

Gaussian Processes for Regression, Classification and Dimensionality Reduction

Advanced Machine Learning
École Polytechnique Fédérale de Lausanne, Switzerland



6 May 2020

Learning Outcomes

- Probabilistic Linear Regression
- Maximum Likelihood and Maximum a posteriori estimations
- Bayesian Model Selection
- Non-linear regression with GPs
- Classification with GPs
- Dimensionality reduction with GPLVM

Probabilistic Linear Regression

Creating the model

Observations y are modeled as linear combination of weights w and inputs x contaminated with noise ϵ

$$\hat{y} = f(x) = w^T x + \epsilon \quad (1)$$

The probabilistic regression model is created by handling noise as random variable (RV):

- $\epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$

Gaussian properties:

- The addition of a Normal RV to a constant variable yields a Normal random variable

$$p(\hat{y}) = \mathcal{N}(\hat{w}^T x, \sigma_\epsilon^2) \quad (2)$$

Probabilistic Linear Regression

Maximum Likelihood Estimation

Given training set: $\{\mathbf{y}, \mathbf{X}\} = \{y_i, \mathbf{x}_i\}_{i=1}^M$ and σ_ϵ is assumed known (Hyperparameter), we can define the likelihood as:

$$p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{w}}, \sigma_\epsilon^2) = \prod_{i=1}^M \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{y_i - \hat{\mathbf{w}}^T \mathbf{x}_i}{2\sigma_\epsilon^2}\right) \quad (3)$$

$$\hat{\mathbf{w}}_{\text{MLE}}, \sigma_\epsilon^2 = \arg \max_{\hat{\mathbf{w}}, \sigma_\epsilon^2} p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{w}}, \sigma_\epsilon^2 | M)$$

Closed form solution

$$\hat{\mathbf{w}}_{\text{MLE}} = (\mathbf{X}\mathbf{X}^T)^{-1} \mathbf{X}\mathbf{y} \rightarrow \text{Same with the OLS estimator}$$

$$\hat{\sigma}_{\text{MLE}}^2 = \frac{1}{N} \left(\hat{\mathbf{w}}_{\text{MLE}}^T \mathbf{X} - \mathbf{y} \right)^T \left(\hat{\mathbf{w}}_{\text{MLE}}^T \mathbf{X} - \mathbf{y} \right) \rightarrow \text{Average of squared deviations}$$

Probabilistic Linear Regression

Maximum a Posteriori (MAP) estimation

Weights are also handled as RVs

Prior distribution over weights: $\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{w}}, \boldsymbol{\Sigma}_{\mathbf{w}})$

Model Likelihood: $p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{w}}, \sigma_{\epsilon}^2 \mathbf{I}_M)$

Bayes Rule: $\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}$

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{w}}, \sigma_{\epsilon}^2) p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})} \quad (4)$$

Prior and Likelihood are conjugate distributions¹. The posterior has a closed form solution

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{w}|\mathbf{X}, \mathbf{y}}, \mathbf{A}^{-1})$$

¹See Conjugate Bayesian analysis of the Gaussian distribution by Murphy

Probabilistic Linear Regression

Maximum a Posteriori (MAP) estimation

MAP estimation derives from the expected value of the posterior normal distribution

$$\begin{aligned}\hat{\mathbf{w}}_{\text{MAP}} &= \mathbb{E} \{p(\mathbf{w}|\mathbf{x}, \mathbf{Y})\} \\ \hat{\Sigma}_{\text{MAP}}^{-1} &= \mathbf{A} = \frac{1}{\sigma_{\epsilon}^2} \mathbf{X}\mathbf{X}^T + \Sigma_{\mathbf{w}}^{-1} \\ \hat{\mathbf{w}}_{\text{MAP}} &= \mu_{\mathbf{w}|\mathbf{x}, \mathbf{y}} = \mathbf{A}^{-1} \Sigma_{\mathbf{w}}^{-1} \mu_{\mathbf{w}} + \frac{1}{\sigma_{\epsilon}^2} \mathbf{A}^{-1} \mathbf{X}\mathbf{y}\end{aligned}\tag{5}$$

For the special of $\mu_{\mathbf{w}} = \mathbf{0}$ and $\Sigma_{\mathbf{w}} = \tau \mathbf{I}$ the Bayesian regression reduces to ridge regression with $\lambda = \frac{\sigma_{\epsilon}^2}{\tau}$

Probabilistic Linear Regression

Posterior Predictive Distribution

The predictive distribution is a distribution over output y and allows to make predictions given a set of testing data \mathbf{x}_*

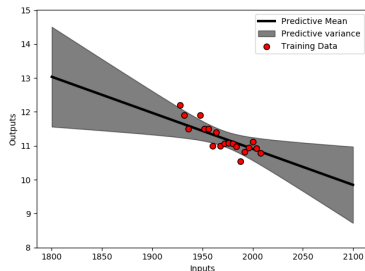
$$p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \underbrace{p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{w}}, \sigma_\epsilon^2 \mathbf{I}_M)}_{\text{Likelihood}} \underbrace{p(\mathbf{w}|\mathbf{X}, \mathbf{y})}_{\text{Posterior}} d\mathbf{w} \quad (6)$$

The integral has a closed form solution in the case of Normal distributions

$$p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) \sim \mathcal{N}(\mu_{y^*}, \sigma_{y^*}^2)$$
$$\mu_{y^*} = \mathbf{x}_*^T \hat{\mathbf{w}}_{\text{MAP}} = \frac{1}{\sigma_\epsilon^2} \mathbf{x}_*^T \mathbf{A}^{-1} \mathbf{X} \mathbf{y} \quad \swarrow \text{assuming } \mu_{\mathbf{w}} = \mathbf{0}$$
$$\sigma_{y^*}^2 = \mathbf{x}_*^T \mathbf{A}^{-1} \mathbf{x}_*$$
(7)

Probabilistic Linear Regression

Posterior Predictive Distribution



$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \sim \mathcal{N}(\mu_{y_*}, \sigma_{y_*}^2)$$

$$\mu_{y_*}(\mathbf{x}_*) = \mathbf{x}_*^T \hat{\mathbf{w}}_{\text{MAP}} = \frac{1}{\sigma_\epsilon^2} \mathbf{x}_*^T \mathbf{A}^{-1} \mathbf{X} \mathbf{y}$$

$$\sigma_{y_*}^2(\mathbf{x}_*) = \mathbf{x}_*^T \mathbf{A}^{-1} \mathbf{x}_*$$

The variance gives a measure of the uncertainty of the prediction

Probabilistic Linear Regression

Marginal Likelihood

Bayesian Regression depends on hyper-parameters

The marginal likelihood provides a metric for optimization

Bayes Rule:

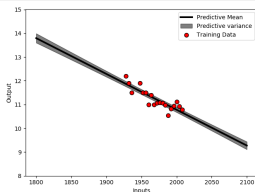
$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{w}}, \sigma_\epsilon^2)p(\mathbf{w})}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{Marginal Likelihood}}}$$

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, \hat{\mathbf{w}}, \sigma_\epsilon^2)p(\mathbf{w}) d\mathbf{w} \quad (8)$$

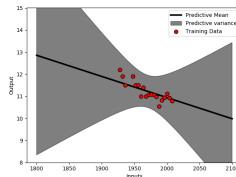
Measures how well our model explains the observed data. For the current case, it has a closed form solution

Probabilistic Linear Regression

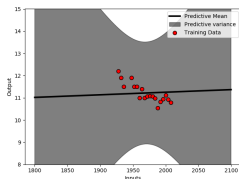
Marginal Likelihood



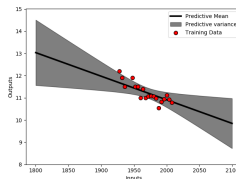
(a) 33.38



(b) 47.19



(c) 67.19



(d) 27.32

Figure: Negative Log-Marginal Likelihood for different models

Gaussian Process for Regression

From probabilistic linear to nonlinear regression

Starting with the predictive distribution of Bayesian linear regression, Apply non-linear mapping $\phi(\mathbf{x})$ to the feature space.

$$\begin{aligned}\mu_{y^*}(\phi_*) &= \frac{1}{\sigma_\epsilon^2} \phi_*^T \mathbf{A}^{-1} \Phi \mathbf{y} \\ p(y_* | \phi_*, \Phi, \mathbf{y}) &\sim \mathcal{N}(\mu_{y^*}, \sigma_{y^*}^2) \quad \sigma_{y^*}^2(\phi_*) = \phi_*^T \mathbf{A}^{-1} \phi_* \\ \mathbf{A} &= \frac{1}{\sigma_\epsilon^2} \Phi \Phi^T + \Sigma_{\mathbf{w}}^{-1}\end{aligned}$$

The $p(y_* | \phi_*, \Phi, \mathbf{y})$ is a Normal distribution over functions.

$$\left. \begin{array}{l} \phi_*^T \mathbf{A}^{-1} \Phi \\ \phi_*^T \mathbf{A}^{-1} \phi_* \end{array} \right\} \text{Inner product in feature space}$$

Gaussian Process for Regression

From probabilistic linear to nonlinear regression

Define kernel as: $k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}(\mathbf{x}) \boldsymbol{\Sigma}_{\mathbf{w}}^{-1} \boldsymbol{\phi}(\mathbf{x}')$ and apply on the mean and variance²

$$\mathbb{E} \{p(y_* | \boldsymbol{\phi}_*, \boldsymbol{\Phi}, \mathbf{y})\} = \sum_{i=1}^M \alpha_i k(\mathbf{x}_i, \mathbf{x}_*)$$

$$\text{where: } \mathbf{a} = [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_y^2 \mathbf{I}]^{-1} \mathbf{y}$$

$$\text{Var} \{p(y_* | \boldsymbol{\phi}_*, \boldsymbol{\Phi}, \mathbf{y})\} = k(\mathbf{x}_*, \mathbf{x}_*) + \mathbf{k}(\mathbf{x}_*, \mathbf{X}) [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_y^2 \mathbf{I}]^{-1} \mathbf{k}(\mathbf{X}, \mathbf{x}_*)$$

In GPs all the training data are used for predictions!

²See supplement material for steps

Gaussian Process for Regression

Kernels

Kernel functions:

- Linear

$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}' + c$ where c the intercept \equiv Bayes regression

- Gaussian (RBF)

$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\ell^2}\right)$ where ℓ the kernel width (lengthscale)

- Polynomial

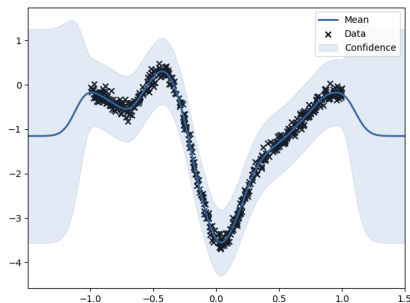
$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + c)^d$ where d is the degree of polynomial

- Periodic

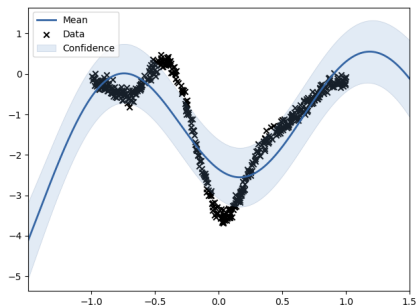
$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{2 \sin^2(\pi \|\mathbf{x} - \mathbf{x}'\| / p)}{\ell^2}\right)$ si relation périodique
where p is the period.

Gaussian Process for Regression

Kernels – RBF



(a) $\ell = 0.1$

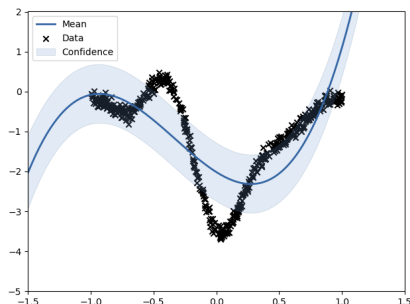


(b) $\ell = 1$

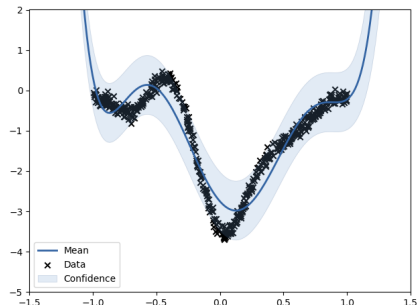
Figure: Impact of the width of RBF kernel

Gaussian Process for Regression

Kernels – Polynomial



(a) $d = 3$

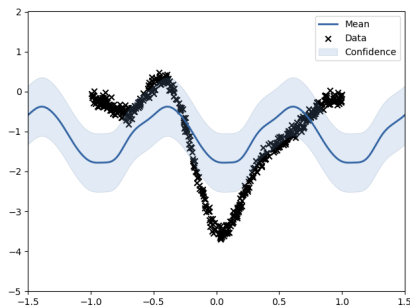


(b) $d = 6$

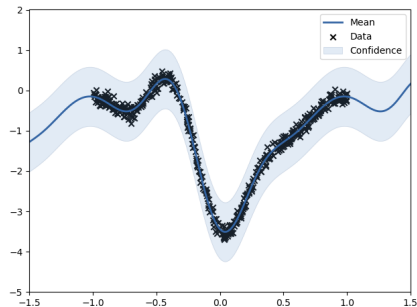
Figure: Impact of the degree of polynomial kernel

Gaussian Process for Regression

Kernels – Periodic



(a) $p = 1$



(b) $p = 2$

Figure: Impact of the period of periodic kernel

Gaussian Process for Regression

Construction of new kernels

Not all kernels are appropriate for the all datasets.

There exists many kernels in literature.

As seen in kernel lecture, new kernels can be created by combination of kernels:

valid kernel: symmetric and positive definite

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \text{ where } c > 0$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \text{ where } f(\cdot) \text{ any function}$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') k_2(\mathbf{x}, \mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}'))$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}_a, \mathbf{x}'_a) + k_2(\mathbf{x}_b, \mathbf{x}'_b) \text{ where } a, b \text{ dimensions of } \mathbf{x}$$

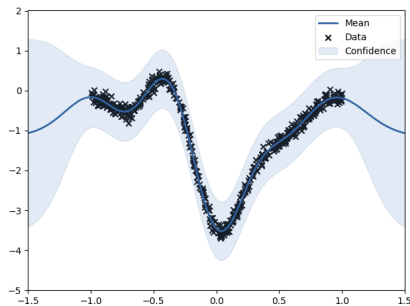
$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}_a, \mathbf{x}'_a) k_2(\mathbf{x}_b, \mathbf{x}'_b)$$

↑ peut avoir diff kernel sur diff dimensions

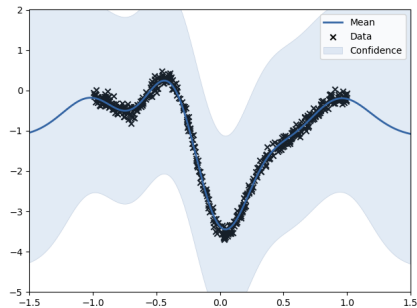
Gaussian Process for Regression

The noise

The larger the noise, the more uncertainty



(a) $\sigma_y^2 = 0.1$



(b) $\sigma_y^2 = 1$

Figure: Impact of noise variance at the confidence interval

Gaussian Process for Regression

Hyper-parameter tuning

Hyper-parameters: kernel's parameters and noise variance

Tuning:

- Cross-validation
- Minimize Marginal-likelihood (*local optima*)
s.t. hyper-parameters, where $\mathbf{K}_{\sigma_y^2} = [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_y^2 \mathbf{I}]$:

$$-\log(y|\mathbf{X}) = 0.5 \underbrace{(\mathbf{y}^T \mathbf{K}_{\sigma_y^2}^{-1} \mathbf{y})}_{\text{Fit}} + \underbrace{\log|\mathbf{K}_{\sigma_y^2}|}_{\text{Complexity}} + \log 2\pi \quad (9)$$

*ex RBF σ petit good fit
max $|\mathbf{K}| \uparrow$*

Automatically provides trade-off between model fit and complexity³

³See supplement material

Gaussian Process for Regression

Zero Mean Prior

Prior expectations

Linear Model:

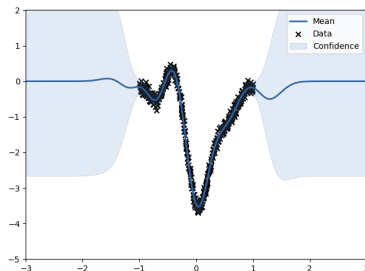
$$\begin{aligned}\mathbb{E}\{y\} &= \mathbb{E}\left\{w^T x + \mathcal{N}(0, \sigma_y^2)\right\} \\ &= \mathbb{E}\left\{w^T\right\} x + 0\end{aligned}$$

Non-Linear Model:

$$\begin{aligned}\mathbb{E}\{y\} &= \mathbb{E}\left\{w^T\right\} \phi(x) + \mathcal{N}(0, \sigma_y^2) \\ &= \mathbb{E}\left\{w^T\right\} \phi(x) + 0\end{aligned}$$

Both models have a zero mean Gaussian prior.

GPs with RBF kernel, predict zero far from the data



Gaussian Process for Regression

Mean function

Instead of zero mean one can use a mean function $m(\cdot) : \mathbf{x} \mapsto y$.

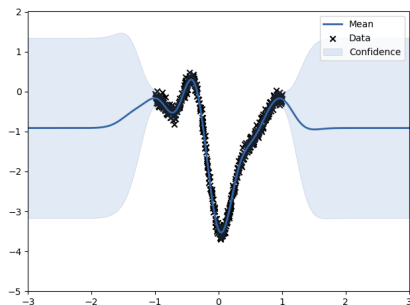
$$\mathbb{E} \{p(y_* | \phi_*, \Phi, \mathbf{y})\} = \sum_{i=1}^M m(\mathbf{x}_*) + \alpha_i k(\mathbf{x}_i, \mathbf{x}_*) \quad (10)$$

$$\text{where: } \mathbf{a} = [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_y^2 \mathbf{I}]^{-1} (\mathbf{y} - m(\mathbf{x}_*))$$

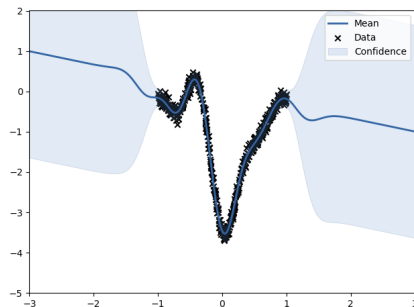
The parameters of the mean can be auto-tuned based on the marginal likelihood

Gaussian Process for Regression

Impact of mean function



(a) Constant mean



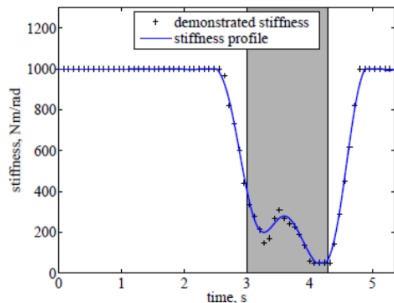
(b) Linear Mean

Figure: GPs with different mean functions

Gaussian Process for Regression

Applications – Learning from demonstration

Learning Compliant Manipulation through Kinesthetic and Tactile Human-Robot Interaction⁴



The stiffness profile is encoded as a time-varying input using GPR. The shaded area corresponds to the striking phase.

⁴ Kronander, K. and Billard, A. (2013) Learning Compliant Manipulation through Kinesthetic and Tactile Human-Robot Interaction. IEEE Transactions on Haptics. 10.1109/TOH.2013.54.

Gaussian Process for Regression

Applications – Surface models

GPR can be used to model the shape of objects⁵

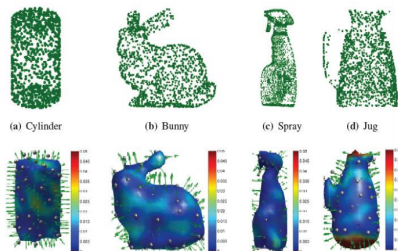


Figure: Top: 3D points sampled either from a camera or from tactile sensing. Bottom: 3D shape reconstructed by GPR. The arrows represent the predicted normals at the surface

⁵ El Khoury, S., Li, M., and Billard, A. (2013). On the generation of a variety of grasps. *Robotics and Autonomous Systems*, 61(12):1335–34

Gaussian Process for Classification

GPs from the Bayesian perspective

Given training data \mathbf{X}, \mathbf{y} , infer continuous function f such as:

$$y = f(\mathbf{x}) + \epsilon, \text{ where: } \epsilon \sim \mathcal{N}(0, \sigma_y^2)$$

GPs are jointly normal distributions over the outputs of the latent function $p(f) \sim \mathcal{GP}(m(x), k(x, x'))$.

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_M \end{bmatrix} = \mathcal{N} \left(\begin{bmatrix} m(x_1) \\ m(x_2) \\ \vdots \\ m(x_M) \end{bmatrix}, \begin{bmatrix} k(x_1, x_1) + \sigma_y^2 & k(x_1, x_2) & \dots & k(x_1, x_M) \\ k(x_2, x_1) & k(x_2, x_2) + \sigma_y^2 & \dots & k(x_2, x_M) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_M, x_1) & k(x_M, x_2) & \dots & k(x_M, x_M) + \sigma_y^2 \end{bmatrix} \right)$$

The prediction for f_* for a novel input x_* can also be made by conditioning at the joint Normal $p(f_* | \mathbf{f}, \mathbf{X}, \mathbf{y})$

Gaussian Process for Classification

GPs from the Bayesian perspective

Given training data \mathbf{X}, \mathbf{y} , infer a continuous function f such as:

$$y = f(\mathbf{x}) + \epsilon, \text{ where: } \epsilon \sim \mathcal{N}(0, \sigma_y^2)$$

GPs are probability distributions over a latent function $p(f)$.

Bayes rule for GPs:

$$p(f|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{f}, \mathbf{X}) p(\mathbf{f})}{p(\mathbf{y}|\mathbf{X})}$$

Prior : $p(\mathbf{f}) \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$

Likelihood: $p(\mathbf{y}|\mathbf{f}, \mathbf{X})$ Normal for regression

Posterior: $p(\mathbf{f}|\mathbf{X}, \mathbf{y}) \sim \mathcal{GP}(m_{\text{post}}, k_{\text{post}}(\mathbf{x}, \mathbf{x}'))$

For Normal likelihood and a GP prior the posterior has a closed form solution

Gaussian Process for Classification

GPs from the Bayesian perspective

Given data $\mathbf{X}, \mathbf{y} \in [0; 1]$, for a binary classification problem, infer a continuous function f such as:

$$y = \lambda(f(\mathbf{x}) + \epsilon) \text{ , where: } \epsilon \sim \mathcal{N}(0, \sigma_y^2) \text{ and } \lambda(f(\mathbf{x}) + \epsilon) = \frac{1}{1 + \exp(y_i f_i + \epsilon)}$$

$$p(\mathbf{f} | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{f}, \mathbf{X}) p(\mathbf{f})}{p(\mathbf{y} | \mathbf{X})}$$

Prior : $p(\mathbf{f}) \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'))$

Likelihood: $p(\mathbf{y} | \mathbf{f}, \mathbf{X}) \sim \lambda(y_i f_i + \epsilon)$ if $y_i = 1$ otherwise $1 - \lambda(y_i f_i + \epsilon)$

Posterior: $p(\mathbf{f} | \mathbf{X}, \mathbf{y})$ No closed form solution!!!

The logistic likelihood and the GP prior are not conjugate

Gaussian Process for Classification

Approximation Methods

Many methods exist to approximate the posterior:

- Laplace approximation
Taylor series expansion about mode of posterior
- Expectation propagation
Minimize divergence between approximated and true distribution
- Markov Chain Monte Carlo
Take a large amount of samples from posterior (more computationally expensive)

Gaussian Process for Classification

Laplace Approximation

Approximate as a Normal distribution:

$$\tilde{p}(\mathbf{f}|\mathbf{X}, \mathbf{y}) \sim \mathcal{N}(\tilde{\mathbf{f}}, \mathbf{C}^{-1})$$

where: $\tilde{\mathbf{f}} = \arg \max_{\mathbf{f}} p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ (11)

$$\mathbf{C} = -\nabla^2 \log p(\mathbf{f}|\mathbf{X}, \mathbf{y})|_{\mathbf{f}=\tilde{\mathbf{f}}}$$

\uparrow
curvature

Gaussian Process for Classification

Laplace approximation of the posterior

Find $\tilde{\mathbf{f}}$ by maximizing log-posterior :

$$\log p(\mathbf{f}|\mathbf{X}, \mathbf{y}) = \log p(\mathbf{y}|\mathbf{f}, \mathbf{X}) + \log p(\mathbf{f})$$

Differentiating w.r.t \mathbf{f}

$$\begin{aligned}\nabla \log p(\mathbf{f}|\mathbf{X}, \mathbf{y}) &= \nabla \log p(\mathbf{y}|\mathbf{f}, \mathbf{X}) - \mathbf{K}^{-1}\mathbf{f} \\ \nabla^2 \log p(\mathbf{f}|\mathbf{X}, \mathbf{y}) &= \nabla^2 \log p(\mathbf{y}|\mathbf{f}, \mathbf{X}) - \mathbf{K}^{-1} \\ &= -\mathbf{D} - \mathbf{K}^{-1}\end{aligned}$$

$$\tilde{p}(\mathbf{f}|\mathbf{X}, \mathbf{y}) \sim \mathcal{N}\left(\tilde{\mathbf{f}}, (\mathbf{D} + \mathbf{K}^{-1})^{-1}\right) \quad (12)$$

Gaussian Process for Classification

Making predictions

Posterior predictive distribution:

$$\tilde{p}(\mathbf{f}_*|\mathbf{x}_*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{f}_*|\mathbf{x}_*, \mathbf{X}, \mathbf{Y}, \mathbf{f}) \tilde{p}(\mathbf{f}|\mathbf{X}, \mathbf{Y}) d\mathbf{f}$$

$$\mathbb{E}_{\tilde{p}}[\mathbf{f}_*|\mathbf{X}, \mathbf{Y}, \mathbf{x}_*] = \sum_{i=1}^M \alpha_i \mathbf{k}(\mathbf{x}_*, \mathbf{x}_i), \text{ where: } \mathbf{a} = [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_y^2 \mathbf{I}]^{-1} \tilde{\mathbf{f}}$$

$$\text{Var}_{\tilde{p}}[\mathbf{f}_*|\mathbf{X}, \mathbf{Y}, \mathbf{x}_*] = k(\mathbf{x}_* \mathbf{x}_*) - \mathbf{k}(\mathbf{X} \mathbf{x}_*)^T (\mathbf{D} + \mathbf{K}^{-1})^{-1} \mathbf{k}(\mathbf{X}, \mathbf{x}_*)$$

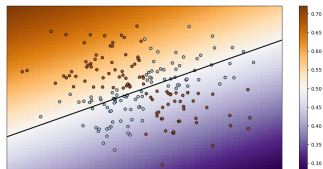
Class predictions:

$$p(y_* = 1|\mathbf{x}_*, \mathbf{X}_*, \mathbf{Y}) = \int p(y_* = 1|\mathbf{f}_*) \tilde{p}(\mathbf{f}_*|\mathbf{x}_*, \mathbf{X}, \mathbf{Y}) d\mathbf{f}_*$$

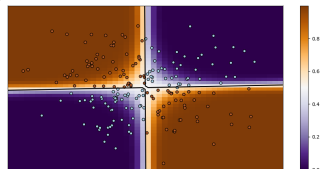
No closed form solution, has to be approximated

Gaussian Process for Classification

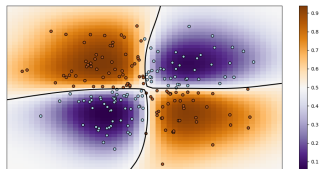
Toy Examples



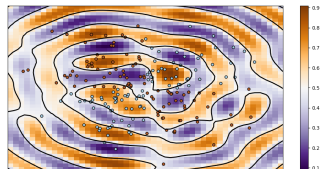
(a) Linear



(b) Second degree polynomial



(c) RBF

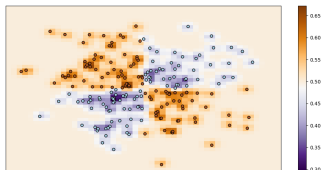


(d) Periodic

Figure: GP classification of a XOR dataset for different kernels

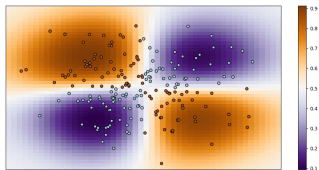
Gaussian Process for Classification

Impact of kernel width

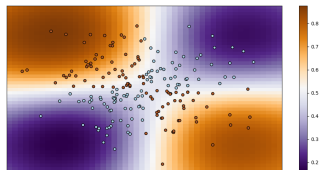


kernel on every points

(a) width = 0.1



(b) width = 1



(c) width = 2

Figure: GP classification of a XOR dataset using RBF kernel with various widths

GPLVM

Outline

- GPLVM is an unsupervised variation of GPs
- Probabilistic dimensionality reduction

Derivation:

- PCA to Probabilistic PCA
- Dual Formulation of Probabilistic PCA
- GPLVM through kernel trick

GPLVM

From PCA to Probabilistic PCA

M data $\mathbf{Y} \in \mathbb{R}^d$ at the original space and \mathbf{X} their projection to the latent space.

PCA: $\mathbf{y} = \mathbf{A}\mathbf{x}$

Probabilistic PCA

- Introduce stochasticity by adding noise variable $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$
 $\mathbf{y} = \mathbf{A}\mathbf{x} + \epsilon$

- Likelihood:

$$p(\mathbf{y}|\mathbf{A}, \mathbf{x}) \sim \mathcal{N}(\mathbf{A}\mathbf{x}, \sigma^2 \mathbf{I})$$

- Integrate over \mathbf{x} to get marginal likelihood

$$p(\mathbf{y}|\mathbf{A}) = \int p(\mathbf{y}|\mathbf{A}, \mathbf{x}) p(\mathbf{x})$$

- Requires prior over \mathbf{x} :

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

← un informative et besoin GP

1/

GPLVM

From PCA to Probabilistic PCA

Marginal likelihood has a closed form solution:

$$p(\mathbf{y}|\mathbf{A}) = \mathcal{N}(\mathbf{0}, \mathbf{A}\mathbf{A}^T + \sigma^2\mathbf{X})$$

Assuming \mathbf{Y} are i.i.d:

$$p(\mathbf{Y}|\mathbf{A}) = \prod_{i=1}^M p(\mathbf{y}_i|\mathbf{A})$$

Maximizing marginal likelihood wrt to \mathbf{A} yields unique optimal solution⁶

⁶Tipping and Bishop. "Probabilistic principal component analysis."

GPLVM

Dual Probabilistic PCA

Probabilistic PCA: Integrate over \mathbf{x} and optimize w.r.t to \mathbf{A}

$$p(\mathbf{y}|\mathbf{A}) = \int p(\mathbf{y}|\mathbf{A}, \mathbf{x}) p(\mathbf{x})$$

Dual Probabilistic PCA Integrate over \mathbf{A} and optimize w.r.t to \mathbf{x}

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}|\mathbf{A}, \mathbf{x}) p(\mathbf{A})$$

Prior over \mathbf{A} :

$$p(\mathbf{A}) = \prod_{d=1}^D \mathcal{N}(\mathbf{a}_d; \mathbf{0}, \mathbf{I})$$

GPLVM

Dual Probabilistic PCA

Marginal likelihood:

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

GPLVM:

- $\mathbf{X}^T \mathbf{X} = \mathbf{K}(\mathbf{x}, \mathbf{x}')$ Inner product allows to apply kernel trick
- $\mathcal{N}(\mathbf{y}_d; \mathbf{0}, \mathbf{K}(\mathbf{x}, \mathbf{x}') + \sigma^2 \mathbf{I})$ A Gaussian process on each dimension of \mathbf{Y}
- Minimize negative log-marginal likelihood w.r.t \mathbf{X}

$$-\log p(\mathbf{y}_d|\mathbf{X}) = \frac{1}{2} \left(DN \log 2\pi + D \log |\mathbf{K}| + \text{tr}(\mathbf{K}^{-1} \mathbf{Y}^T \mathbf{Y}) \right)$$

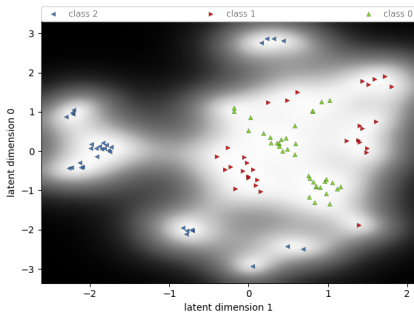
GPLVM

Comparison with kPCA

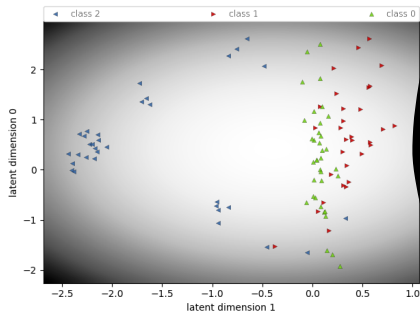
Method	Proximity	Mapping	Non-Linear	Probabilistic
GPLVM	YES	$\mathbf{X} \rightarrow \mathbf{Y}$	YES	YES
kPCA	YES	$\mathbf{Y} \rightarrow \mathbf{X}$	YES	NO

GPLVM

Examples



(a) RBF kernel



(b) Linear Kernel

Figure: GPLVM with different kernels

2/

GPLVM

Choosing latent dimensions

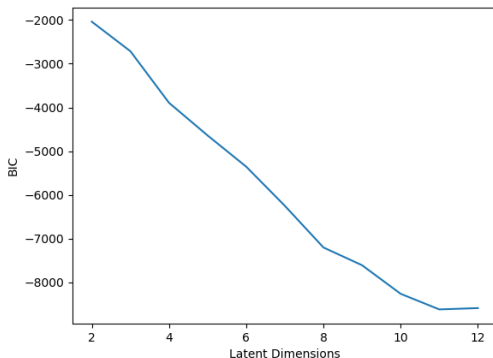


Figure: BIC as function of the number of latent dimensions

$$\text{BIC} = M \log(k) - 2 \log(p(\mathbf{Y}|\mathbf{X})) \quad (13)$$

GPLVM

Choosing latent dimensions

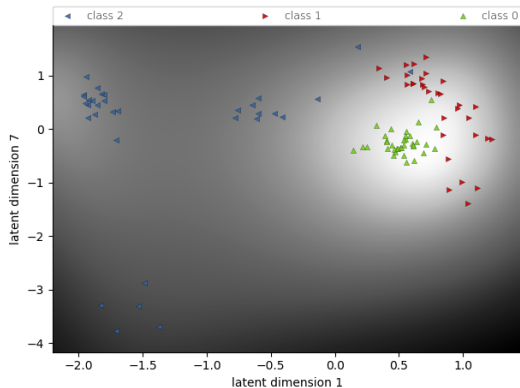


Figure: Combination of latent dimensions which increase separability

Gaussian Processes

Discussion

- Advantages
 - Generalization
 - Accuracy
 - Metric of Uncertainty
 - Auto-tuning of hyper-parameters
- Disadvantages
 - Computational complexity
- Recommended packages
 - GPy
 - scikit-learn