

Human Motion Modelling with Gaussian Processes

Neil D. Lawrence

7th February 2008

- 1 Probabilistic Dimensionality Reduction
- 2 Examples
- 3 Model Extensions
- 4 Conclusions

1 Probabilistic Dimensionality Reduction

2 Examples

3 Model Extensions

4 Conclusions

A Probabilistic Non-linear PCA

- PCA has a probabilistic interpretation [Tipping and Bishop, 1999].
- It is difficult to 'non-linearise'.

Dual Probabilistic PCA

- We present a new probabilistic interpretation of PCA [Lawrence, 2005].
- This interpretation can be made non-linear.
- The result is non-linear probabilistic PCA.

q — dimension of latent/embedded space

d — dimension of data space

n — number of data points

centred data, $\mathbf{Y} = [\mathbf{y}_{1,:}, \dots, \mathbf{y}_{n,:}]^T = [\mathbf{y}_{:,1}, \dots, \mathbf{y}_{:,d}] \in \mathbb{R}^{n \times d}$

latent variables, $\mathbf{X} = [\mathbf{x}_{1,:}, \dots, \mathbf{x}_{n,:}]^T = [\mathbf{x}_{:,1}, \dots, \mathbf{x}_{:,q}] \in \mathbb{R}^{n \times q}$

mapping matrix, $\mathbf{W} \in \mathbb{R}^{d \times q}$

$\mathbf{a}_{i,:}$ is a vector from the i th row of a given matrix \mathbf{A}

$\mathbf{a}_{:,j}$ is a vector from the j th row of a given matrix \mathbf{A}

\mathbf{X} and \mathbf{Y} are design matrices

Covariance given by $n^{-1}\mathbf{Y}^T\mathbf{Y}$.

Inner product matrix given by $\mathbf{Y}\mathbf{Y}^T$.

Linear Latent Variable Model

- Represent data, \mathbf{Y} , with a lower dimensional set of latent variables \mathbf{X} .

Assume a linear relationship of the form

$$\mathbf{y}_{i,:} = \mathbf{W}\mathbf{x}_{i,:} + \boldsymbol{\eta}_{i,:},$$

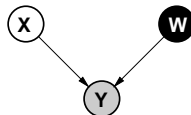
where

$$\boldsymbol{\eta}_{i,:} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}).$$

Probabilistic PCA

- Define *linear-Gaussian relationship* between latent variables and data.
- **Standard** Latent variable approach:

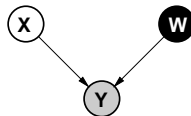
- ▶ Define Gaussian prior over *latent space*, \mathbf{X} .
- ▶ Integrate out *latent variables*.



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

Probabilistic PCA

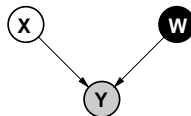
- Define *linear-Gaussian relationship* between latent variables and data.
- **Standard** Latent variable approach:
 - ▶ Define Gaussian prior over *latent space*, \mathbf{X} .
 - ▶ Integrate out *latent variables*.



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

Probabilistic PCA

- Define *linear-Gaussian relationship* between latent variables and data.
- **Standard** Latent variable approach:
 - ▶ Define Gaussian prior over *latent space*, \mathbf{X} .
 - ▶ Integrate out *latent variables*.

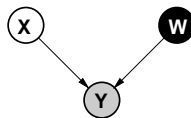


$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^n N(\mathbf{y}_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

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

Probabilistic PCA

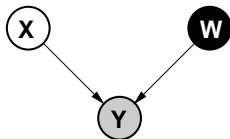
- Define *linear-Gaussian relationship* between latent variables and data.
- **Standard** Latent variable approach:
 - ▶ Define Gaussian prior over *latent space*, \mathbf{X} .
 - ▶ Integrate out *latent variables*.



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

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

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



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

$$p(\mathbf{Y}|\mathbf{W}) = \prod_{i=1}^n N(\mathbf{y}_{i,:} | \mathbf{0}, \mathbf{C}), \quad \mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$$

$$\log p(\mathbf{Y}|\mathbf{W}) = -\frac{n}{2} \log |\mathbf{C}| - \frac{1}{2} \text{tr}(\mathbf{C}^{-1} \mathbf{Y}^T \mathbf{Y}) + \text{const.}$$

If \mathbf{U}_q are first q principal eigenvectors of $n^{-1} \mathbf{Y}^T \mathbf{Y}$ and the corresponding eigenvalues are Λ_q ,

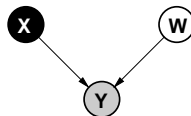
$$\mathbf{W} = \mathbf{U}_q \mathbf{L} \mathbf{V}^T, \quad \mathbf{L} = (\Lambda_q - \sigma^2 \mathbf{I})^{\frac{1}{2}}$$

where \mathbf{V} is an arbitrary rotation matrix.

Dual Probabilistic PCA

- Define *linear-Gaussian relationship* between latent variables and data.
- **Novel** Latent variable approach:

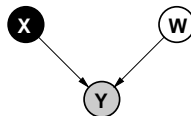
- ▶ Define Gaussian prior over *parameters*, \mathbf{W} .
- ▶ Integrate out *parameters*.



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

Dual Probabilistic PCA

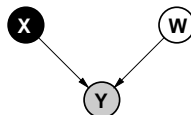
- Define *linear-Gaussian relationship* between latent variables and data.
- **Novel** Latent variable approach:
 - ▶ Define Gaussian prior over *parameters*, \mathbf{W} .
 - ▶ Integrate out *parameters*.



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

Dual Probabilistic PCA

- Define *linear-Gaussian relationship* between latent variables and data.
- **Novel** Latent variable approach:
 - ▶ Define Gaussian prior over *parameters*, \mathbf{W} .
 - ▶ Integrate out *parameters*.

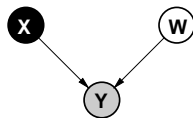


$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^n N(\mathbf{y}_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{W}) = \prod_{i=1}^d N(\mathbf{w}_{i,:} | \mathbf{0}, \mathbf{I})$$

Dual Probabilistic PCA

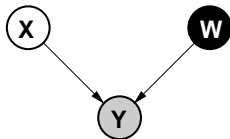
- Define *linear-Gaussian relationship* between latent variables and data.
- **Novel** Latent variable approach:
 - ▶ Define Gaussian prior over *parameters*, \mathbf{W} .
 - ▶ Integrate out *parameters*.



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

$$p(\mathbf{W}) = \prod_{i=1}^d N(\mathbf{w}_{i,:} | \mathbf{0}, \mathbf{I})$$

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^d N(\mathbf{y}_{:,j} | \mathbf{0}, \mathbf{X}\mathbf{X}^T + \sigma^2 \mathbf{I})$$



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

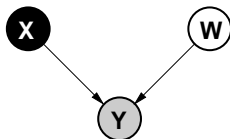
$$p(\mathbf{Y}|\mathbf{W}) = \prod_{i=1}^n N(\mathbf{y}_{i,:} | \mathbf{0}, \mathbf{C}), \quad \mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$$

$$\log p(\mathbf{Y}|\mathbf{W}) = -\frac{n}{2} \log |\mathbf{C}| - \frac{1}{2} \text{tr}(\mathbf{C}^{-1} \mathbf{Y}^T \mathbf{Y}) + \text{const.}$$

If \mathbf{U}_q are first q principal eigenvectors of $n^{-1} \mathbf{Y}^T \mathbf{Y}$ and the corresponding eigenvalues are Λ_q ,

$$\mathbf{W} = \mathbf{U}_q \mathbf{L} \mathbf{V}^T, \quad \mathbf{L} = (\Lambda_q - \sigma^2 \mathbf{I})^{\frac{1}{2}}$$

where \mathbf{V} is an arbitrary rotation matrix.



$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^d N(\mathbf{y}_{:,j} | \mathbf{0}, \mathbf{X}\mathbf{X}^T + \sigma^2 \mathbf{I})$$

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^d N(\mathbf{y}_{:,j} | \mathbf{0}, \mathbf{K}), \quad \mathbf{K} = \mathbf{X}\mathbf{X}^T + \sigma^2 \mathbf{I}$$

$$\log p(\mathbf{Y}|\mathbf{X}) = -\frac{d}{2} \log |\mathbf{K}| - \frac{1}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y}\mathbf{Y}^T) + \text{const.}$$

If \mathbf{U}'_q are first q principal eigenvectors of $d^{-1} \mathbf{Y}\mathbf{Y}^T$ and the corresponding eigenvalues are Λ_q ,

$$\mathbf{X} = \mathbf{U}'_q \mathbf{L} \mathbf{V}^T, \quad \mathbf{L} = (\Lambda_q - \sigma^2 \mathbf{I})^{\frac{1}{2}}$$

where \mathbf{V} is an arbitrary rotation matrix.

The Eigenvalue Problems are equivalent

- Solution for Probabilistic PCA (solves for the mapping)

$$\mathbf{Y}^T \mathbf{Y} \mathbf{U}_q = \mathbf{U}_q \Lambda_q \quad \mathbf{W} = \mathbf{U}_q \mathbf{L} \mathbf{V}^T$$

- Solution for Dual Probabilistic PCA (solves for the latent positions)

$$\mathbf{Y} \mathbf{Y}^T \mathbf{U}'_q = \mathbf{U}'_q \Lambda_q \quad \mathbf{X} = \mathbf{U}'_q \mathbf{L} \mathbf{V}^T$$

- Equivalence is from

$$\mathbf{U}_q = \mathbf{Y}^T \mathbf{U}'_q \Lambda_q^{-\frac{1}{2}}$$

Prior for Functions

- Probability Distribution over Functions
- Functions are infinite dimensional.
 - ▶ Prior distribution over *instantiations* of the function: finite dimensional objects.
- Can prove by induction that GP is 'consistent'.
- Mean and Covariance Functions
 - ▶ Instead of mean and covariance matrix, GP is defined by mean function and covariance function.
 - ▶ Mean function often taken to be zero or constant.
 - ▶ Covariance function must be *positive definite*.
 - ▶ Class of valid covariance functions is the same as the class of *Mercer kernels*.

Zero mean Gaussian Process

- A (zero mean) Gaussian process likelihood is of the form

$$p(\mathbf{y}|\mathbf{X}) = N(\mathbf{y}|\mathbf{0}, \mathbf{K}),$$

where \mathbf{K} is the covariance function or *kernel*.

- The *linear kernel* with noise has the form

$$\mathbf{K} = \mathbf{X}\mathbf{X}^T + \sigma^2\mathbf{I}$$

- Priors over non-linear functions are also possible.
 - ▶ To see what functions look like, we can sample from the prior process.

demCovFuncSample

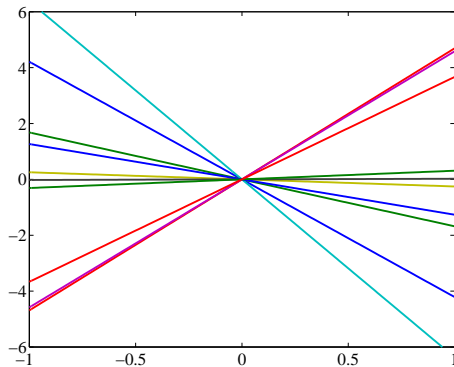


Figure: linear kernel, $\mathbf{K} = \mathbf{X}\mathbf{X}^T$

demCovFuncSample

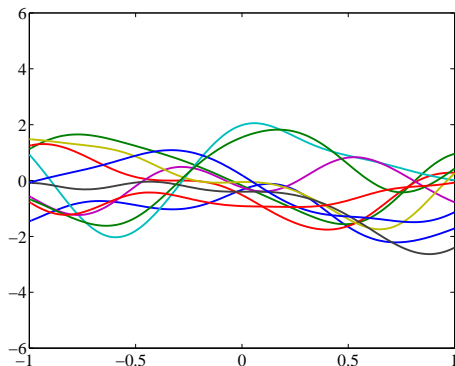


Figure: RBF kernel, $k_{i,j} = \alpha \exp\left(-\frac{1}{2l} \|\mathbf{x}_i - \mathbf{x}_j\|^2\right)$, with $l = 0.32$, $\alpha = 1$

Covariance Samples

demCovFuncSample

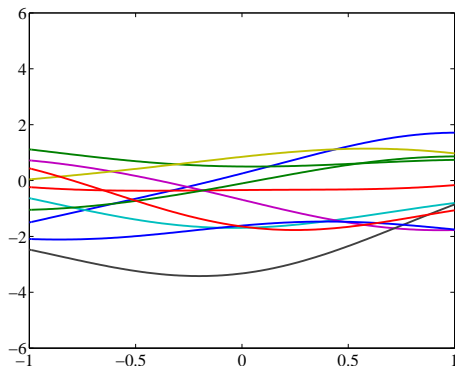


Figure: RBF kernel, $k_{i,j} = \alpha \exp\left(-\frac{1}{2l} \|\mathbf{x}_i - \mathbf{x}_j\|^2\right)$, with $l = 1$, $\alpha = 1$

demCovFuncSample

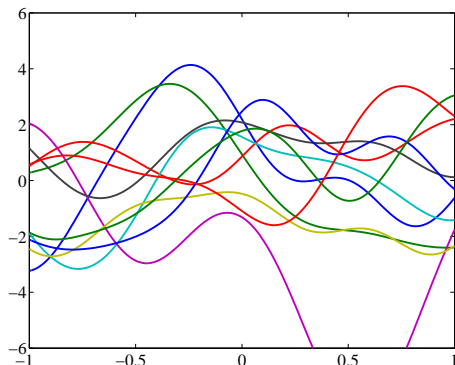


Figure: RBF kernel, $k_{i,j} = \alpha \exp\left(-\frac{1}{2l} \|\mathbf{x}_i - \mathbf{x}_j\|^2\right)$, with $l = 0.3$, $\alpha = 4$

Covariance Samples

demCovFuncSample

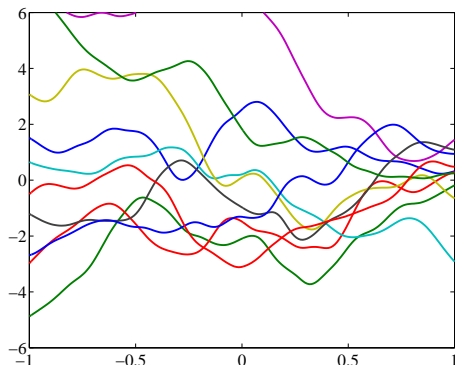


Figure: MLP kernel, $k_{i,j} = \alpha \sin^{-1} \left(\frac{w \mathbf{x}_i^T \mathbf{x}_j + b}{\sqrt{(w \mathbf{x}_i^T \mathbf{x}_i + b + 1)(w \mathbf{x}_j^T \mathbf{x}_j + b + 1)}} \right)$, with $\alpha = 8$, $w = 100$ and $b = 100$

Covariance Samples

demCovFuncSample

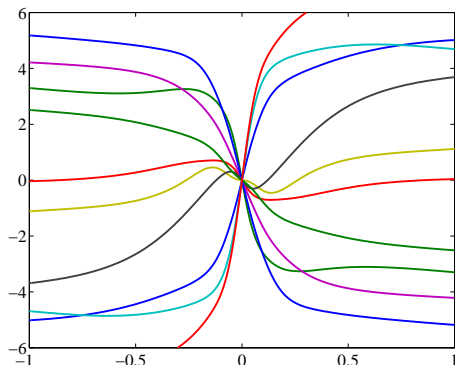


Figure: MLP kernel, $k_{i,j} = \alpha \sin^{-1} \left(\frac{w \mathbf{x}_i^T \mathbf{x}_j + b}{\sqrt{(w \mathbf{x}_i^T \mathbf{x}_i + b + 1)(w \mathbf{x}_j^T \mathbf{x}_j + b + 1)}} \right)$, with $\alpha = 8$, $b = 0$ and $w = 100$

demCovFuncSample

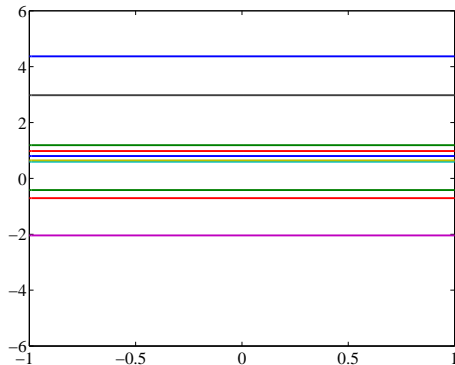


Figure: bias 'kernel', $k_{i,j} = \alpha$, with $\alpha = 1$ and

demCovFuncSample

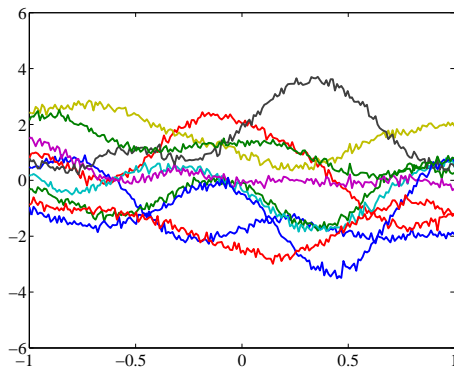


Figure: summed combination of: RBF kernel, $\alpha = 1$, $l = 0.3$; bias kernel, $\alpha = 1$; and white noise kernel, $\beta = 100$

Posterior Distribution over Functions

- Gaussian processes are often used for regression.
- We are given a known inputs \mathbf{X} and targets \mathbf{Y} .
- We assume a prior distribution over functions by selecting a kernel.
- Combine the prior with data to get a *posterior* distribution over functions.

Gaussian Process Regression

demRegression

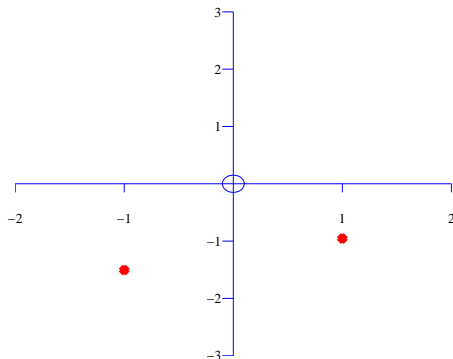


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

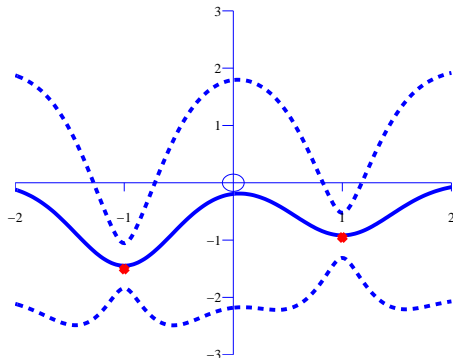


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

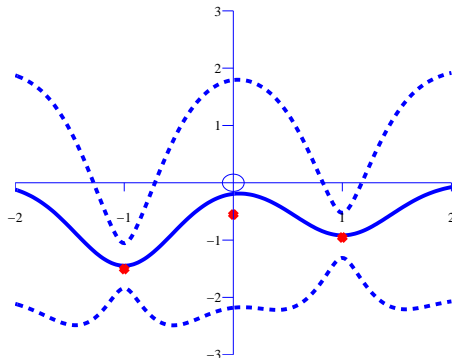


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

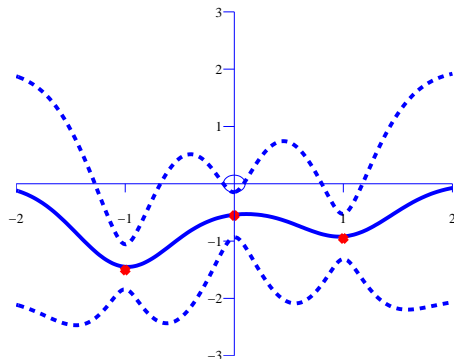


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

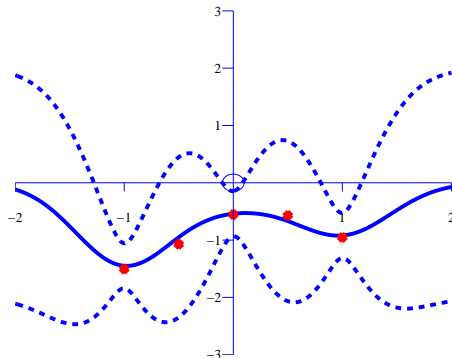


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

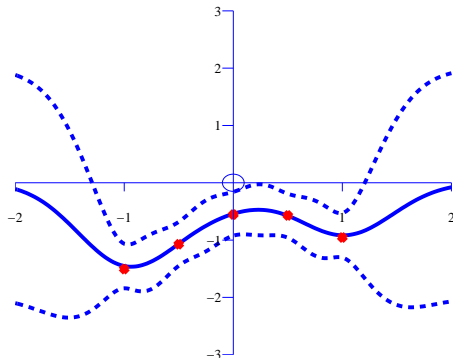


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

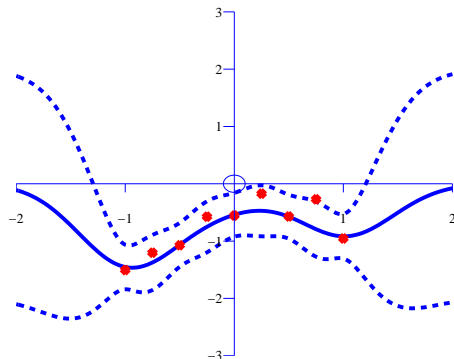


Figure: Examples include WiFi localization, C14 calibration curve.

Gaussian Process Regression

demRegression

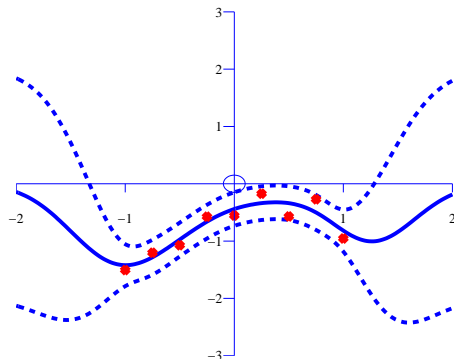
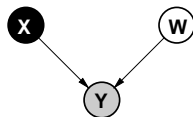


Figure: Examples include WiFi localization, C14 calibration curve.

Dual Probabilistic PCA

- Define *linear-Gaussian relationship* between latent variables and data.
- **Novel** Latent variable approach:
 - ▶ Define Gaussian prior over *parameters*, \mathbf{W} .
 - ▶ Integrate out *parameters*.



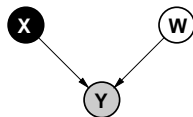
$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^n N(\mathbf{y}_{i,:} | \mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{W}) = \prod_{i=1}^d N(\mathbf{w}_{i,:} | \mathbf{0}, \mathbf{I})$$

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^d N(\mathbf{y}_{:,j} | \mathbf{0}, \mathbf{X}\mathbf{X}^T + \sigma^2 \mathbf{I})$$

Dual Probabilistic PCA

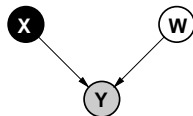
- Inspection of the marginal likelihood shows ...
 - ▶ The covariance matrix is a covariance function.
 - ▶ We recognise it as the 'linear kernel'.
 - ▶ We call this the Gaussian Process Latent Variable model (GP-LVM).



$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^d \mathcal{N}(\mathbf{y}_{:,j} | \mathbf{0}, \mathbf{X}\mathbf{X}^T + \sigma^2 \mathbf{I})$$

Dual Probabilistic PCA

- Inspection of the marginal likelihood shows ...
 - ▶ The covariance matrix is a covariance function.
 - ▶ We recognise it as the 'linear kernel'.
 - ▶ We call this the Gaussian Process Latent Variable model (GP-LVM).

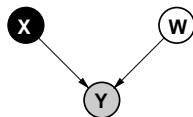


$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^d \mathcal{N}(\mathbf{y}_{:,j} | \mathbf{0}, \mathbf{K})$$

$$\mathbf{K} = \mathbf{X}\mathbf{X}^T + \sigma^2\mathbf{I}$$

Dual Probabilistic PCA

- Inspection of the marginal likelihood shows ...
 - ▶ The covariance matrix is a covariance function.
 - ▶ We recognise it as the 'linear kernel'.
 - ▶ We call this the Gaussian Process Latent Variable model (GP-LVM).



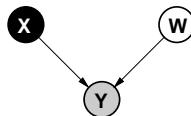
$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^d \mathcal{N}(\mathbf{y}_{:,j} | \mathbf{0}, \mathbf{K})$$

$$\mathbf{K} = \mathbf{X}\mathbf{X}^T + \sigma^2\mathbf{I}$$

This is a product of Gaussian processes with linear kernels.

Dual Probabilistic PCA

- Inspection of the marginal likelihood shows ...
 - ▶ The covariance matrix is a covariance function.
 - ▶ We recognise it as the 'linear kernel'.
 - ▶ We call this the Gaussian Process Latent Variable model (GP-LVM).



$$p(\mathbf{Y}|\mathbf{X}) = \prod_{j=1}^d N(\mathbf{y}_{:,j} | \mathbf{0}, \mathbf{K})$$

$$\mathbf{K} = ?$$

Replace linear kernel with non-linear kernel for non-linear model.

RBF Kernel

- The RBF kernel has the form $k_{ij} = k(\mathbf{x}_{i,:}, \mathbf{x}_{j,:})$, where

$$k(\mathbf{x}_{i,:}, \mathbf{x}_{j,:}) = \alpha \exp \left(-\frac{(\mathbf{x}_{i,:} - \mathbf{x}_{j,:})^T (\mathbf{x}_{i,:} - \mathbf{x}_{j,:})}{2l^2} \right).$$

- No longer possible to optimise wrt \mathbf{X} via an eigenvalue problem.
- Instead find gradients with respect to \mathbf{X} , α , l and σ^2 and optimise using conjugate gradients.

Traditional Model

- For a parameteric model maximise the marginal.

$$\log p(\mathbf{Y}|\mathbf{W}) = \log \int p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) p(\mathbf{X}) d\mathbf{X} = \prod_{i=1}^N \mathcal{N}(\mathbf{y}_{i,:} | \mathbf{0}, \mathbf{W}\mathbf{W}^T + \mathbf{I}\sigma^2)$$

leading to $\hat{\mathbf{W}}$.

- Evaluate the log likelihood of a new, test point \mathbf{y}_* ,

$$\log p(\mathbf{y}_* | \hat{\mathbf{W}}) = \log \mathcal{N}(\mathbf{y}_* | \mathbf{0}, \hat{\mathbf{W}}\hat{\mathbf{W}}^T + \mathbf{I}\sigma^2).$$

New Model

- For non-parameteric model maximise the marginal

$$\log p(\mathbf{Y}|\mathbf{X}) = \log \int p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) p(\mathbf{W}) d\mathbf{W} = \log \prod_{j=1}^d \mathcal{N}(\mathbf{y}_{:,j} | \mathbf{0}, \mathbf{K})$$

leading to $\hat{\mathbf{X}}$.

- Evaluate the log likelihood of a new, test point \mathbf{y}_* ??

$$\log p(\mathbf{y}_*, \mathbf{Y} | \mathbf{x}_*, \hat{\mathbf{X}}) = \log p(\mathbf{y}_* | \mathbf{Y}, \mathbf{x}_*, \hat{\mathbf{X}}) + \log p(\mathbf{Y} | \hat{\mathbf{X}}).$$

Maximise New Latent Variable

- Maximise log likelihood with respect to \mathbf{x}_* ,

$$\log p(\mathbf{y}_* | \mathbf{Y}, \mathbf{x}_*, \hat{\mathbf{X}}) = \sum_{j=1}^d \log \mathcal{N}(y_{j*} | \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{y}_{:,j}, k_{*,*} - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_*)$$

where the new covariance (of the joint process) is partitioned as

$$\mathbf{K}' = \begin{bmatrix} \mathbf{K} & \mathbf{k}_* \\ \mathbf{k}_*^T & k_{*,*} \end{bmatrix}$$

Multiple Test Data

- What if we are given two data points though? \mathbf{y}_* and \mathbf{y}_{**}

$$\mathbf{K}'' = \begin{bmatrix} \mathbf{K} & \mathbf{K}_* \\ \mathbf{K}_* & \mathbf{K}_{*,*} \end{bmatrix}$$

$$\log p(\mathbf{y}_*, \mathbf{y}_{**} | \mathbf{Y}, \mathbf{x}_*, \mathbf{x}_{**}, \hat{\mathbf{X}}) = \sum_{j=1}^d \log \mathcal{N}([y_{j*} \ y_{j**}]^T | \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{y}_{:,j}, \mathbf{K}_{*,*} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_*)$$

- But this is a joint covariance!!! It is not the same as maximising the individual data twice:

$$\begin{aligned} \log p(\mathbf{y}_* | \mathbf{Y}, \mathbf{x}_*, \hat{\mathbf{X}}) p(\mathbf{y}_{**} | \mathbf{Y}, \mathbf{x}_{**}, \hat{\mathbf{X}}) &= \log \sum_{j=1}^d \mathcal{N}(y_{j*} | \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{y}_{:,j}, k_{*,*} - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_*) \\ &\quad + \log \sum_{j=1}^d \mathcal{N}(y_{j**} | \mathbf{k}_{**}^T \mathbf{K}^{-1} \mathbf{y}_{:,j}, k_{**,**} - \mathbf{k}_{**}^T \mathbf{K}^{-1} \mathbf{k}_{**}) \end{aligned}$$

- What does it mean?

Possible Explanation?

Comment from John yesterday:

- For Gaussians maximising $\log p(x, y | \Sigma)$ wrt Σ is equivalent to $\log p(x|y, \Sigma) + \log p(y|x, \Sigma)$ wrt Σ .
- So perhaps maximising

$$\log p(\mathbf{y}_*, \mathbf{y}_{**} | \mathbf{Y}, \mathbf{x}_*, \mathbf{x}_{**}, \hat{\mathbf{X}}) \equiv \log p(\mathbf{y}_*, \mathbf{y}_{**} | \mathbf{Y}, \mathbf{K}'')$$

is equivalent to maximising

$$\log p(\mathbf{y}_* | \mathbf{y}_{**}, \mathbf{Y}, \mathbf{K}'') + \log p(\mathbf{y}_{**} | \mathbf{y}_*, \mathbf{Y}, \mathbf{K}'')$$

- *cf*

$$\log p(\mathbf{y}_* | \mathbf{Y}, \mathbf{K}'') + \log p(\mathbf{y}_{**} | \mathbf{Y}, \mathbf{K}'')$$

- 1 Probabilistic Dimensionality Reduction
- 2 Examples
- 3 Model Extensions
- 4 Conclusions

Generalization with less Data than Dimensions

- Powerful uncertainty handling of GPs leads to surprising properties.
- Non-linear models can be used where there are fewer data points than dimensions *without overfitting*.
- Example: Modelling a stick man in 102 dimensions with 55 data points!

demStick1

Figure: The latent space for the stick man motion capture data.

demStick1

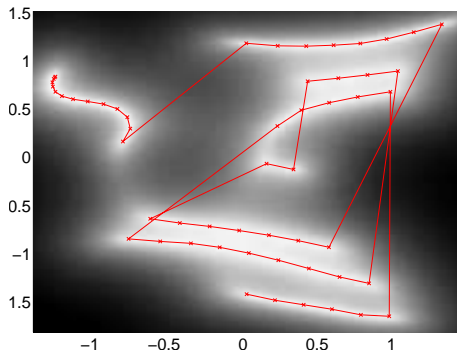


Figure: The latent space for the stick man motion capture data.

Style Based Inverse Kinematics

- Facilitating animation through modelling human motion with the GP-LVM [Grochow et al., 2004]

Tracking

- Tracking using models of human motion learnt with the GP-LVM [Urtasun et al., 2005, 2006]

.

Outline

- 1 Probabilistic Dimensionality Reduction
- 2 Examples
- 3 Model Extensions**
- 4 Conclusions

Local Distance Preservation [Lawrence and Quiñonero Candela, 2006]

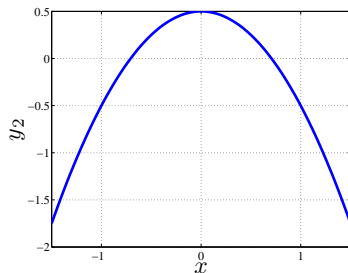
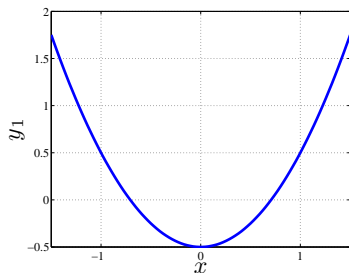
- Most dimensional reduction techniques preserve local distances.
- The GP-LVM does not.
- GP-LVM maps smoothly from latent to data space.
- Points close in latent space are close in data space.
 - ▶ This does not imply points close in data space are close in latent space.
- Many methods map smoothly from data to latent space.
 - ▶ Points close in data space are close in latent space.
 - ▶ This does not imply points close in latent space are close in data space.

Back Constraints II

Forward Mapping (demBackMapping in oxford toolbox)

- Mapping from 1-D latent space to 2-D data space.

$$y_1 = x^2 - 0.5, \quad y_2 = -x^2 + 0.5$$

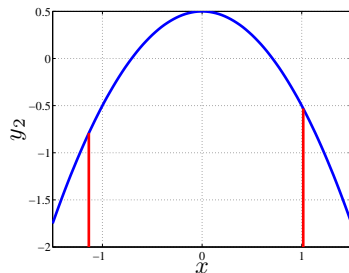
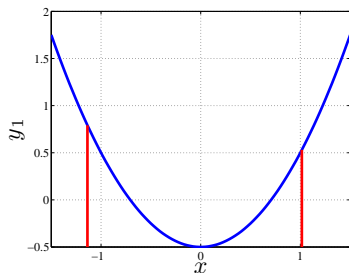


Back Constraints II

Forward Mapping (demBackMapping in oxford toolbox)

- Mapping from 1-D latent space to 2-D data space.

$$y_1 = x^2 - 0.5, \quad y_2 = -x^2 + 0.5$$

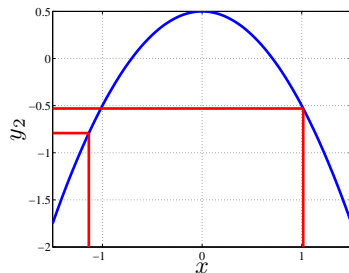
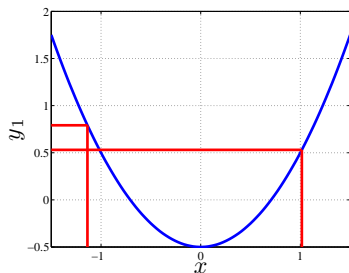


Back Constraints II

Forward Mapping (demBackMapping in oxford toolbox)

- Mapping from 1-D latent space to 2-D data space.

$$y_1 = x^2 - 0.5, \quad y_2 = -x^2 + 0.5$$

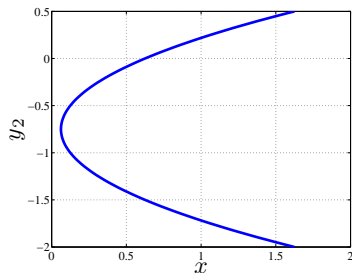
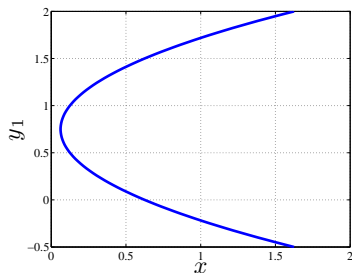


Back Constraints II

Backward Mapping (demBackMapping in oxford toolbox)

- Mapping from 2-D data space to 1-D latent.

$$x = 0.5 (y_1^2 + y_2^2 + 1)$$

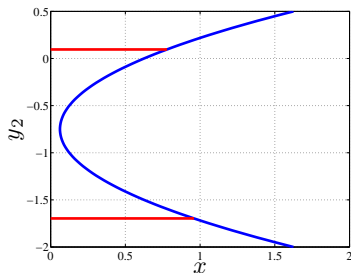
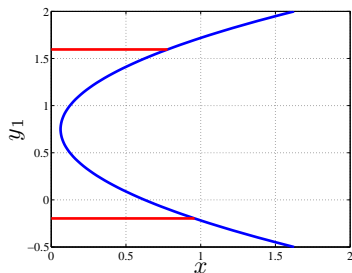


Back Constraints II

Backward Mapping (demBackMapping in oxford toolbox)

- Mapping from 2-D data space to 1-D latent.

$$x = 0.5 (y_1^2 + y_2^2 + 1)$$

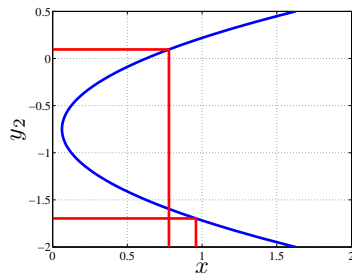
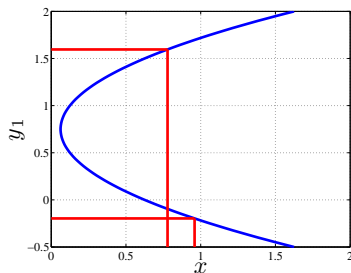


Back Constraints II

Backward Mapping (demBackMapping in oxford toolbox)

- Mapping from 2-D data space to 1-D latent.

$$x = 0.5 (y_1^2 + y_2^2 + 1)$$



Multi-Dimensional Scaling with a Mapping

- Lowe and Tipping [1997] made latent positions a function of the data.

$$x_{ij} = f_j(\mathbf{y}_i; \mathbf{w})$$

- Function was either multi-layer perceptron or a radial basis function network.
- Their motivation was different from ours:
- They wanted to add the advantages of a true mapping to multi-dimensional scaling.

Back Constraints

- We can use the same idea to force the GP-LVM to respect local distances.[Lawrence and Quiñero Candela, 2006]
 - ▶ By constraining each \mathbf{x}_i to be a 'smooth' mapping from \mathbf{y}_i local distances can be respected.
 - ▶ This works because in the GP-LVM we maximise wrt latent variables, we don't integrate out.
- Can use any 'smooth' function:
 - 1 Neural network.
 - 2 RBF Network.
 - 3 Kernel based mapping.

Computing Gradients

- GP-LVM normally proceeds by optimising

$$L(\mathbf{X}) = \log p(\mathbf{Y}|\mathbf{X})$$

with respect to \mathbf{X} using $\frac{dL}{d\mathbf{X}}$.

- The back constraints are of the form

$$x_{ij} = f_j(\mathbf{y}_{i,:}; \mathbf{B})$$

where \mathbf{B} are parameters.

- We can compute $\frac{dL}{d\mathbf{B}}$ via chain rule and optimise parameters of mapping.

demStick1 and demStick3

Figure: The latent space for the motion capture data with (*right*) and without (*left*) dynamics. The dynamics use a Gaussian process with an RBF kernel.

.

Motion Capture Results

demStick1 and demStick3

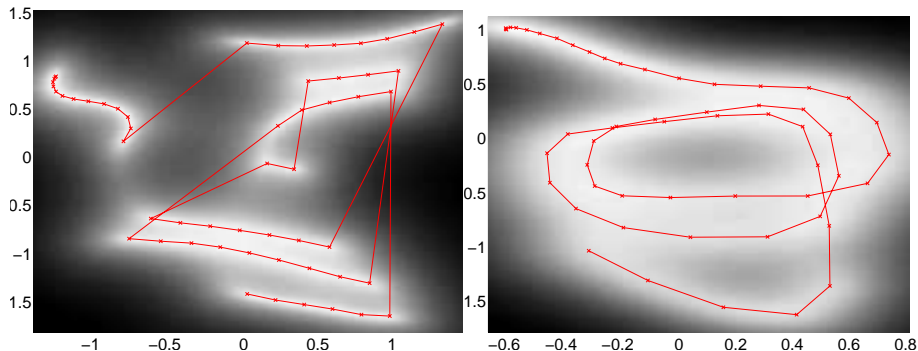
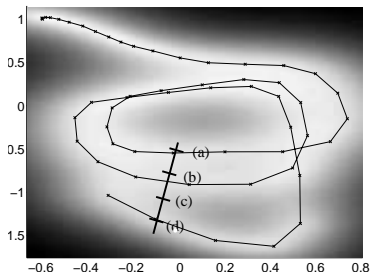


Figure: The latent space for the motion capture data with (*right*) and without (*left*) dynamics. The dynamics use a Gaussian process with an RBF kernel.

Stick Man Results

demStickResults



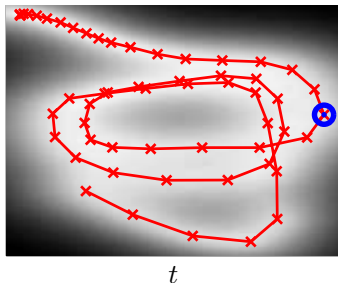
Projection into data space from four points in the latent space. The inclination of the runner changes becoming more upright.

MAP Solutions for Dynamics Models

- Data often has a temporal ordering.
 - ▶ Markov-based dynamics are often used.
- For the GP-LVM
 - ▶ Marginalising such dynamics is intractable.
 - ▶ But: MAP solutions are trivial to implement.
 - ▶ Many choices: Kalman filter, Markov chains *etc.*.
 - ▶ Wang et al. [2006] suggest using a Gaussian Process.

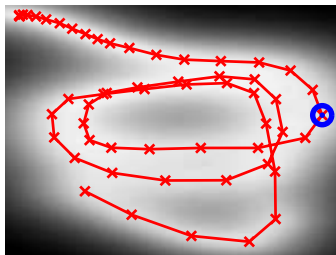
GP-LVM with Dynamics

- Autoregressive Gaussian process mapping in latent space between time points.

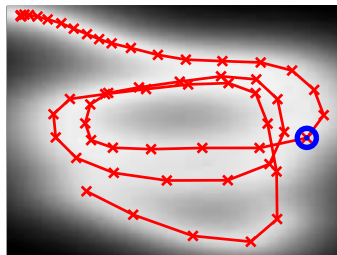


GP-LVM with Dynamics

- Autoregressive Gaussian process mapping in latent space between time points.



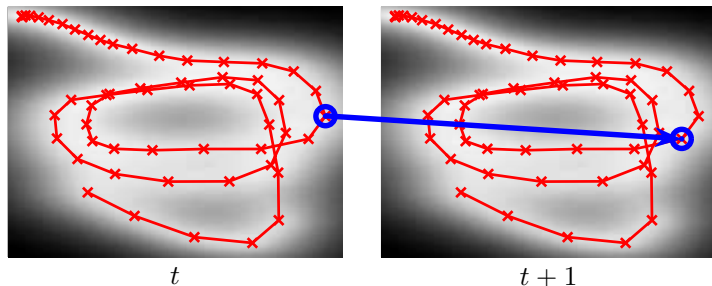
t



$t + 1$

GP-LVM with Dynamics

- Autoregressive Gaussian process mapping in latent space between time points.



demStick1 and demStick2

Figure: The latent space for the motion capture data without dynamics (*left*), with auto-regressive dynamics (*right*) based on an RBF kernel.

Motion Capture Results

demStick1 and demStick2

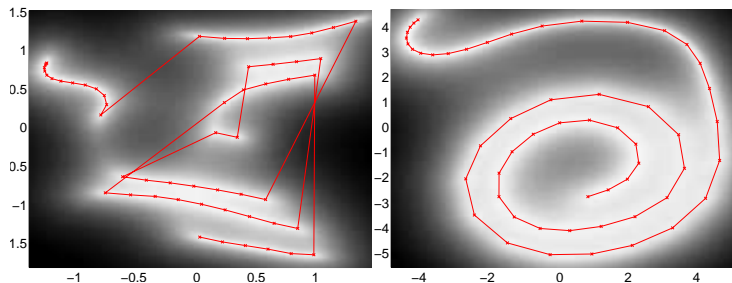
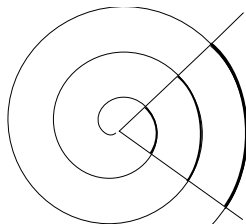


Figure: The latent space for the motion capture data without dynamics (*left*), with auto-regressive dynamics (*right*) based on an RBF kernel.

Inner Groove Distortion

- Autoregressive unimodal dynamics, $p(\mathbf{x}_t | \mathbf{x}_{t-1})$.
 - ▶ Forces spiral visualisation.
 - ▶ Poorer model due to inner groove distortion.



Direct use of Time Variable

- Instead of auto-regressive dynamics, consider regressive dynamics.
- Take \mathbf{t} as an input, use a prior $p(\mathbf{X}|\mathbf{t})$.
- User a Gaussian process prior for $p(\mathbf{X}|\mathbf{t})$.
- Also allows us to consider variable sample rate data.

Motion Capture Results

demStick1, demStick2 and demStick5

Figure: The latent space for the motion capture data without dynamics (*left*), with auto-regressive dynamics (*middle*) and with regressive dynamics (*right*) based on an RBF kernel.

Motion Capture Results

demStick1, demStick2 and demStick5

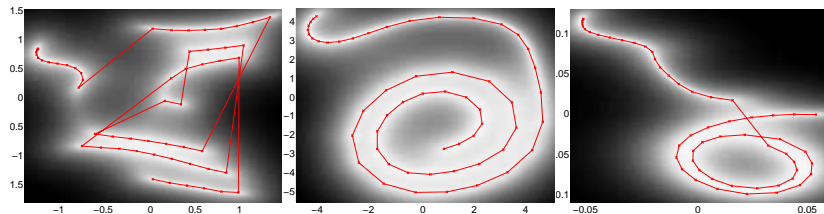


Figure: The latent space for the motion capture data without dynamics (*left*), with auto-regressive dynamics (*middle*) and with regressive dynamics (*right*) based on an RBF kernel.

MAP Solutions for Dynamics Models

- Autoregressive Gaussian Processes. Wang et al. [2006]

Force the Model to Respect Local Distances

- Back constrained GP-LVM.

Developments Made Under Pump Priming Grant

- Sparse Approximations for Large Data Sets
- Hierarchical Models for Subject Decomposition
- Three Dimensional Pose Reconstruction from Images

Stacking Gaussian Processes

- Regressive dynamics provides a simple hierarchy.
- The input space of the GP is governed by another GP.
- By stacking GPs we can consider more complex hierarchies.
- Ideally we should marginalise latent spaces
 - ▶ In practice we seek MAP solutions.

Two Correlated Subjects

demHighFive1

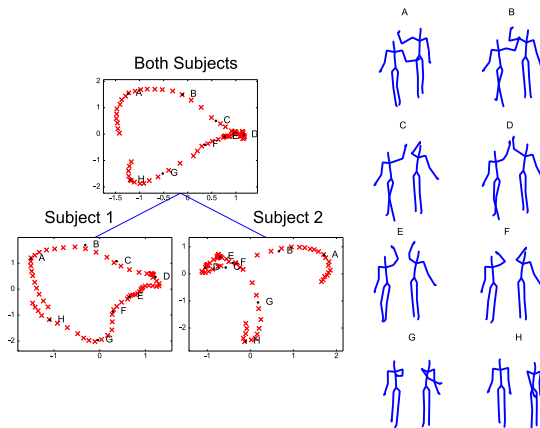


Figure: Hierarchical model of a 'high five'.

Decomposition of Body

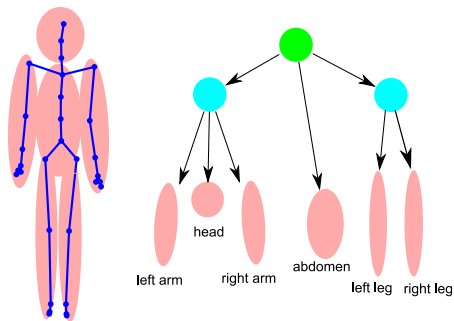


Figure: Decomposition of a subject.

Single Subject Run/Walk

demRunWalk1

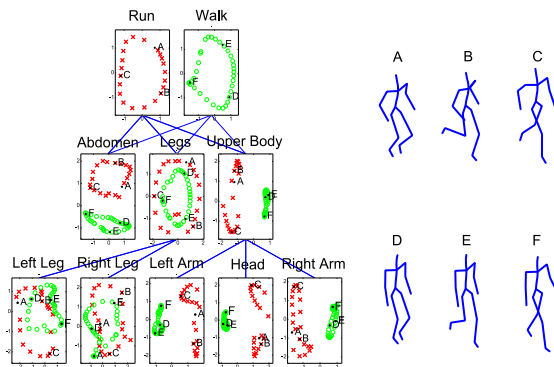


Figure: Hierarchical model of a walk and a run.

- Gaussian processes inherently
 - ▶ $O(N^3)$ complexity,
 - ▶ $O(N^2)$ storage.
- Sparse Gaussian processes normally give
 - ▶ $O(k^2N)$ complexity,
 - ▶ $O(kN)$ storage
- FITC Approximation [Snelson and Ghahramani, 2006, Quiñonero Candela and Rasmussen, 2005, Presented/Developed at PASCAL workshop!].

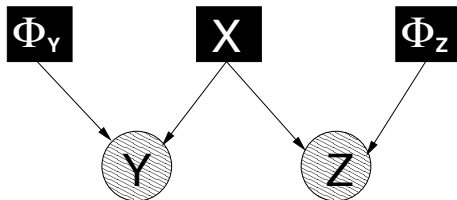
- Recreate results of Taylor et al. [2007] on human motion capture data set.
- Data was walking and running motions from subject 35 in the CMU Mocap data base.
- Used dynamical refinement of the GP-LVM proposed by Wang et al. [2006]
- Taylor et al. [2007] applied their binary latent variable model to two missing data problems
 - ▶ right leg was removed from the test sequence
 - ▶ upper body was removed.
- Reconstruction obtained compared with nearest neighbour.

- Used the FITC approximation with 100 inducing points.
- The models were back constrained [Lawrence and Quiñonero Candela, 2006] .
- The data set size was 2613 frames.

Root mean squared angle error results on test data.

Data	Leg	Body
GP-LVM ($q = 3$)	3.40	2.49
GP-LVM ($q = 4$)	3.38	2.72
GP-LVM ($q = 5$)	4.25	2.78
NN (s)	4.44	2.62
NN	4.11	3.20

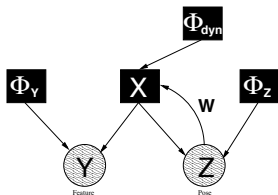
Table: NN: nearest neighbour, NN (s): nearest neighbour in scaled space, GP-LVM (latent dimension): the GP-LVM with different latent dimensions, q .



- Learn two separate kernels from a single shared latent representation \mathbf{X} [Shon et al., 2006]
- Objective

$$p(\mathbf{Y}, \mathbf{Z} | \mathbf{X}, \Phi_Y, \Phi_Z) = p(\mathbf{Y} | \mathbf{X}, \Phi_Y) p(\mathbf{Z} | \mathbf{X}, \Phi_Z)$$

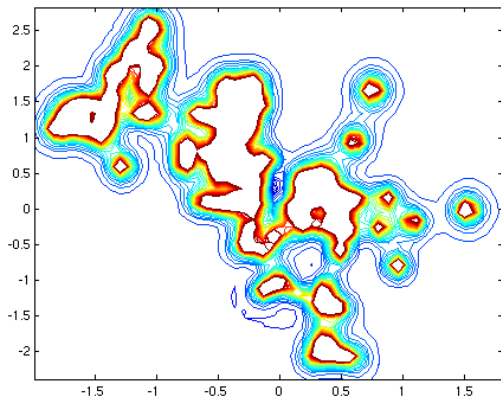
Shared GP-LVM Experiments¹



- Silhouette Features: $\mathbf{y}_i \in \mathbb{R}^{100}$, Pose Parameters: $\mathbf{z}_i \in \mathbb{R}^{54}$
- **Back constraints:** force bijective mapping between latent space and pose [Lawrence and Quiñero Candela, 2006].
- **Dynamics:** add GP auto regressive dynamics to latent space [Wang et al., 2006].
- **Artificially generated training data:** from Agarwal and Triggs [2006].

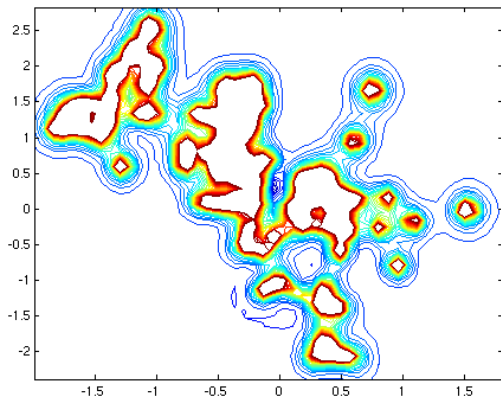
¹Ek et al. [2007]

Shared GP-LVM Experiments

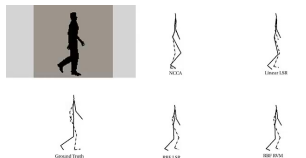


- Highly multimodal latent space given silhouette.

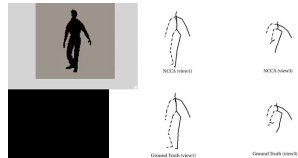
Shared GP-LVM Experiments



- Highly multimodal latent space given silhouette.

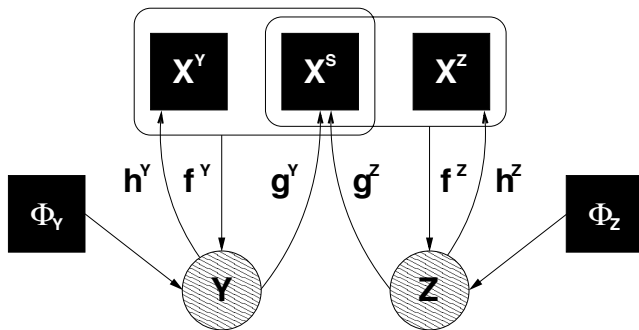


`runcca_all.sh`



`runcca_only.sh`

Modified Model



Shared Latent space by kernel CCA:

- Find directions $\{\mathbf{W}_Y, \mathbf{W}_Z\}$ in each feature space maximizing the correlation
- Canonical variate
$$\begin{cases} \mathbf{a}_Y &= \mathbf{Y}\mathbf{W}_Y \\ \mathbf{a}_Z &= \mathbf{Z}\mathbf{W}_Z \end{cases}$$
- Solution through Eigenvalue problem.

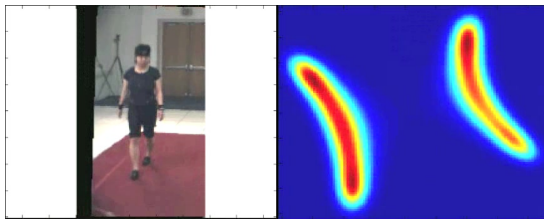
Non Shared Latent Space

- Find further directions *orthogonal* to CCA directions of maximum variance.
- We named these non-consolidating components.
- Solution through eigenvalue problem.

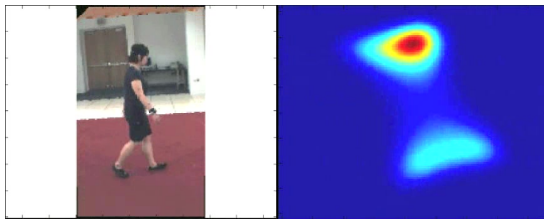
Feature Spaces:

- Many possible choices of feature space
 - 1 Linear Kernel
 - 2 RBF
 - 3 Maximum Variance Unfolding, Isomap
- Choose between them using GP-LVM likelihood [Harmeling, 2007].

Results of Initialisation

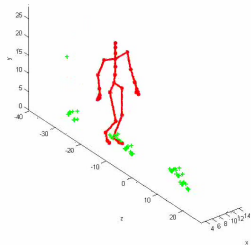


`runspectral_test.sh`

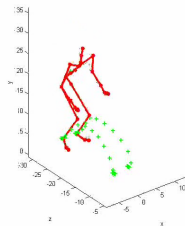


`runspectral_test.sh`

Topological Constraints



`runtopology_walk.sh`



`runtopology_jump.sh`

Outline

- 1 Probabilistic Dimensionality Reduction
- 2 Examples
- 3 Model Extensions
- 4 Conclusions

- GP-LVM is a Probabilistic Non-Linear Generalisation of PCA.
- Works Effectively as a Probabilistic Model in High Dimensional Spaces.
- Back constraints can be introduced to force local distance preservation.
- Dynamics can be introduced for modelling data with a temporal structure.
- Hierarchical models can encode conditional independencies.
- Topologically constraints can be imposed.



References

- A. Agarwal and B. Triggs. Recovering 3D human pose from monocular images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(1), 2006. doi: 10.1109/TPAMI.2006.21.
- C. H. Ek, P. H. Torr, and N. D. Lawrence. Gaussian process latent variable models for human pose estimation. In *MLMI 2007*, 2007. To appear.
- K. Grochow, S. L. Martin, A. Hertzmann, and Z. Popovic. Style-based inverse kinematics. In *ACM Transactions on Graphics (SIGGRAPH 2004)*, pages 522–531, 2004. doi: 10.1145/1186562.1015755.
- S. Harmeling. Exploring model selection techniques for nonlinear dimensionality reduction. Technical Report EDI-INF-RR-0960, University of Edinburgh, 2007.
- N. D. Lawrence. Gaussian process models for visualisation of high dimensional data. In S. Thrun, L. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems*, volume 16, pages 329–336, Cambridge, MA, 2004. MIT Press.
- N. D. Lawrence. Probabilistic non-linear principal component analysis with Gaussian process latent variable models. *Journal of Machine Learning Research*, 6:1783–1816, 11 2005.
- N. D. Lawrence and J. Quiñero Candela. Local distance preservation in the GP-LVM through back constraints. In W. Cohen and A. Moore, editors, *Proceedings of the International Conference in Machine Learning*, volume 23, pages 513–520. Omnipress, 2006. ISBN 1-59593-383-2. doi: 10.1145/1143844.1143909.
- D. Lowe and M. E. Tipping. Neuroscale: Novel topographic feature extraction with radial basis function networks. In M. C. Mozer, M. I. Jordan, and T. Petsche, editors, *Advances in Neural Information Processing Systems*, volume 9, pages 543–549, Cambridge, MA, 1997. MIT Press.
- J. Quiñero Candela and C. E. Rasmussen. A unifying view of sparse approximate Gaussian process regression. *Journal of Machine Learning Research*, 6:1939–1959, 2005.
- A. P. Shon, K. Grochow, A. Hertzmann, and R. P. N. Rao. Learning shared latent structure for image synthesis and robotic imitation. In Y. Weiss, B. Schölkopf, and J. C. Platt, editors, *Advances in Neural Information Processing Systems*, volume 18, Cambridge, MA, 2006. MIT Press.
- E. Snelson and Z. Ghahramani. Sparse Gaussian processes using pseudo-inputs. In Y. Weiss, B. Schölkopf, and J. C. Platt, editors, *Advances in Neural Information Processing Systems*, volume 18, Cambridge, MA, 2006. MIT Press.
- G. W. Taylor, G. E. Hinton, and S. Roweis. Modeling human motion using binary latent variables. In B. Schölkopf, J. C. Platt, and T. Hofmann, editors, *Advances in Neural Information Processing Systems*, volume 19, Cambridge, MA, 2007. MIT Press.
- M. E. Tipping and C. M. Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society, B*, 6(3): 611–622, 1999.
- R. Urtasun, D. J. Fleet, A. Hertzmann, and P. Fua. Priors for people tracking from small training sets. In *IEEE International Conference on Computer Vision (ICCV)*, pages 403–410, Beijing, China, 17–21 Oct. 2005. IEEE Computer Society Press.
- R. Urtasun, D. J. Fleet, and P. Fua. 3D people tracking with Gaussian process dynamical models. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*, pages 238–245, New York, U.S.A., 17–22 Jun. 2006. IEEE