

CS-E4075 Special course on Gaussian processes: Session #5 Latent variable models

Charles Gadd

Aalto University

charles.gadd@aalto.fi

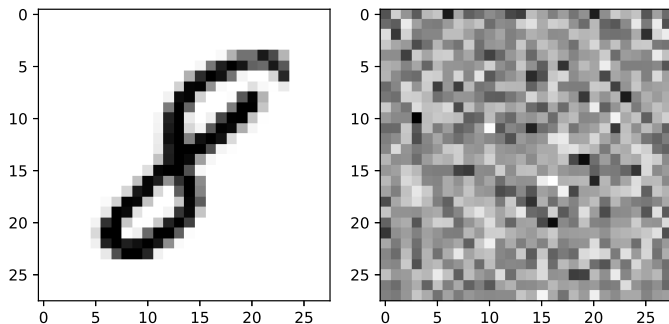
Monday 25.01.2021

Agenda for today

- Introduction
 - Why are LVMs useful?
 - Definition of LVMs
- Gaussian process latent variable models
 - Principal Component Analysis
 - Probabilistic PCA
 - Dual probabilistic PCA
 - GPLVM
- Multi-output models
 - Intrinsic Model of Coregionalisation
 - Semiparametric Latent Factor Model
 - Linear Model of Coregionalisation

Why are LVMs useful?

- Data has structure



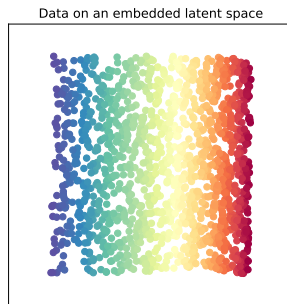
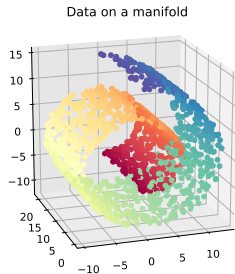
- ... the dimension of this space is large ($D = 784$)
- ... you would never sample this digit randomly

Why are LVMs useful?

- These samples lie on a **very** narrow manifold in $\mathbb{R}^{28 \times 28}$
- We should only require enough dimensions to describe the digit sufficiently
 - e.g. shape and distortions (rotation, translation, stretching)
- The number of these dimensions is called the *intrinsic dimensionality* and is often significantly smaller than the number of features.
- Its often far easier to perform your inference on this embedded manifold

Swiss roll example

Moving from \mathbb{R}^3 to \mathbb{R}^2



Some notation: Features $Y \in \mathbb{R}^{n \times D}$, latent variables $X \in \mathbb{R}^{n \times d}$

Definition of LVMs

Definition: Dimensionality reduction - Learning a projection onto a lower dimensional embedding.

Definition: Manifold learning - Learning this embedding **and** a *pre-image map* $g : x \rightarrow y$

A latent variable model is of the form:

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

Often such models make assumptions of

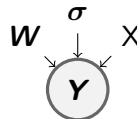
- Independence across latent samples
- Conditional independence across features **given** latent samples

The Gaussian process latent variable model (GPLVM)

Principal Component Analysis (Recap)

Formulating GPLVM

$$\mathbf{y}_i = \mathbf{W} \mathbf{x}_i + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$$

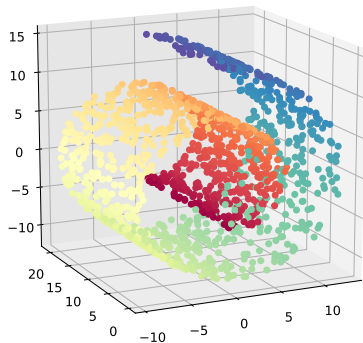


- Linear projection
- ... that projects to a new co-ordinate system
- ... such that the new basis is comprised of *principal components*
- ... which span the directions of greatest variance

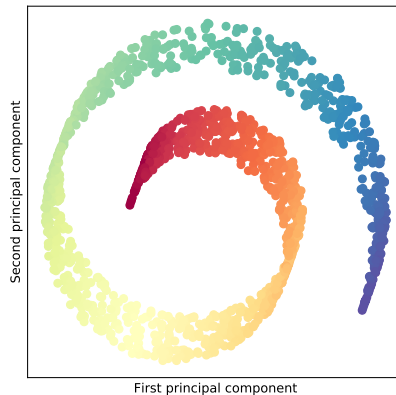
Only works well if data lies on a **plane** in high dimensional space

No representation of uncertainty

Principal Component Analysis

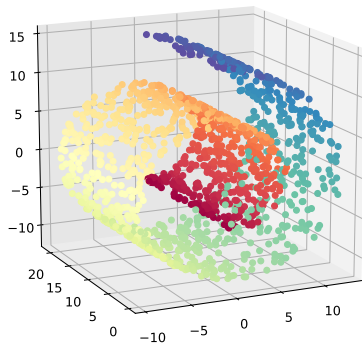


(a) Swiss roll data

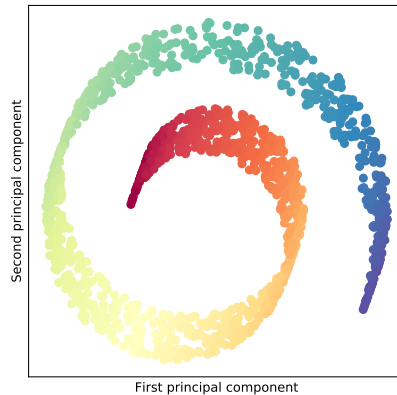


(b) PCA embedding

Principal Component Analysis



(a) Swiss roll data



(b) PCA embedding

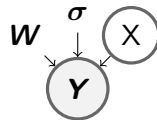
This is not optimal. This embedding does not capture all of the variance!

Probabilistic PCA

Formulating GPLVM

Likelihood

$$p(\mathbf{Y}|\mathbf{W}, \mathbf{X}, \sigma) \sim \prod_{i=1}^n \mathcal{N}(\mathbf{y}_i | \mathbf{W} \mathbf{x}_i, \sigma^2 \mathbf{I})$$

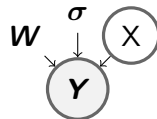


Probabilistic PCA

Formulating GPLVM

Likelihood

$$p(\mathbf{Y} | \mathbf{W}, \mathbf{X}, \sigma) \sim \prod_{i=1}^n \mathcal{N}(\mathbf{y}_i | \mathbf{W} \mathbf{x}_i, \sigma^2 \mathbf{I})$$



Place a conjugate Gaussian **prior** over latent space \mathcal{Z}

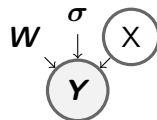
$$p(\mathbf{X}) \sim \prod_{i=1}^n \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I})$$

Probabilistic PCA

Formulating GPLVM

Likelihood

$$p(\mathbf{Y}|\mathbf{W}, \mathbf{X}, \sigma) \sim \prod_{i=1}^n \mathcal{N}(\mathbf{y}_i | \mathbf{W} \mathbf{x}_i, \sigma^2 \mathbf{I})$$



Place a conjugate Gaussian **prior** over latent space \mathcal{Z}

$$p(\mathbf{X}) \sim \prod_{i=1}^n \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I})$$

... and integrate over latent variables \mathbf{z} to obtain the **marginal likelihood**

$$p(\mathbf{Y}|\mathbf{W}, \sigma) = \prod_{i=1}^n \int p(\mathbf{y}_i | \mathbf{W}, \mathbf{x}_i, \sigma) p(\mathbf{x}_i) d\mathbf{x}$$

Probabilistic PCA

Formulating GPLVM

Marginal likelihood

$$p(\mathbf{Y}|\mathbf{W}, \sigma) = \prod_{i=1}^n \int p(\mathbf{y}_i|\mathbf{W}, \mathbf{x}_i, \sigma) p(\mathbf{x}_i) d\mathbf{x}$$

Using scaling and summation results

$$\alpha \mathcal{N}(\mu, \Sigma^2) = \mathcal{N}(\alpha\mu, \alpha^2\Sigma^2)$$

$$\sum_i \mathcal{N}(\mu_i, \Sigma_i^2) = \mathcal{N}\left(\sum_i \mu_i, \sum_i \Sigma_i^2\right)$$

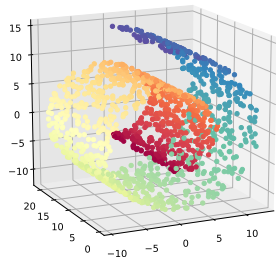
we can derive this

$$p(\mathbf{Y}|\mathbf{W}, \sigma) \sim \prod_{i=1}^n \int \mathcal{N}(\mathbf{y}_i|\mathbf{W}\mathbf{x}_i, \sigma^2\mathbf{I}) \mathcal{N}(\mathbf{x}_i|0, \mathbf{I}) d\mathbf{x}_i$$

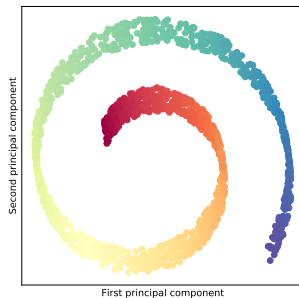
$$p(\mathbf{Y}|\mathbf{W}, \sigma) \sim \prod_{i=1}^n \mathcal{N}(\mathbf{y}_i|0, \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})$$

Probabilistic PCA - Swiss roll example

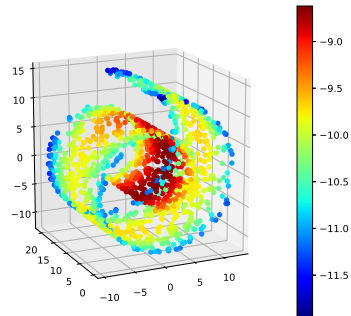
Formulating GPLVM



(a) Swiss roll data



(b) PCA embedding



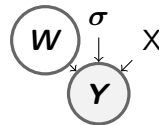
(c) Log-likelihood

Dual PPCA

Formulating GPLVM

Likelihood

$$p(\mathbf{Y}|\mathbf{W}, \mathbf{X}, \sigma) \sim \prod_{d=1}^D \mathcal{N}(\mathbf{y}_d | \mathbf{W} \mathbf{x}_d, \sigma^2 \mathbf{I})$$

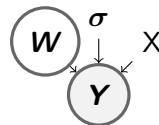


Dual PPCA

Formulating GPLVM

Likelihood

$$p(\mathbf{Y}|\mathbf{W}, \mathbf{X}, \sigma) \sim \prod_{d=1}^D \mathcal{N}(\mathbf{y}_d | \mathbf{W} \mathbf{x}_d, \sigma^2 \mathbf{I})$$



Place a conjugate Gaussian **prior** over the space of linear transformations

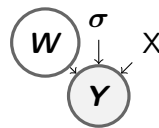
$$p(\mathbf{W}) \sim \prod_{i=1}^D \mathcal{N}(\mathbf{w}_i | 0, \mathbf{I})$$

Dual PPCA

Formulating GPLVM

Likelihood

$$p(\mathbf{Y}|\mathbf{W}, \mathbf{X}, \sigma) \sim \prod_{d=1}^D \mathcal{N}(\mathbf{y}_d | \mathbf{W} \mathbf{x}_d, \sigma^2 \mathbf{I})$$



Place a conjugate Gaussian **prior** over the space of linear transformations

$$p(\mathbf{W}) \sim \prod_{i=1}^D \mathcal{N}(\mathbf{w}_i | 0, \mathbf{I})$$

... and integrate over transformation matrix \mathbf{W} to obtain the **marginal likelihood**

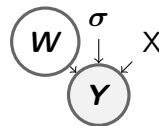
$$p(\mathbf{Y}|\mathbf{X}, \sigma) = \prod_{d=1}^D \int p(\mathbf{y}_d | \mathbf{W}, \mathbf{x}_d, \sigma) p(\mathbf{W}) d\mathbf{W}$$

Dual PPCA

Formulating GPLVM

Likelihood

$$p(\mathbf{Y}|\mathbf{W}, \mathbf{X}, \sigma) \sim \prod_{d=1}^D \mathcal{N}(\mathbf{y}_d | \mathbf{W} \mathbf{x}_d, \sigma^2 \mathbf{I})$$



Place a conjugate Gaussian **prior** over the space of linear transformations

$$p(\mathbf{W}) \sim \prod_{i=1}^D \mathcal{N}(\mathbf{w}_i | 0, \mathbf{I})$$

... and integrate over transformation matrix \mathbf{W} to obtain the **marginal likelihood**

$$p(\mathbf{Y} | \mathbf{X}, \sigma) = \prod_{i=1}^D \mathcal{N}(\mathbf{y}_i | 0, \mathbf{X} \mathbf{X}^T + \sigma^2 \mathbf{I})$$

The kernel

Formulating GPLVM

The dual PPCA marginal likelihood

$$p(\mathbf{Y}|\mathbf{X}, \sigma) = \prod_{i=1}^D \mathcal{N}(\mathbf{y}_i | 0, \mathbf{X}\mathbf{X}^T + \sigma^2 \mathbf{I})$$

The kernel

Formulating GPLVM

The dual PPCA marginal likelihood

$$p(\mathbf{Y}|\mathbf{X}, \sigma) = \prod_{i=1}^D \mathcal{N}(\mathbf{y}_i | 0, \mathbf{X}\mathbf{X}^T + \sigma^2 \mathbf{I})$$

The linear kernel

$$k(\mathbf{x}, \mathbf{x}') = \theta_b^2 + \theta_v^2 (\mathbf{x} - \mathbf{c})(\mathbf{x} - \mathbf{c})^T$$

The kernel

Formulating GPLVM

The dual PPCA marginal likelihood

$$p(\mathbf{Y} | \mathbf{X}, \sigma) = \prod_{i=1}^D \mathcal{N}(\mathbf{y}_i | 0, \mathbf{X} \mathbf{X}^T + \sigma^2 \mathbf{I})$$

The linear kernel

$$k(\mathbf{x}, \mathbf{x}') = \theta_b^2 + \theta_v^2 (\mathbf{x} - \mathbf{c})(\mathbf{x} - \mathbf{c})^T$$

The marginal likelihood for DPPCA is a product of D independent Gaussian processes with a linear kernel

The kernel

Formulating GPLVM

The dual PPCA marginal likelihood

$$p(\mathbf{Y}|\mathbf{X}, \sigma) = \prod_{i=1}^D \mathcal{N}(\mathbf{y}_i | 0, \mathbf{X}\mathbf{X}^T + \sigma^2 \mathbf{I})$$

The linear kernel

$$k(\mathbf{x}, \mathbf{x}') = \theta_b^2 + \theta_v^2 (\mathbf{x} - \mathbf{c})(\mathbf{x} - \mathbf{c})^T$$

The marginal likelihood for DPPCA is a product of D independent Gaussian processes with a linear kernel

We can change this kernel and obtain the GPLVM class of models

- DPPCA is a special case of GPLVM with a linear kernel
- Each dimension of the marginal can be interpreted as an independent GP
- Each dimension is *a priori* assumed independent, and identically distributed

Analytic solutions

- For PCA, PPCA and DPPCA, an analytic solution exists via solving an eigenvalue problem

Analytic solutions

- For PCA, PPCA and DPPCA, an analytic solution exists via solving an eigenvalue problem

Maximum likelihood (ML)

- Once we use a non-linear kernel analytical solutions often become intractable
- Instead we may resort to gradient based optimisation for (X, θ, σ) .¹

$$\hat{X}, \hat{\theta}, \hat{\sigma} = \arg \max_{X, \theta, \sigma} \{\log(\mathbf{Y} | X, \theta, \sigma)\} \propto \arg \max_{X, \theta, \sigma} \left\{ -\frac{D}{2} \log |\mathbf{K}| - \frac{1}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T) \right\}$$

¹where θ are the kernel hyper-parameters and, for example, $\mathbf{K} = \mathbf{X} \mathbf{X}^T + \sigma^2 \mathbf{I}$ in the linear kernel case

Maximum a-posteriori (MAP)

- Place a prior on latent variables \mathbf{Z}

$$\hat{\mathbf{X}}, \hat{\theta}, \hat{\sigma} = \arg \max_{\mathbf{X}, \theta, \sigma} \{ \log(\mathbf{Y} | \mathbf{X}, \theta, \sigma) + \log p(\mathbf{X}) \}$$

- We again use gradient based optimisation
- This prior acts to regularise the latent variables

These both optimise over a huge space, but don't do as poorly as you'd expect!

- Optimisation is *very* non-convex.
 - Multiple restarts, initialisation. How do we initialise?
- What is the latent dimensionality?
- Computational cost.

- We may also want to place a prior on \mathbf{Z} and integrate it out.²

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{i=1}^D \mathcal{N}(\mathbf{y}_i | \mathbf{0}, \mathbf{K})$$

$$p(\mathbf{Y}) = \int p(\mathbf{Y}|\mathbf{X}) p(\mathbf{X}) d\mathbf{X} \quad \text{and introduce} \quad p(\mathbf{X}) = \prod_{i=1}^n \mathcal{N}(\mathbf{x}_i | \mathbf{0}, \mathbf{I}_d)$$

- This is intractable as \mathbf{X} appears non-linearly in the inverse of the kernel
- Lets try and apply the standard *variational Bayes* approach

²Dropping dependence on θ and σ for clarity

Introduce a variational distribution

$$p(\mathbf{X}|\mathbf{Y}) \approx q(\mathbf{X}) = \prod_{i=1}^n \mathcal{N}(\mathbf{x}_m | \mu_n, S_n)$$

And compute the Jensen's lower bound

$$\log p(\mathbf{Y}) \geq \underbrace{\sum_{i=1}^D \int q(\mathbf{X}) \log p(\mathbf{y}_i | \mathbf{X}) d\mathbf{X}}_{\text{this remains intractable}} - \underbrace{\int q(\mathbf{X}) \log \frac{q(\mathbf{X})}{p(\mathbf{X})} d\mathbf{X}}_{KL(q||p)}$$

Introduce a variational distribution

$$p(\mathbf{X}|\mathbf{Y}) \approx q(\mathbf{X}) = \prod_{i=1}^n \mathcal{N}(\mathbf{x}_m | \mu_n, S_n)$$

And compute the Jensen's lower bound

$$\log p(\mathbf{Y}) \geq \underbrace{\sum_{i=1}^D \int q(\mathbf{X}) \log p(\mathbf{y}_i | \mathbf{X}) d\mathbf{X}}_{\text{this remains intractable}} - \underbrace{\int q(\mathbf{X}) \log \frac{q(\mathbf{X})}{p(\mathbf{X})} d\mathbf{X}}_{KL(q||p)}$$

Lets apply the variational sparse methodology of [Titsias(2009)] that we learnt last lecture! ³

³Slide 23 of session 4

Bayesian GPLVM (revisiting the variational sparse approach)

Lets expand our intractable integral term and augment inducing variables⁴

$$p(\mathbf{y}_i, \mathbf{f}_i, \mathbf{u}_i | \mathbf{X}, \mathbf{Z}) = p(\mathbf{y}_i | \mathbf{f}_i) p(\mathbf{f}_i | \mathbf{u}_i, \mathbf{X}, \mathbf{Z}) p(\mathbf{u}_i | \mathbf{Z})$$

⁴**Notation reminder:** inducing inputs \mathbf{Z} , inducing outputs \mathbf{U} .

Bayesian GPLVM (revisiting the variational sparse approach)

Lets expand our intractable integral term and augment inducing variables⁴

$$p(\mathbf{y}_i, \mathbf{f}_i, \mathbf{u}_i | \mathbf{X}, \mathbf{Z}) = p(\mathbf{y}_i | \mathbf{f}_i) p(\mathbf{f}_i | \mathbf{u}_i, \mathbf{X}, \mathbf{Z}) p(\mathbf{u}_i | \mathbf{Z})$$

Derive a variational approximation for the posterior

$$p(\mathbf{f}_i, \mathbf{u}_i | \mathbf{Y}, \mathbf{X}, \mathbf{Z}) \approx q(\mathbf{f}_i, \mathbf{u}_i) = p(\mathbf{f}_i | \mathbf{u}_i, \mathbf{X}, \mathbf{Z}) q(\mathbf{u}_i)$$

⁴**Notation reminder:** inducing inputs \mathbf{Z} , inducing outputs \mathbf{U} .

Bayesian GPLVM (revisiting the variational sparse approach)

Lets expand our intractable integral term and augment inducing variables⁴

$$p(\mathbf{y}_i, \mathbf{f}_i, \mathbf{u}_i | \mathbf{X}, \mathbf{Z}) = p(\mathbf{y}_i | \mathbf{f}_i) p(\mathbf{f}_i | \mathbf{u}_i, \mathbf{X}, \mathbf{Z}) p(\mathbf{u}_i | \mathbf{Z})$$

Derive a variational approximation for the posterior

$$p(\mathbf{f}_i, \mathbf{u}_i | \mathbf{Y}, \mathbf{X}, \mathbf{Z}) \approx q(\mathbf{f}_i, \mathbf{u}_i) = p(\mathbf{f}_i | \mathbf{u}_i, \mathbf{X}, \mathbf{Z}) q(\mathbf{u}_i)$$

⁴**Notation reminder:** inducing inputs \mathbf{Z} , inducing outputs \mathbf{U} .

Bayesian GPLVM (revisiting the variational sparse approach)

Lets expand our intractable integral term and augment inducing variables⁴

$$p(\mathbf{y}_i, \mathbf{f}_i, \mathbf{u}_i | \mathbf{X}, \mathbf{Z}) = p(\mathbf{y}_i | \mathbf{f}_i) p(\mathbf{f}_i | \mathbf{u}_i, \mathbf{X}, \mathbf{Z}) p(\mathbf{u}_i | \mathbf{Z})$$

Derive a variational approximation for the posterior

$$p(\mathbf{f}_i, \mathbf{u}_i | \mathbf{Y}, \mathbf{X}, \mathbf{Z}) \approx q(\mathbf{f}_i, \mathbf{u}_i) = p(\mathbf{f}_i | \mathbf{u}_i, \mathbf{X}, \mathbf{Z}) q(\mathbf{u}_i)$$

And we can use this to derive a new lower bound

$$\int q(\mathbf{X}) \log p(\mathbf{y}_i | \mathbf{X}) d\mathbf{X} \geq \int q(\mathbf{X}) q(\mathbf{f}_i, \mathbf{u}_i) \log \frac{p(\mathbf{y}_i, \mathbf{f}_i, \mathbf{u}_i | \mathbf{X}, \mathbf{Z})}{q(\mathbf{f}_i, \mathbf{u}_i)} d\mathbf{f}_i d\mathbf{u}_i d\mathbf{X}$$

⁴**Notation reminder:** inducing inputs \mathbf{Z} , inducing outputs \mathbf{U} .

Bayesian GPLVM (revisiting the variational sparse approach)

Lets expand likelihood term of our intractable integral and augment inducing variables⁴

$$p(\mathbf{y}_i, \mathbf{f}_i, \mathbf{u}_i | \mathbf{X}, \mathbf{Z}) = p(\mathbf{y}_i | \mathbf{f}_i) \cancel{p(\mathbf{f}_i | \mathbf{u}_i, \mathbf{X}, \mathbf{Z})} p(\mathbf{u}_i | \mathbf{Z})$$

Derive a variational approximation for the posterior

$$p(\mathbf{f}_i, \mathbf{u}_i | \mathbf{Y}, \mathbf{X}, \mathbf{Z}) \approx \mathbf{q}(\mathbf{f}_i, \mathbf{u}_i) = \cancel{p(\mathbf{f}_i | \mathbf{u}_i, \mathbf{X}, \mathbf{Z})} q(\mathbf{u}_i)$$

And we can use this to derive a new lower bound

$$\int q(\mathbf{X}) \log p(\mathbf{y}_i | \mathbf{X}) d\mathbf{X} \geq \int q(\mathbf{X}) q(\mathbf{f}_i, \mathbf{u}_i) \log \frac{p(\mathbf{y}_i, \mathbf{f}_i, \mathbf{u}_i | \mathbf{X}, \mathbf{Z})}{\mathbf{q}(\mathbf{f}_i, \mathbf{u}_i)} d\mathbf{f}_i d\mathbf{u}_i d\mathbf{X}$$

⁴**Notation reminder:** inducing inputs \mathbf{Z} , inducing outputs \mathbf{U} .

And we can use this to derive a new lower bound

$$\int q(\mathbf{X}) \log p(\mathbf{y}_i | \mathbf{X}) d\mathbf{X} \geq \int q(\mathbf{X}) q(\mathbf{f}_i, \mathbf{u}_i) \log \frac{p(\mathbf{y}_i | \mathbf{f}_i) p(\mathbf{u}_i | \mathbf{Z})}{q(\mathbf{u}_i)} d\mathbf{f}_i d\mathbf{u}_i d\mathbf{X}$$

And we can use this to derive a new lower bound

$$\int q(\mathbf{X}) \log p(\mathbf{y}_i | \mathbf{X}) d\mathbf{X} \geq \int q(\mathbf{X}) q(\mathbf{f}_i, \mathbf{u}_i) \log \frac{p(\mathbf{y}_i | \mathbf{f}_i) p(\mathbf{u}_i | \mathbf{Z})}{q(\mathbf{u}_i)} d\mathbf{f}_i d\mathbf{u}_i d\mathbf{X}$$

We can now split this into an integral and a tractable KL term

$$\begin{aligned} \int q(\mathbf{X}) \log p(\mathbf{y}_i | \mathbf{X}) d\mathbf{X} &\geq \int q(\mathbf{X}) q(\mathbf{f}_i, \mathbf{u}_i) \log p(\mathbf{y}_i | \mathbf{f}_i) d\mathbf{f}_i d\mathbf{u}_i d\mathbf{X} \\ &\quad - \text{KL}(q(\mathbf{u}_i) || p(\mathbf{u}_i | \mathbf{Z})) \end{aligned}$$

And we can use this to derive a new lower bound

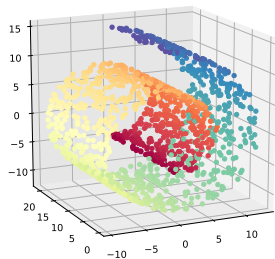
$$\int q(\mathbf{X}) \log p(\mathbf{y}_i | \mathbf{X}) d\mathbf{X} \geq \int q(\mathbf{X}) q(\mathbf{f}_i, u_i) \log \frac{p(\mathbf{y}_i | \mathbf{f}_i) p(u_i | Z)}{q(u_i)} d\mathbf{f}_i du_i d\mathbf{X}$$

We can now split this into an integral and a tractable KL term

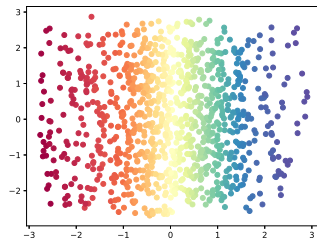
$$\begin{aligned} \int q(\mathbf{X}) \log p(\mathbf{y}_i | \mathbf{X}) d\mathbf{X} &\geq \int q(\mathbf{X}) q(\mathbf{f}_i, u_i) \log p(\mathbf{y}_i | \mathbf{f}_i) d\mathbf{f}_i du_i d\mathbf{X} \\ &\quad - \text{KL}(q(u_i) || p(u_i | Z)) \end{aligned}$$

- Swap order of integration
- This integral is now tractable (for some kernels) as \mathbf{X} no longer needs to be pushed through the kernel
- For more details see [Damianou(2015)] or [Titsias and Lawrence(2010)]

Bayesian GPLVM - example



(a) Swiss roll data



(b) BGPLVM embedded mean

The many flavours of GPLVM

- Shared GPLVM - map from a shared latent space to separate observation spaces
- Back constrained GPLVM - preserving locality in the image map
- Dynamic GPLVM (or GP dynamical model) - add a dynamic prior for supervised learning
- 'Deep' GPs - add a GP prior onto Z .

... and many more!

The applications of GPLVM

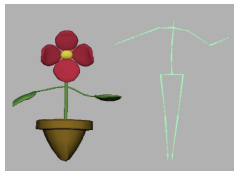


Figure: Shared GPLVM: Disney research ([link](#))

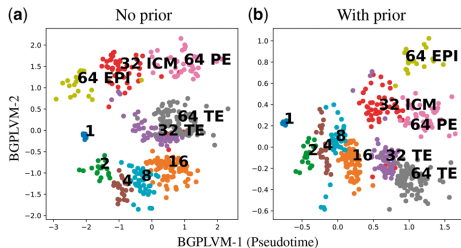


Figure: BGPLVM for single cell data

Models

- [Lawrence(2005)] - PCA \rightarrow PPCA \rightarrow DPPCA \rightarrow GPLVM derivation details
- [Titsias and Lawrence(2010)] - Bayesian GPLVM paper
- [Damianou(2015)] - Bayesian GPLVM thesis
- [Lawrence and Quiñonero-Candela(2006)] - Back constrained GPLVM
- [Ek and Lawrence(2009)] - Shared GPLVM

Applications

- [Yamane et al.(2010)] Yamane, Ariki, and Hodgins] - Disney research Shared GPLVM for animation
- [Ahmed et al.(2019)] Ahmed, Rattray, and Boukouvalas] - BGPLVM for single cell data

Multi-output Gaussian processes

Pre-requisites: The Kronecker product

Take two matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{p \times q}$

$$A = \begin{bmatrix} a_{1,1} & \dots & a_{1,n} \\ \vdots & & \vdots \\ a_{m,1} & \dots & a_{m,n} \end{bmatrix}$$

$$A \otimes B = \begin{bmatrix} a_{1,1}B & \dots & a_{1,n}B \\ \vdots & & \vdots \\ a_{m,1}B & \dots & a_{m,n}B \end{bmatrix}$$

Pre-requisites: The Kronecker product

Take two matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{p \times q}$

$$A = \begin{bmatrix} a_{1,1} & \dots & a_{1,n} \\ \vdots & & \vdots \\ a_{m,1} & \dots & a_{m,n} \end{bmatrix}$$

$$A \otimes B = \begin{bmatrix} a_{1,1}B & \dots & a_{1,n}B \\ \vdots & & \vdots \\ a_{m,1}B & \dots & a_{m,n}B \end{bmatrix}$$

So what is the dimension of $A \otimes B$?

Pre-requisites: The Kronecker product

Take two matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{p \times q}$

$$A = \begin{bmatrix} a_{1,1} & \dots & a_{1,n} \\ \vdots & & \vdots \\ a_{m,1} & \dots & a_{m,n} \end{bmatrix}$$

$$A \otimes B = \begin{bmatrix} a_{1,1}B & \dots & a_{1,n}B \\ \vdots & & \vdots \\ a_{m,1}B & \dots & a_{m,n}B \end{bmatrix}$$

So what is the dimension of $A \otimes B$?

$$A \otimes B \in \mathbb{R}^{mp \times nq}$$

Pre-requisites: The Kronecker product

Take two matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{p \times q}$

$$A = \begin{bmatrix} a_{1,1} & \dots & a_{1,n} \\ \vdots & & \vdots \\ a_{m,1} & \dots & a_{m,n} \end{bmatrix} \quad A \otimes B = \begin{bmatrix} a_{1,1}B & \dots & a_{1,n}B \\ \vdots & & \vdots \\ a_{m,1}B & \dots & a_{m,n}B \end{bmatrix}$$

So what is the dimension of $A \otimes B$?

$$A \otimes B \in \mathbb{R}^{mp \times nq}$$

The inversion rule:

$$(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$$

And is invertible *if and only if* both A and B are invertible.

Motivation: Multiple processes

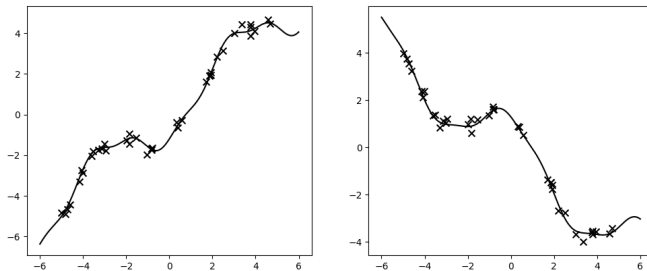


Figure: Two linearly correlated processes

Motivation: Multiple processes

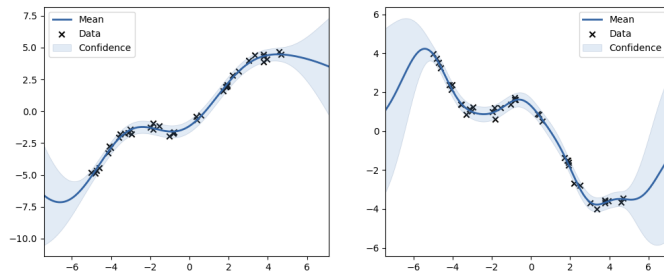


Figure: Two independent Gaussian process fits

$$f_1(\mathbf{x}) \sim \mathcal{GP}(0, k_1(\mathbf{x}, \mathbf{x}'))$$

$$f_1 \sim \mathcal{N}(0, K_1)$$

$$f_2(\mathbf{x}) \sim \mathcal{GP}(0, k_2(\mathbf{x}, \mathbf{x}'))$$

$$f_2 \sim \mathcal{N}(0, K_2)$$

Motivation: Multiple processes

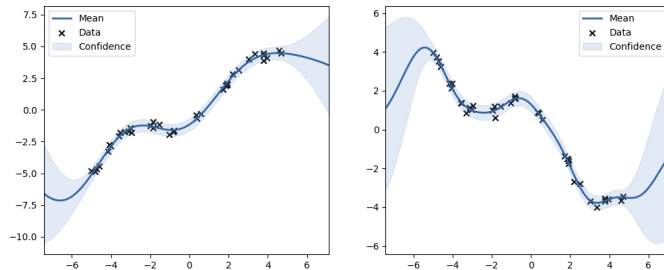


Figure: Two independent Gaussian process fits

$$\begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K_1 & 0 \\ 0 & K_2 \end{bmatrix} \right)$$

Motivation: Multiple processes

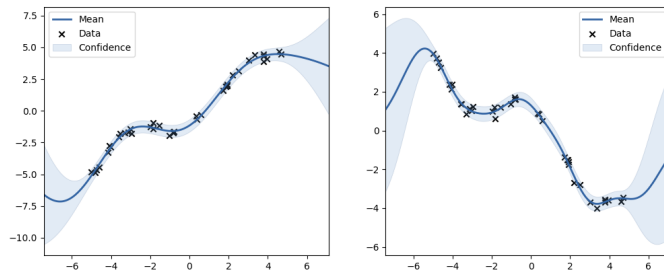


Figure: Two independent Gaussian process fits

$$\begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K_1 & ? \\ ? & K_2 \end{bmatrix} \right)$$

Intrinsic Model of Coregionalisation (IMC)

General case

Sample S functions i.i.d. from the shared underlying process $\mathbf{u}^{(s)} \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$.

Intrinsic Model of Coregionalisation (IMC)

General case

Sample S functions i.i.d. from the shared underlying process $\mathbf{u}^{(s)} \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$.

The vector valued process is then defined as a weighted sum,

$$\mathbf{f}(\mathbf{x}) = \sum_{s=1}^S \mathbf{a}^{(s)} \mathbf{u}^{(s)}(\mathbf{x}).$$

where

$$\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_D(\mathbf{x})]^T, \quad \mathbf{a}^{(s)} = [a_1^{(s)}, a_2^{(s)}, \dots, a_D^{(s)}]^T,$$

Intrinsic Model of Coregionalisation (IMC)

General case

Sample S functions i.i.d. from the shared underlying process $u^{(s)} \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$.

The vector valued process is then defined as a weighted sum,

$$\mathbf{f}(\mathbf{x}) = \sum_{s=1}^S \mathbf{a}^{(s)} u^{(s)}(\mathbf{x}).$$

where

$$\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_D(\mathbf{x})]^T, \quad \mathbf{a}^{(s)} = [a_1^{(s)}, a_2^{(s)}, \dots, a_D^{(s)}]^T,$$

$$\begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K_1 & ? \\ ? & K_2 \end{bmatrix}\right)$$

Intrinsic Model of Coregionalisation (IMC)

General case

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \mathbb{E}[\mathbf{f}(\mathbf{x}) \mathbf{f}(\mathbf{x}')^T] - \mathbb{E}[\mathbf{f}(\mathbf{x})] \mathbb{E}[\mathbf{f}(\mathbf{x}')]^T$$

$$\mathbf{f}(\mathbf{x}) = \mathbf{a}^{(1)} u^{(1)}(\mathbf{x}) + \mathbf{a}^{(2)} u^{(2)}(\mathbf{x}) + \dots$$

$$\begin{aligned} \text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) &= \mathbf{a}^{(1)} \mathbf{a}^{(1)T} \text{cov}(u^{(1)}(\mathbf{x}), u^{(1)}(\mathbf{x}')) + \mathbf{a}^{(2)} \mathbf{a}^{(2)T} \text{cov}(u^{(2)}(\mathbf{x}), u^{(2)}(\mathbf{x}')) + \dots \\ &= \mathbf{a}^{(1)} \mathbf{a}^{(1)T} k(\mathbf{x}, \mathbf{x}') + \mathbf{a}^{(2)} \mathbf{a}^{(2)T} k(\mathbf{x}, \mathbf{x}') + \dots \\ &= \left[\mathbf{a}^{(1)} \mathbf{a}^{(1)T} + \mathbf{a}^{(2)} \mathbf{a}^{(2)T} + \dots \right] k(\mathbf{x}, \mathbf{x}') \end{aligned}$$

Intrinsic Model of Coregionalisation (IMC)

General case

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \mathbb{E}[\mathbf{f}(\mathbf{x}) \mathbf{f}(\mathbf{x}')^T] - \mathbb{E}[\mathbf{f}(\mathbf{x})] \mathbb{E}[\mathbf{f}(\mathbf{x}')]^T$$

$$\mathbf{f}(\mathbf{x}) = \mathbf{a}^{(1)} u^{(1)}(\mathbf{x}) + \mathbf{a}^{(2)} u^{(2)}(\mathbf{x}) + \dots$$

$$\begin{aligned} \text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) &= \mathbf{a}^{(1)} \mathbf{a}^{(1)T} \text{cov}(u^{(1)}(\mathbf{x}), u^{(1)}(\mathbf{x}')) + \mathbf{a}^{(2)} \mathbf{a}^{(2)T} \text{cov}(u^{(2)}(\mathbf{x}), u^{(2)}(\mathbf{x}')) + \dots \\ &= \mathbf{a}^{(1)} \mathbf{a}^{(1)T} k(\mathbf{x}, \mathbf{x}') + \mathbf{a}^{(2)} \mathbf{a}^{(2)T} k(\mathbf{x}, \mathbf{x}') + \dots \\ &= \underbrace{\left[\mathbf{a}^{(1)} \mathbf{a}^{(1)T} + \mathbf{a}^{(2)} \mathbf{a}^{(2)T} + \dots \right]}_{\hat{\mathbf{B}} \in \mathbb{R}^{D \times D}} k(\mathbf{x}, \mathbf{x}') \end{aligned}$$

Intrinsic Model of Coregionalisation (IMC)

General case

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \mathbb{E}[\mathbf{f}(\mathbf{x}) \mathbf{f}(\mathbf{x}')^T] - \mathbb{E}[\mathbf{f}(\mathbf{x})] \mathbb{E}[\mathbf{f}(\mathbf{x}')]^T$$

$$\mathbf{f}(\mathbf{x}) = \mathbf{a}^{(1)} u^{(1)}(\mathbf{x}) + \mathbf{a}^{(2)} u^{(2)}(\mathbf{x}) + \dots$$

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \mathbf{a}^{(1)} \mathbf{a}^{(1)T} \text{cov}(u^{(1)}(\mathbf{x}), u^{(1)}(\mathbf{x}')) + \mathbf{a}^{(2)} \mathbf{a}^{(2)T} \text{cov}(u^{(2)}(\mathbf{x}), u^{(2)}(\mathbf{x}')) + \dots$$

$$= \mathbf{a}^{(1)} \mathbf{a}^{(1)T} k(\mathbf{x}, \mathbf{x}') + \mathbf{a}^{(2)} \mathbf{a}^{(2)T} k(\mathbf{x}, \mathbf{x}') + \dots$$

$$= \underbrace{\left[\mathbf{a}^{(1)} \mathbf{a}^{(1)T} + \mathbf{a}^{(2)} \mathbf{a}^{(2)T} + \dots \right]}_{\hat{\mathbf{B}} \in \mathbb{R}^{D \times D}} k(\mathbf{x}, \mathbf{x}')$$

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \hat{\mathbf{B}} k(\mathbf{x}, \mathbf{x}')$$

Intrinsic Model of Coregionalisation (IMC)

General case

$$\hat{\mathbf{B}} = \begin{bmatrix} b_{1,1} & b_{1,2} \\ b_{2,1} & b_{2,2} \end{bmatrix}, \quad \text{if } D = 2$$

$$\begin{aligned} \begin{bmatrix} f_1(x) \\ f_2(x) \end{bmatrix} &\sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} b_{1,1}K & b_{1,2}K \\ b_{2,1}K & b_{2,2}K \end{bmatrix} \right) \\ &\sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \hat{\mathbf{B}} \otimes K \right) \end{aligned}$$

Semiparametric latent factor model (SLFM)

General case

Sample Q functions from separate processes $u_q \sim \mathcal{GP}(0, k_q(\mathbf{x}, \mathbf{x}'))$.

Semiparametric latent factor model (SLFM)

General case

Sample Q functions from separate processes $u_q \sim \mathcal{GP}(0, k_q(\mathbf{x}, \mathbf{x}'))$.

The vector valued process is then defined as a weighted sum,

$$\mathbf{f}(\mathbf{x}) = \sum_{q=1}^Q \mathbf{a}_q u_q(\mathbf{x}).$$

where

$$\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_D(\mathbf{x})]^T, \quad \mathbf{a}_q = [a_{q,1}, a_{q,2}, \dots, a_{q,D}]^T,$$

Semiparametric latent factor model (SLFM)

General case

Sample Q functions from separate processes $u_q \sim \mathcal{GP}(0, k_q(\mathbf{x}, \mathbf{x}'))$.

The vector valued process is then defined as a weighted sum,

$$\mathbf{f}(\mathbf{x}) = \sum_{q=1}^Q \mathbf{a}_q u_q(\mathbf{x}).$$

where

$$\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_D(\mathbf{x})]^T, \quad \mathbf{a}_q = [a_{q,1}, a_{q,2}, \dots, a_{q,D}]^T,$$

$$\begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K_1 & ? \\ ? & K_2 \end{bmatrix}\right)$$

Semiparametric latent factor model (SLFM)

General case

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \mathbb{E}[\mathbf{f}(\mathbf{x}) \mathbf{f}(\mathbf{x}')^T] - \mathbb{E}[\mathbf{f}(\mathbf{x})] \mathbb{E}[\mathbf{f}(\mathbf{x}')]^T$$

$$\mathbf{f}(\mathbf{x}) = \mathbf{a}_1 u_1(\mathbf{x}) + \mathbf{a}_2 u_2(\mathbf{x}) + \dots$$

$$\begin{aligned} \text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) &= \mathbf{a}_1 \mathbf{a}_1^T \text{cov}(u_1(\mathbf{x}), u_1(\mathbf{x}')) + \mathbf{a}_2 \mathbf{a}_2^T \text{cov}(u_2(\mathbf{x}), u_2(\mathbf{x}')) + \dots \\ &= \underbrace{\mathbf{a}_1 \mathbf{a}_1^T}_{\tilde{\mathbf{B}}_1 \in \mathbb{R}^{D \times D}} k_1(\mathbf{x}, \mathbf{x}') + \underbrace{\mathbf{a}_2 \mathbf{a}_2^T}_{\tilde{\mathbf{B}}_2 \in \mathbb{R}^{D \times D}} k_2(\mathbf{x}, \mathbf{x}') + \dots \end{aligned}$$

Semiparametric latent factor model (SLFM)

General case

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \mathbb{E}[\mathbf{f}(\mathbf{x}) \mathbf{f}(\mathbf{x}')^T] - \mathbb{E}[\mathbf{f}(\mathbf{x})] \mathbb{E}[\mathbf{f}(\mathbf{x}')]^T$$

$$\mathbf{f}(\mathbf{x}) = \mathbf{a}_1 u_1(\mathbf{x}) + \mathbf{a}_2 u_2(\mathbf{x}) + \dots$$

$$\begin{aligned} \text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) &= \mathbf{a}_1 \mathbf{a}_1^T \text{cov}(u_1(\mathbf{x}), u_1(\mathbf{x}')) + \mathbf{a}_2 \mathbf{a}_2^T \text{cov}(u_2(\mathbf{x}), u_2(\mathbf{x}')) + \dots \\ &= \underbrace{\mathbf{a}_1 \mathbf{a}_1^T}_{\tilde{\mathbf{B}}_1 \in \mathbb{R}^{D \times D}} k_1(\mathbf{x}, \mathbf{x}') + \underbrace{\mathbf{a}_2 \mathbf{a}_2^T}_{\tilde{\mathbf{B}}_2 \in \mathbb{R}^{D \times D}} k_2(\mathbf{x}, \mathbf{x}') + \dots \end{aligned}$$

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \tilde{\mathbf{B}}_1 k_1(\mathbf{x}, \mathbf{x}') + \tilde{\mathbf{B}}_2 k_2(\mathbf{x}, \mathbf{x}') + \dots$$

$$\begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \sum_{q=1}^Q \tilde{\mathbf{B}}_q \otimes K_q \right), \quad \text{with } D = 2$$

Linear Model of Coregionalisation (LMC)

General case

Sample S_q functions from Q separate processes $u_q^{(s)} \sim \mathcal{GP}(0, k_q(\mathbf{x}, \mathbf{x}'))$.

Linear Model of Coregionalisation (LMC)

General case

Sample S_q functions from Q separate processes $u_q^{(s)} \sim \mathcal{GP}(0, k_q(\mathbf{x}, \mathbf{x}'))$.

The vector valued process is then defined as a weighted sum,

$$\mathbf{f}(\mathbf{x}) = \sum_{q=1}^Q \sum_{s=1}^S \mathbf{a}_q^{(s)} u_q^{(s)}(\mathbf{x}).$$

where

$$\mathbf{f}(\mathbf{x}) = [f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_D(\mathbf{x})]^T, \quad \mathbf{a}_q^{(s)} = [a_{q,1}^{(s)}, a_{q,2}^{(s)}, \dots, a_{q,D}^{(s)}]^T$$

Linear Model of Coregionalisation (LMC)

General case

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \mathbb{E}[\mathbf{f}(\mathbf{x}) \mathbf{f}(\mathbf{x}')^T] - \mathbb{E}[\mathbf{f}(\mathbf{x})] \mathbb{E}[\mathbf{f}(\mathbf{x}')]^T$$

Intrinsic Model of Coregionalization

$$\mathbf{f}(\mathbf{x}) = \left[\mathbf{a}_1^{(1)} u_1^{(1)}(\mathbf{x}) + \mathbf{a}_1^{(2)} u_1^{(2)}(\mathbf{x}) + \dots \right] \\ + \left[\mathbf{a}_2^{(1)} u_2^{(1)}(\mathbf{x}) + \mathbf{a}_2^{(2)} u_2^{(2)}(\mathbf{x}) + \dots \right] + \dots$$

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \underbrace{\left[\mathbf{a}_1^{(1)} \mathbf{a}_1^{(1)T} + \mathbf{a}_1^{(2)} \mathbf{a}_1^{(2)T} \right]}_{\mathbf{B}_1 \in \mathbb{R}^{D \times D}} k_1(\mathbf{x}, \mathbf{x}') \\ + \underbrace{\left[\mathbf{a}_2^{(1)} \mathbf{a}_2^{(1)T} + \mathbf{a}_2^{(2)} \mathbf{a}_2^{(2)T} \right]}_{\mathbf{B}_2 \in \mathbb{R}^{D \times D}} k_2(\mathbf{x}, \mathbf{x}') + \dots$$

Linear Model of Coregionalisation (LMC)

General case

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \mathbb{E}[\mathbf{f}(\mathbf{x}) \mathbf{f}(\mathbf{x}')^T] - \mathbb{E}[\mathbf{f}(\mathbf{x})] \mathbb{E}[\mathbf{f}(\mathbf{x}')]^T$$

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \underbrace{\left[\mathbf{a}_1^{(1)} \mathbf{a}_1^{(1)T} + \mathbf{a}_1^{(2)} \mathbf{a}_1^{(2)T} \right]}_{\mathbf{B}_1 \in \mathbb{R}^{D \times D}} k_1(\mathbf{x}, \mathbf{x}') + \underbrace{\left[\mathbf{a}_2^{(1)} \mathbf{a}_2^{(1)T} + \mathbf{a}_2^{(2)} \mathbf{a}_2^{(2)T} \right]}_{\mathbf{B}_2 \in \mathbb{R}^{D \times D}} k_2(\mathbf{x}, \mathbf{x}') + \dots$$

$$\text{cov}(\mathbf{f}(\mathbf{x}), \mathbf{f}(\mathbf{x}')) = \mathbf{B}_1 k_1(\mathbf{x}, \mathbf{x}') + \mathbf{B}_2 k_2(\mathbf{x}, \mathbf{x}') + \dots$$

$$\begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \sum_{q=1}^Q \mathbf{B}_q \otimes K_q \right), \quad \text{with } D = 2$$

Linear Model of Coregionalisation (LMC)

Example

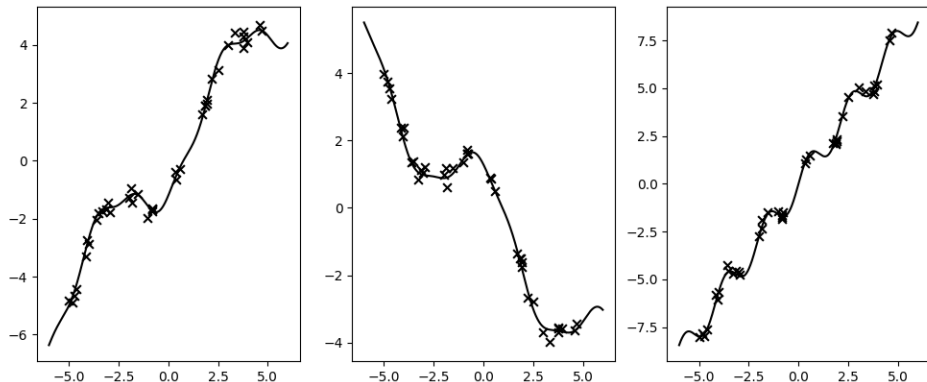


Figure: A third process

Linear Model of Coregionalisation (LMC)

Example

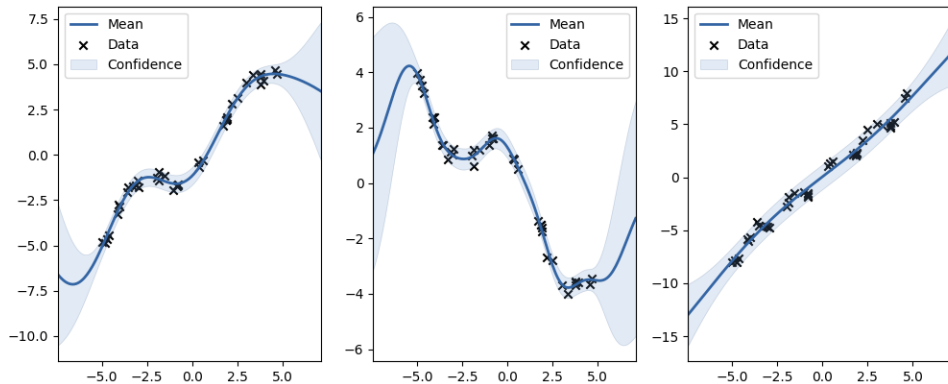


Figure: Independent Gaussian process fits

Linear Model of Coregionalisation (LMC)

Example

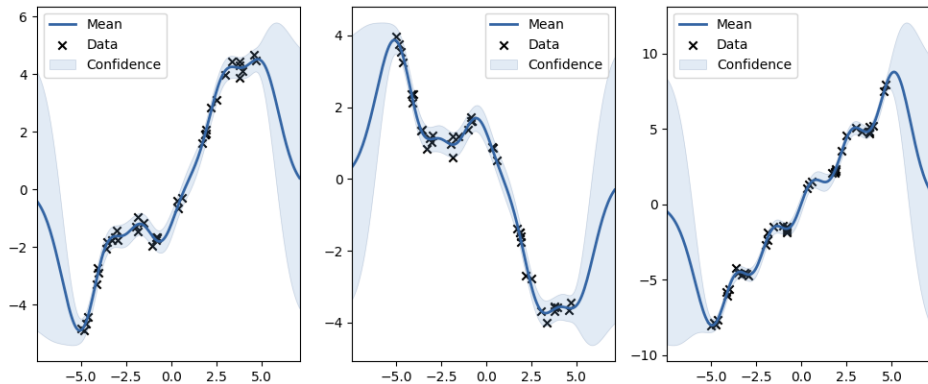


Figure: Linear Model of Coregionalisation fit

Next lecture on kernel learning

- Matérn kernel
- Multiple kernel learning
- Non-stationary RBF kernel
- Spectral kernels



Sumon Ahmed, Magnus Rattray, and Alexis Boukouvalas.

Grandprix: scaling up the bayesian gplvm for single-cell data.

Bioinformatics, 35(1):47–54, 2019.



Andreas Damianou.

Deep Gaussian processes and variational propagation of uncertainty.

PhD thesis, University of Sheffield, 2015.



Carl Henrik Ek and PHTND Lawrence.

Shared Gaussian process latent variable models.

PhD thesis, Citeseer, 2009.



Neil Lawrence.

Probabilistic non-linear principal component analysis with gaussian process latent variable models.

Journal of machine learning research, 6(Nov):1783–1816, 2005.

Bibliography II



Neil D Lawrence and Joaquin Quiñonero-Candela.

Local distance preservation in the gp-lvm through back constraints.

In *Proceedings of the 23rd international conference on Machine learning*, pages 513–520, 2006.



Michalis Titsias.

Variational learning of inducing variables in sparse gaussian processes.

In *Artificial Intelligence and Statistics*, pages 567–574, 2009.



Michalis Titsias and Neil D Lawrence.

Bayesian gaussian process latent variable model.

In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 844–851. JMLR Workshop and Conference Proceedings, 2010.



Katsu Yamane, Yuka Ariki, and Jessica Hodgins.

Animating non-humanoid characters with human motion data.

In *Proceedings of the 2010 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pages 169–178, 2010.