Latent variable and deep modeling with Gaussian processes; application to system identification

Andreas Damianou

Department of Computer Science, University of Sheffield, UK

Brown University, 17 Feb. 2016

Outline

Part 1: Introduction Recap from yesterday

Part 2: Incorporating latent variables

Unsupervised GP learning Structure in the latent space - representation learning Bayesian Automatic Relevance Determination Deep GP learning Multi-view GP learning

Part 3: Dynamical systems

Policy learning Autoregressive Dynamics Going deeper: Deep Recurrent Gaussian Process Regressive dynamics with deep GPs

Outline

Part 1: Introduction Recap from yesterday

Part 2: Incorporating latent variables

Unsupervised GP learning Structure in the latent space - representation learning Bayesian Automatic Relevance Determination Deep GP learning Multi-view GP learning

Part 3: Dynamical systems

Policy learning Autoregressive Dynamics Going deeper: Deep Recurrent Gaussian Process Regressive dynamics with deep GPs

- Gaussian processes as infinite dimensional Gaussian distributions
- \blacktriangleright \Rightarrow can be used as priors over functions
- Non-parametric: training data act as parameters
- Uncertainty Quantification
- Learning from scarce data

Notation and graphical models

Graphical models



- White nodes: Observed variables
- Shaded nodes: Unobserved, or *latent* variables.
- Convention: Sometimes the latent function will be placed next to the arrow; sometimes I'll explicitly include the collection of instantiations f

Outline

Part 1: Introduction Recap from yesterday

Part 2: Incorporating latent variables

Unsupervised GP learning

Structure in the latent space - representation learning

Bayesian Automatic Relevance Determination Deep GP learning Multi-view GP learning

Part 3: Dynamical systems

Policy learning Autoregressive Dynamics Going deeper: Deep Recurrent Gaussian Process Regressive dynamics with deep GPs

Unsupervised learning: GP-LVM



If X is unobserved, treat it as a parameter and optimize over it.

- ► GP-LVM is interpreted as non-linear, non-parametric probabilistic PCA (Lawrence, JMLR 2015).
- Objective (likelihood function) is similar to GPs:

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}|\mathbf{x}) d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, K_{ff} + \sigma^2 \mathbf{I})$$

but now x's are optimized too.

Neil Lawrence, JMLR, 2005.

Unsupervised learning: GP-LVM



- If X is unobserved, treat it as a parameter and optimize over it.
- ► GP-LVM is interpreted as non-linear, non-parametric probabilistic PCA (Lawrence, JMLR 2015).
- Objective (likelihood function) is similar to GPs:

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}|\mathbf{x}) d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, K_{ff} + \sigma^2 \mathbf{I})$$

but now x's are optimized too.

Neil Lawrence, JMLR, 2005.

Unsupervised learning: GP-LVM



- If X is unobserved, treat it as a parameter and optimize over it.
- ► GP-LVM is interpreted as non-linear, non-parametric probabilistic PCA (Lawrence, JMLR 2015).
- Objective (likelihood function) is similar to GPs:

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}|\mathbf{x}) d\mathbf{f} = \mathcal{N}(\mathbf{y}|\mathbf{0}, K_{ff} + \sigma^2 \mathbf{I})$$

but now x's are optimized too.

Neil Lawrence, JMLR, 2005.

Fitting the GP-LVM



Fitting the GP-LVM

Figure credits: C. H. Ek



Fitting the GP-LVM

Figure credits: C. H. Ek



- Additional difficulty: x's are also missing!
- Improvement: Invoke the Bayesian methodology to find x's too.

Additionally integrate out the latent space.

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{f}) p(\mathbf{f}|\mathbf{x}) p(\mathbf{x}) d\mathbf{f} d\mathbf{x}$$

where

$$p(\mathbf{x}) \sim \mathcal{N}(0, \mathbf{I})$$

Titsias and Lawrence 2010, AISTATS; Damianou et al. 2015, JMLR

Automatic dimensionality detection

- ▶ In general, $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ with $\mathbf{x}^{(n)} \in \Re^Q$.
- Automatic dimenionality detection by employing automatic relevance determination (ARD) priors for the mapping f.
- $f \sim \mathcal{GP}(\mathbf{0}, k_f)$ with:

$$k_f\left(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}\right) = \sigma^2 \exp\left(-\frac{1}{2}\sum_{q=1}^{Q} w^{(q)} \left(x^{(i,q)} - x^{(j,q)}\right)^2\right)$$

Example:





Outline

Part 1: Introduction Recap from yesterday

Part 2: Incorporating latent variables

Unsupervised GP learning Structure in the latent space - representation learning Bayesian Automatic Relevance Determination Deep GP learning Multi-view GP learning

Part 3: Dynamical systems

Policy learning Autoregressive Dynamics Going deeper: Deep Recurrent Gaussian Process Regressive dynamics with deep GPs

Sampling from a deep GP



Model

$$\mathbf{x} = \text{given}$$

$$\mathbf{f}_{h} = f_{h}(x) = \mathcal{GP}(\mathbf{0}, k_{h}(x, x))$$

$$\mathbf{h} = f_{h}(x) + \epsilon_{h}$$

$$\mathbf{f}_{y} = f_{y}(h) = \mathcal{GP}(\mathbf{0}, k_{f}(h, h))$$

$$\mathbf{y} = f_{y}(h) + \epsilon_{y}$$
So:
$$\mathbf{y} = f_{y}(f_{h}(\mathbf{x}) + \epsilon_{h}) + \epsilon_{y}$$

Objective:

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}|\mathbf{f}_y) p(\mathbf{f}_y|\mathbf{h}) p(\mathbf{h}|\mathbf{f}_h) p(\mathbf{f}_h|\mathbf{x}) d\mathbf{f}_y d\mathbf{f}_h d\mathbf{h}$$

Damianou and Lawrence, AISTATS 2013; Damianou, PhD Thesis, 2015

Deep GP: Step function (credits for idea to J. Hensman, R. Calandra, M. Deisenroth)





-0.5

-1

0

1.5 2

1

Standard GP

Non-linear feature learning



- Successive warping creates "knots" which act as features.
- Features discovered in one layer remain in the next ones (ie knots are not un-tied)

 $f_1(f_2(f_3(x)))$

▶ With linear warpings (stacked PCA) we can't achieve this effect:

 $W_1(W_2(W_3x))) = W'x$

Deep GP: MNIST example



Outline

Part 1: Introduction Recap from yesterday

Part 2: Incorporating latent variables

Unsupervised GP learning Structure in the latent space - representation learning Bayesian Automatic Relevance Determination Deep GP learning Multi-view GP learning

Part 3: Dynamical systems

Policy learning Autoregressive Dynamics Going deeper: Deep Recurrent Gaussian Process Regressive dynamics with deep GPs

Multi-view: Manifold Relevance Determination (MRD)

Multi-view data arise from multiple information sources. These sources naturally contain some overlapping, or *shared* signal (since they describe the same "phenomenon"), but also have some *private* signal.

MRD: Model such data using overlapping sets of latent variables



Damianou et al., ICML 2012; Damianou, PhD Thesis, 2015

Multi-view: Manifold Relevance Determination (MRD)

Multi-view data arise from multiple information sources. These sources naturally contain some overlapping, or *shared* signal (since they describe the same "phenomenon"), but also have some *private* signal.

MRD: Model such data using overlapping sets of latent variables



Multi-view: Manifold Relevance Determination (MRD)

Multi-view data arise from multiple information sources. These sources naturally contain some overlapping, or *shared* signal (since they describe the same "phenomenon"), but also have some *private* signal.

MRD: Model such data using overlapping sets of latent variables



Deep GPs: Another multi-view example



Automatic structure discovery summary: ARD and MRD

Tools:

- ► ARD: Eliminate uncessary nodes/connections
- MRD: Conditional independencies
- Approximating evidence: Number of layers (?)



Automatic structure discovery summary: ARD and MRD

Tools:

- ► ARD: Eliminate uncessary nodes/connections
- MRD: Conditional independencies
- ► Approximating evidence: Number of layers (?)



More deepGP representation learning examples...

▶ https://youtu.be/s4zATH1TjG8

iCub interaction

https://youtu.be/uSQ0vxLcfVU

Faces - autoencoder

Outline

Part 1: Introduction Recap from yesterday

Part 2: Incorporating latent variables

Unsupervised GP learning Structure in the latent space - representation learning Bayesian Automatic Relevance Determination Deep GP learning Multi-view GP learning

Part 3: Dynamical systems

Policy learning Autoregressive Dynamics Going deeper: Deep Recurrent Gaussian Process Regressive dynamics with deep GPs Through a model, we wish to learn to:

- perform free simulation by learning patterns coming from the latent generation process (a mechanistic system we do not know)
- perform inter/extrapolation in time-series data which are very high-dimensional (e.g. video)
- detect outliers in data coming from a dynamical system
- optimize policies for control based on a model of the data.

Model-based policy learning

https://www.youtube.com/watch?v=XiigTGKZfks (Cart-pole)

http://mlg.eng.cam.ac.uk/?portfolio=andrew-mchutchon (Unicycle)

Work by: Marc Deisenroth, Andrew McHutchon, Carl Rasmussen

NARX model

A standard NARX model considers an input vector $\mathbf{x}_i \in \mathbb{R}^D$ comprised of L_y past observed outputs $y_i \in \mathbb{R}$ and L_u past exogenous inputs $u_i \in \mathbb{R}$:

$$\mathbf{x}_i = [y_{i-1}, \cdots, y_{i-L_y}, u_{i-1}, \cdots, u_{i-L_u}]^\top,$$

$$y_i = f(\mathbf{x}_i) + \epsilon_i^{(y)}, \qquad \epsilon_i^{(y)} \sim \mathcal{N}(\epsilon_i^{(y)}|0, \sigma_y^2),$$

Latent auto-regressive GP model:

$$x_{i} = f(x_{i-1}, \cdots, x_{i-L_{x}}u_{i-1}, \cdots, u_{i-L_{u}}) + \epsilon_{i}^{(x)},$$

$$y_{i} = x_{i} + \epsilon_{i}^{(y)},$$

<u>Contribution 1:</u> Simultaneous auto-regressive and representation learning. <u>Contribution 2:</u> Latents avoid the feedback of possibly corrupted observations into the dynamics.



Mattos, Damianou, Barreto, Lawrence, 2016

NARX model

A standard NARX model considers an input vector $\mathbf{x}_i \in \mathbb{R}^D$ comprised of L_y past observed outputs $y_i \in \mathbb{R}$ and L_u past exogenous inputs $u_i \in \mathbb{R}$:

$$\mathbf{x}_i = [y_{i-1}, \cdots, y_{i-L_y}, u_{i-1}, \cdots, u_{i-L_u}]^\top,$$

$$y_i = f(\mathbf{x}_i) + \epsilon_i^{(y)}, \qquad \epsilon_i^{(y)} \sim \mathcal{N}(\epsilon_i^{(y)}|0, \sigma_y^2),$$

Latent auto-regressive GP model:

$$x_{i} = f(x_{i-1}, \cdots, x_{i-L_{x}}u_{i-1}, \cdots, u_{i-L_{u}}) + \epsilon_{i}^{(x)},$$

$$y_{i} = x_{i} + \epsilon_{i}^{(y)},$$

<u>Contribution 1:</u> Simultaneous auto-regressive and representation learning.

<u>Contribution 2:</u> Latents avoid the feedback of possibly corrupted observations into the dynamics.



Mattos, Damianou, Barreto, Lawrence, 2016

Robustness to outliers

Latent auto-regressive GP model:

$$\begin{aligned} x_i &= f(x_{i-1}, \cdots, x_{i-L_x} u_{i-1}, \cdots, u_{i-L_u}) + \epsilon_i^{(x)}, \\ y_i &= x_i + \epsilon_i^{(y)}, \\ \epsilon_i^{(x)} &\sim \mathcal{N}(\epsilon_i^{(x)} | 0, \sigma_x^2), \\ \epsilon_i^{(y)} &\sim \mathcal{N}(\epsilon_i^{(y)} | 0, \tau_i^{-1}), \quad \tau_i \sim \Gamma(\tau_i | \alpha, \beta), \end{aligned}$$

<u>Contribution 3:</u> "Switching-off" outliers by including the above Student-t likelihood for the noise.

Mattos, Damianou, Barreto, Lawrence, DYCOPS 2016

Robust GP autoregressive model: demonstration



Figure: RMSE values for free simulation on test data with different levels of contamination by outliers.



Going deeper: Deep Recurrent Gaussian Process



Figure 1: RGP graphical model with H hidden layers.

 \tilde{x} is the lagged latent function values augmented with the lagged exogenous inputs.

Mattos, Dai, Damianou, Barreto, Lawrence, ICLR 2016



Inference is tricky...

$$\begin{split} \log p(\mathbf{y}) &\geq -\frac{N-L}{2} \sum_{h=1}^{H+1} \log 2\pi \sigma_h^2 - \frac{1}{2\sigma_{H+1}^2} \Big(\mathbf{y}^\top \mathbf{y} + \Psi_0^{(H+1)} \\ &- \operatorname{Tr} \left(\left(\mathbf{K}_z^{(H+1)} \right)^{-1} \Psi_2^{(H+1)} \right) \Big) + \frac{1}{2} \left| \mathbf{K}_z^{(H+1)} \right| - \frac{1}{2} \left| \mathbf{K}_z^{(H+1)} + \frac{1}{\sigma_{H+1}^2} \Psi_2^{(H+1)} \right| \\ &+ \frac{1}{2(\sigma_{H+1}^2)^2} \mathbf{y}^\top \Psi_1^{(H+1)} \left(\mathbf{K}_z^{(H+1)} + \frac{1}{\sigma_{H+1}^2} \Psi_2^{(H+1)} \right)^{-1} \left(\Psi_1^{(H+1)} \right)^\top \mathbf{y} \\ &+ \sum_{h=1}^{H} \left\{ -\frac{1}{2\sigma_h^2} \left(\sum_{i=L+1}^N \lambda_i^{(h)} + \left(\mu^{(h)} \right)^\top \mu^{(h)} + \Psi_0^{(h)} - \operatorname{Tr} \left(\left(\mathbf{K}_z^{(h)} \right)^{-1} \Psi_2^{(h)} \right) \right) \\ &+ \frac{1}{2} \left| \mathbf{K}_z^{(h)} \right| - \frac{1}{2} \left| \mathbf{K}_z^{(h)} + \frac{1}{\sigma_h^2} \Psi_2^{(h)} \right| \\ &+ \frac{1}{2(\sigma_h^2)^2} \left(\mu^{(h)} \right)^\top \Psi_1^{(h)} \left(\mathbf{K}_z^{(h)} + \frac{1}{\sigma_h^2} \Psi_2^{(h)} \right)^{-1} \left(\Psi_1^{(h)} \right)^\top \mu^{(h)} \\ &- \sum_{i=L+1}^N \int_{x_i^{(h)}} q \left(x_i^{(h)} \right) \log q \left(x_i^{(h)} \right) + \sum_{i=1}^L \int_{x_i^{(h)}} q \left(x_i^{(h)} \right) \log p \left(x_i^{(h)} \right) \right\}. \end{split}$$

Results in nonlinear systems identification:

- 1. artificial dataset
- 2. "drive" dataset: by a system with two electric motors that drive a pulley using a flexible belt.
 - input: the sum of voltages applied to the motors
 - output: speed of the belt.





Avatar control



Figure: The generated motion with a step function signal, starting with walking (blue), switching to running (red) and switching back to walking (blue).

Videos:

https://youtu.be/FR-oeGxV6yY Switching between learned speeds
 https://youtu.be/AT0HMtoPgic Interpolating (un)seen speed
 https://youtu.be/FuF-uZ83VMw Constant unseen speed

Regressive dynamics with deep GPs



Instead of coupling f's by encoding the Markov property, we couple them by **coupling the** f's inputs through another **GP** with time as input.

$$y = f(x) + \epsilon$$
$$x \sim \mathcal{GP}(0, k_x(t, t))$$
$$f \sim \mathcal{GP}(0, k_f(x, x))$$

Damianou et al., NIPS 2011; Damianou, Titsias and Lawrence, JMLR, 2015

Dynamics

- Dynamics are encoded in the covariance matrix $\mathbf{K} = k(\mathbf{t}, \mathbf{t})$.
- We can consider special forms for K.



Model individual sequences



Model periodic data



https://www.youtube.com/watch?v=fHDWloJtgk8 (mocap)

- Data-driven, model-based approach to real and control problems
- Gaussian processes: Uncertainty quantification / propagation gives an advantage
- Deep Gaussian processes: Representation learning + dynamics learning
- Future work: Deep Gaussian processes + mechanistic information; consider "real" applications.

Thanks!

Thanks to: Neil Lawrence, Michalis Titsias, Carl Henrik Ek, James Hensman, Cesar Lincoln Mattos, Zhenwen Dai, Javier Gonzalez, Tony Prescott, Uriel Martinez-Hernandez, Luke Boorman

Thanks to George Karniadakis, Paris Perdikaris for hosting me

BACKUP SLIDES 1: Deep GP Optimization - variational inference

MAP optimisation?





- MAP optimization is extremely problematic because:
 - Dimensionality of hs has to be decided a priori
 - Prone to overfitting, if \boldsymbol{h} are treated as parameters
 - Deep structures are not supported by the model's objective but have to be forced [Lawrence & Moore '07]
- We want:
 - To use the marginal likelihood as the objective: marg. lik. $= \int_{h_2,h_1} p(y|h_2)p(h_2|h_1)p(h_1|x)$
 - Further regularization tools.

Let's try to marginalize out the top layer only:

$$p(\mathbf{h}_2) = \int p(\mathbf{h}_2 | \mathbf{h}_1) p(\mathbf{h}_1) d\mathbf{h}_1$$

=
$$\int \int p(\mathbf{h}_2 | \mathbf{f}_2) p(\mathbf{f}_2 | \mathbf{h}_1) p(\mathbf{h}_1) d\mathbf{f}_2 \mathbf{h}_1$$

=
$$\int p(\mathbf{h}_2 | \mathbf{f}_2) \Big[\underbrace{\int p(\mathbf{f}_2 | \mathbf{h}_1) p(\mathbf{h}_1) d\mathbf{h}_1}_{\text{Intractable!}} \Big] d\mathbf{f}_2$$

Intractability: \mathbf{h}_1 appears non-linearly in $p(\mathbf{f}_2|\mathbf{h}_1)$, inside \mathbf{K}^{-1} (and also the determinant term), where $\mathbf{K} = k(\mathbf{h}_1, \mathbf{h}_1)$.

- Similar issues arise for 1-layer models. Solution was given by Titsias and Lawrence, 2010. A small modification to that solution does the trick in deep GPs too.
- Extend Titsias' method for variational learning of inducing variables in Sparse GPs.
- Analytic variational bound $\mathcal{F} \leq p(y|x)$
- Approximately marginalise out h
- Hence obtain the approximate posterior q(h)

Inducing points: sparseness, tractability and Big Data

h



Inducing points: sparseness, tractability and Big Data



Inducing points: sparseness, tractability and Big Data



- Inducing points originally introduced for faster (sparse) GPs
- But this also induces tractability in our models, due to the conditional independencies assumed
- ► Viewing them as global variables ⇒ extension to Big Data [Hensman et al., UAI 2013]

Bayesian regularization

