti-Fourier transform the mixture.

Consequence: We went from "3D" positive semidefiniteness (time) to "2D" positive semidefiniteness (freq) since the third dimension is guaranteed by the anti-Fourier transform



multioutput spectral mixture

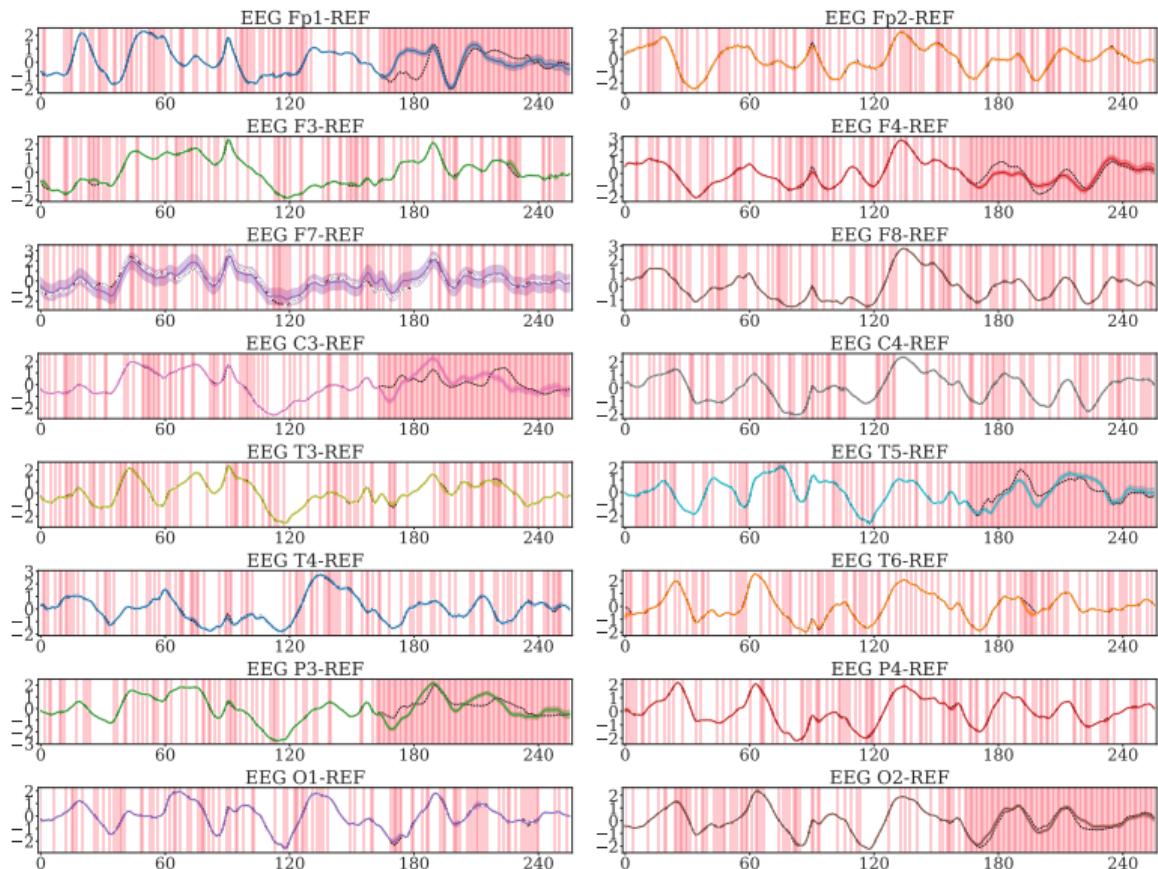
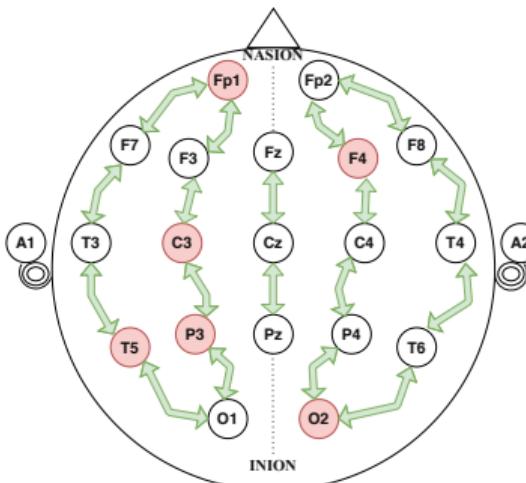


MOSM in action: Electroencephalography

Aim: To detect neonatal epileptic fits, critical for brain development

How: Adjust MOSM, inspect hypers

Findings: Automatic selection of quasistationary windows, change-point detection, nonstationarity



de Wolff, Cuevas & Tobar, MOGPTK
github.com/GAMES-UChile/mogptk

Table of Contents

① Introduction

② Classics

Fundamentals of Bayesian inference

From one to finitely-many Gaussian RVs

Infinitely-many Gaussian RVs: The Gaussian process

Choosing the kernel and learning its parameters

Implementation of a GP

Beyond Gaussian likelihood

③ Contemporary

Sparse approximations

Current trends on kernel design

From GPLVM to Deep GPs

Multioutput GPs

④ Concluding Remarks

Concluding Remarks

- GPs are **Bayesian nonparametric** models, meaning that they:
 - i) naturally combine existing biases and the knowledge extracted from data;
 - ii) provide **probabilistic estimates**;
 - iii) the more data they see, the better they become.
- GPs can be applied to a plethora of settings: denoising, regression, interpolation, extrapolation, classification, optimization, unsupervised learning, hierarchical modelling, etc.
- **Hyperparameter optimization** can be performed via marginal likelihood, and jointly with sparse approximations, if they are required.
- The GP community is very active, come on over and join us!

Additional GP topics

Gaussian Process Summer Schools - GPSS

Superb presentations and workshops on several GP related subjects: <http://gpss.cc/>.

- More thorough discussion on sparse and scalable GPs
 - Bauer, van der Wilk and Rasmussen. “Understanding probabilistic sparse Gaussian process approximations”. NeurIPS, 2016.
 - Bui, Yan and Turner. “A Unifying Framework for Gaussian Process Pseudo-Point Approximations using Power Expectation Propagation”. JMLR, 2017.
 - Liu, *et al.* “When Gaussian process meets big data: A review of scalable GPs”. IEEE TNNLS, 2020.
 - Leibfried, *et al.* “A tutorial on sparse Gaussian processes and variational inference”. arXiv preprint, 2021.
- Bayesian optimisation
 - Frazier, “Bayesian Optimization”. INFORMS Tutorials, 2018.
 - Shahriari, Swersky, Wang, Adams and de Freitas, “Taking the Human Out of the Loop: A Review of Bayesian Optimization”. Proc. of the IEEE, 2016.
 - Snoek, Larochelle and Adams, “Practical Bayesian Optimization of Machine Learning Algorithms”. NeurIPS, 2012.

Additional GP topics

Gaussian Process Summer Schools - GPSS

Superb presentations and workshops on several GP related subjects: <http://gpss.cc/>.

- More thorough discussion on sparse and scalable GPs
 - Bauer, van der Wilk and Rasmussen. “**Understanding probabilistic sparse Gaussian process approximations**”. NeurIPS, 2016.
 - Bui, Yan and Turner. “**A Unifying Framework for Gaussian Process Pseudo-Point Approximations using Power Expectation Propagation**”. JMLR, 2017.
 - Liu, *et al.* “**When Gaussian process meets big data: A review of scalable GPs**”. IEEE TNNLS, 2020.
 - Leibfried, *et al.* “**A tutorial on sparse Gaussian processes and variational inference**”. arXiv preprint, 2021.
- Bayesian optimisation
 - Frazier, “**Bayesian Optimization**”. INFORMS Tutorials, 2018.
 - Shahriari, Swersky, Wang, Adams and de Freitas, “**Taking the Human Out of the Loop: A Review of Bayesian Optimization**”. Proc. of the IEEE, 2016.
 - Snoek, Larochelle and Adams, “**Practical Bayesian Optimization of Machine Learning Algorithms**”. NeurIPS, 2012.

Additional GP topics

Gaussian Process Summer Schools - GPSS

Superb presentations and workshops on several GP related subjects: <http://gpss.cc/>.

- More thorough discussion on sparse and scalable GPs
 - Bauer, van der Wilk and Rasmussen. “**Understanding probabilistic sparse Gaussian process approximations**”. NeurIPS, 2016.
 - Bui, Yan and Turner. “**A Unifying Framework for Gaussian Process Pseudo-Point Approximations using Power Expectation Propagation**”. JMLR, 2017.
 - Liu, *et al.* “**When Gaussian process meets big data: A review of scalable GPs**”. IEEE TNNLS, 2020.
 - Leibfried, *et al.* “**A tutorial on sparse Gaussian processes and variational inference**”. arXiv preprint, 2021.
- Bayesian optimisation
 - Frazier, “**Bayesian Optimization**”. INFORMS Tutorials, 2018.
 - Shahriari, Swersky, Wang, Adams and de Freitas, “**Taking the Human Out of the Loop: A Review of Bayesian Optimization**”. Proc. of the IEEE, 2016.
 - Snoek, Larochelle and Adams, “**Practical Bayesian Optimization of Machine Learning Algorithms**”. NeurIPS, 2012.

Additional GP topics

- Deep GPs and beyond
 - Salimbeni and Deisenroth. “**Doubly stochastic variational inference for deep Gaussian processes**”. NeurIPS, 2017.
 - Salimbeni, *et al.* “**Deep Gaussian processes with importance-weighted variational inference**”. ICML, 2019.
 - Bui, *et al.* “**Deep Gaussian processes for regression using approximate expectation propagation**”. ICML, 2016.
 - Dunlop, *et al.* “**How deep are deep Gaussian processes?**”. JMLR, 2018.
 - Aitchison, Yang and Ober. “**Deep kernel processes**”. ICML, 2021.
- GPs, deep learning and Neural Tangent Kernel (NTK)
 - Matthews, *et al.* “**Gaussian process behaviour in wide deep neural networks**”. ICLR, 2018.
 - Lee, *et al.* “**Deep neural networks as Gaussian processes**”. ICLR, 2018.
 - Jacot, Gabriel and Hongler. “**Neural tangent kernel: Convergence and generalization in neural networks**”. NeurIPS, 2018.
 - Dutordoir, *et al.* “**Deep Neural Networks as Point Estimates for Deep Gaussian Processes**”. arXiv preprint, 2021.

Additional GP topics

- Deep GPs and beyond
 - Salimbeni and Deisenroth. “**Doubly stochastic variational inference for deep Gaussian processes**”. NeurIPS, 2017.
 - Salimbeni, *et al.* “**Deep Gaussian processes with importance-weighted variational inference**”. ICML, 2019.
 - Bui, *et al.* “**Deep Gaussian processes for regression using approximate expectation propagation**”. ICML, 2016.
 - Dunlop, *et al.* “**How deep are deep Gaussian processes?**”. JMLR, 2018.
 - Aitchison, Yang and Ober. “**Deep kernel processes**”. ICML, 2021.
- GPs, deep learning and Neural Tangent Kernel (NTK)
 - Matthews, *et al.* “**Gaussian process behaviour in wide deep neural networks**”. ICLR, 2018.
 - Lee, *et al.* “**Deep neural networks as Gaussian processes**”. ICLR, 2018.
 - Jacot, Gabriel and Hongler. “**Neural tangent kernel: Convergence and generalization in neural networks**”. NeurIPS, 2018.
 - Dutordoir, *et al.* “**Deep Neural Networks as Point Estimates for Deep Gaussian Processes**”. arXiv preprint, 2021.

Questions?

César Lincoln C. Mattos
@cesarlincoln
cesarlincoln@dc.ufc.br

Felipe Tobar
@Felipe_Tobar
ftobar@dim.uchile.cl