

Machine Learning 2

2. Gaussian Process Models (GPs)

Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL)
Institute for Computer Science
University of Hildesheim, Germany

Syllabus

		A. Advanced Supervised Learning
Fri. 12.4.	(1)	A.1 Generalized Linear Models
Fri. 26.4.	(2)	A.2 Gaussian Processes
Fri. 3.5.	(3)	A.2b Gaussian Processes (ctd.)
Fri. 10.5.	(4)	A.3 Advanced Support Vector Machines
		B. Ensembles
Fri. 17.5.	(5)	B.1 Stacking
Fri. 24.5.	(6)	B.2 Boosting
Fri. 31.5.	(7)	B.3 Mixtures of Experts
		C. Sparse Models
Fri. 7.6.	(8)	C.1 Homotopy and Least Angle Regression
Fri. 14.6.	—	— Pentecoste Break —
Fri. 21.6.	(9)	C.2 Proximal Gradients
Fri. 28.6.	(10)	C.3 Laplace Priors
Fri. 29.6.	(11)	C.4 Automatic Relevance Determination
		D. Complex Predictors
Fri. 6.7.	(12)	D.1 Latent Dirichlet Allocation (LDA)
Fri. 12.7.	(13)	Q & A

Outline

1. GPs for Regression
2. GPs for Classification

Outline

1. GPs for Regression

2. GPs for Classification

Gaussian Process Model

Gaussian Processes describe

- ▶ the **vector** $y := (y_1, \dots, y_N)^T$ of all targets
- ▶ as a sample from a **normal distribution**
- ▶ where targets of different instances are **correlated by a kernel** Σ :
- ▶ and thus depend on the **matrix** X of all predictors:

$$y \mid X \sim \mathcal{N}(y \mid \mu(X), \Sigma(X))$$

with

$$\mu(X)_n := m(x_n)$$

$$\Sigma(X)_{n,m} := k(x_n, x_m), \quad n, m \in \{1, \dots, N\}$$

with a **kernel function** k and **mean function** m (often $m = 0$).

Kernels

The kernel k measures how much targets y, y' correlate given their predictors x, x' .

- ▶ $k(x, x')$ is larger the more similar x, x' are
- ▶ esp. $k(x, x) \geq k(x, x') \forall x, x'$

Example: **squared exponential kernel** / **Gaussian kernel**

$$k(x, x') := \sigma_f^2 e^{-\frac{1}{2\ell^2} \|x-x'\|^2}$$

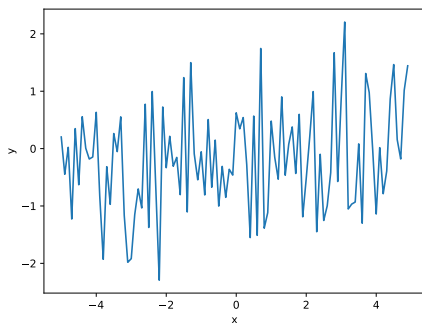
with **kernel (hyper)parameters**

ℓ horizontal length scale (x)

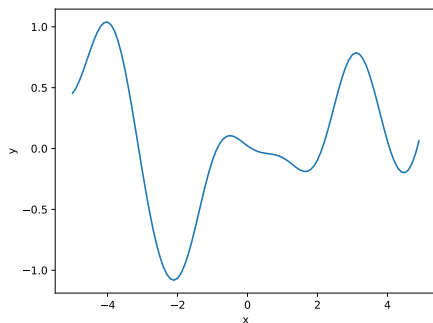
σ_f^2 vertical variation (y)

GPs as Prior on Functions

identity kernel

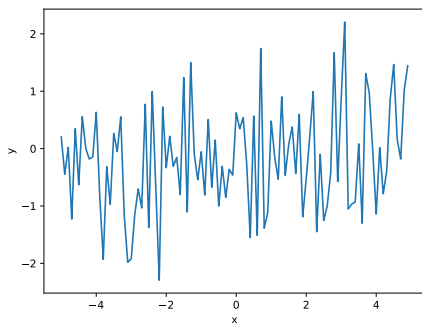


squared exponential kernel

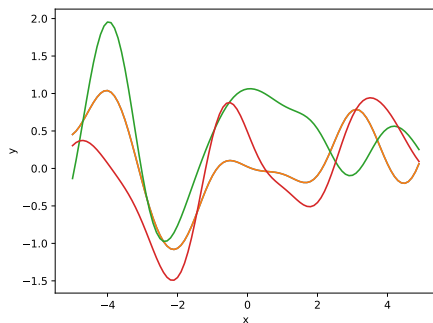


GPs as Prior on Functions

identity kernel



squared exponential kernel



Conditional Distributions of Multivariate Normals

Let y_A, y_B be jointly Gaussian

$$y := \begin{pmatrix} y_A \\ y_B \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} y_A \\ y_B \end{pmatrix} \mid \begin{pmatrix} \mu_A \\ \mu_B \end{pmatrix}, \begin{pmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{pmatrix} \right)$$

then the **conditional distribution** is

$$p(y_B \mid y_A) = \mathcal{N}(y_B \mid \mu_{B|A}, \Sigma_{B|A})$$

with

$$\mu_{B|A} := \mu_B + \Sigma_{BA} \Sigma_{AA}^{-1} (y_A - \mu_A)$$

$$\Sigma_{B|A} := \Sigma_{BB} - \Sigma_{BA} \Sigma_{AA}^{-1} \Sigma_{AB}$$

Predictions w/o Noise

Let y, X be the training data,
 X_* be the test data and
 y_* be the test targets to predict.

$$\begin{pmatrix} y \\ y_* \end{pmatrix} | X, X_* \sim \mathcal{N}\left(\begin{pmatrix} y \\ y_* \end{pmatrix} | \begin{pmatrix} \mu \\ \mu_* \end{pmatrix}, \begin{pmatrix} \Sigma & \Sigma_* \\ \Sigma_*^T & \Sigma_{**} \end{pmatrix}\right)$$

with

$$\begin{aligned} \mu &:= m(X), & \mu_* &:= m(X_*) \\ \Sigma &:= k(X, X), & \Sigma_* &:= k(X, X_*), & \Sigma_{**} &:= k(X_*, X_*) \end{aligned}$$

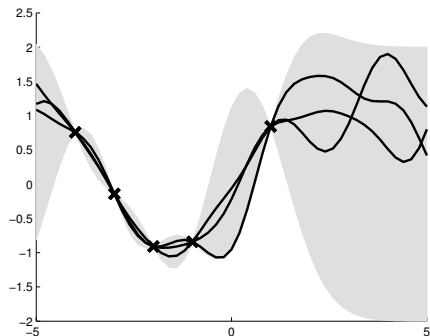
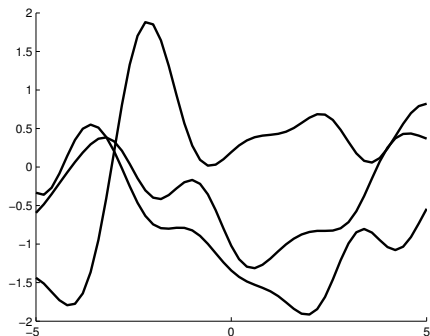
Then

$$p(y_* | y) = \mathcal{N}(y_* | \tilde{\mu}_*, \tilde{\Sigma}_*)$$

$$\tilde{\mu}_* := \mu_* + \Sigma_*^T \Sigma^{-1} (y - \mu)$$

$$\tilde{\Sigma}_* := \Sigma_{**} - \Sigma_*^T \Sigma^{-1} \Sigma_*$$

Example w/o Noise



Without noise the data is interpolated.

[Mur12, fig. 15.2]

Predictions with Noise

No noise:

$$\Sigma := K$$

With noise:

$$\Sigma := K_y = K + \sigma_y^2 I$$

Then as before

$$p(y_* | y) = \mathcal{N}(y_* | \tilde{\mu}_*, \tilde{\Sigma}_*)$$

now with

$$\begin{aligned}\tilde{\mu}_* &:= \mu_* + K_*^T K_y^{-1} (y - \mu) \\ \tilde{\Sigma}_* &:= K_{**} + \sigma_y^2 I - K_*^T K_y^{-1} K_*\end{aligned}$$

where

$$K := k(X, X), \quad K_* := k(X, X_*), \quad K_{**} := k(X_*, X_*)$$

Predictions with Noise, Zero Means

$$p(y_* | y) = \mathcal{N}(y_* | \tilde{\mu}_*, \tilde{\Sigma}_*)$$

with

$$\begin{aligned}\tilde{\mu}_* &:= \mu_* + K_*^T K_y^{-1} (y - \mu) \\ \tilde{\Sigma}_* &:= K_{**} + \sigma_y^2 I - K_*^T K_y^{-1} K_*\end{aligned}$$

With $m = 0$:

$$p(y_* | y) = \mathcal{N}(y_* | \tilde{\mu}_*, \tilde{\Sigma}_*)$$

with

$$\begin{aligned}\tilde{\mu}_* &:= K_*^T K_y^{-1} y \\ \tilde{\Sigma}_* &:= K_{**} + \sigma_y^2 I - K_*^T K_y^{-1} K_*\end{aligned}$$

Prediction for a single instance

$$p(y_* | y) = \mathcal{N}(y_* | \tilde{\mu}_*, \tilde{\Sigma}_*)$$

with

$$\tilde{\mu}_* := K_*^T K_y^{-1} y$$

$$\tilde{\Sigma}_* := K_{**} + \sigma_y^2 I - K_*^T K_y^{-1} K_*$$

Prediction \hat{y} for a single instance x :

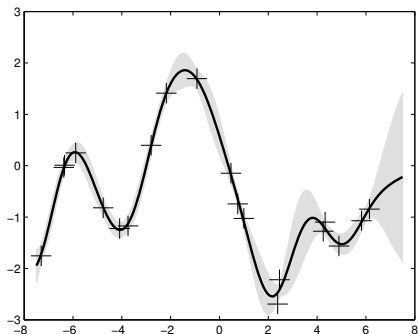
$$\hat{y}(x) := k_*^T K_y^{-1} y = \sum_{n=1}^N \alpha_n k(x_n, x), \quad \alpha := K_y^{-1} y$$

with

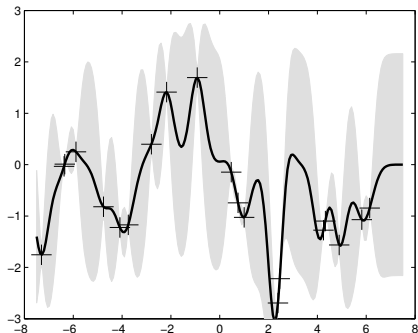
$$k_* := k(X, x)$$

But GPs can provide a joint inference for multiple instances

Example with Noise



$$(\ell, \sigma_f, \sigma_y) = (1, 1, 0.1)$$

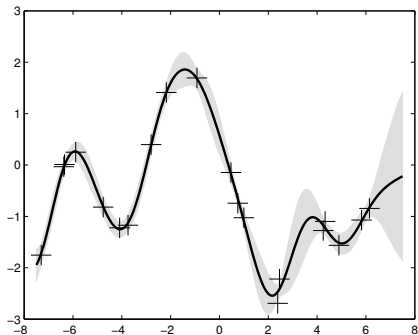


$$(\ell, \sigma_f, \sigma_y) = (0.3, 0.1?, 0.00005)$$

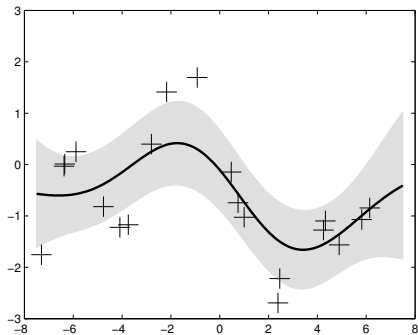
[Mur12, fig. 15.3]



Example with Noise



$$(\ell, \sigma_f, \sigma_y) = (1, 1, 0.1)$$



$$(\ell, \sigma_f, \sigma_y) = (3, 1.16, 0.89)$$

[Mur12, fig. 15.3]



Estimating Kernel Parameters

Either treating them as hyperparameters (grid search, random search) or maximize the marginal likelihood (empirical Bayes; grad. desc.).

Model: $p(y | X, \theta) = \mathcal{N}(y | 0, K_y)$ with $\theta = (\ell, \sigma_f^2, \sigma_y^2)$

Negative log-likelihood:

$$L(\theta) = -\log p(y | X, \theta) = \frac{1}{2} y^T K_y^{-1} y + \frac{1}{2} \log |K_y| + \frac{N}{2} \log(2\pi) \quad (1)$$

Gradients: (via $\partial(X^{-1}) = X^{-1}(\partial X)X^{-1}$ and $\partial \det X = \frac{1}{\det X} \text{tr}((X^{-1})^T \partial X)$)

$$\begin{aligned} \frac{\partial L}{\partial \theta_j} &= -\frac{1}{2} y^T K_y^{-1} \frac{\partial K_y}{\partial \theta_j} K_y^{-1} y + \frac{1}{2} \text{tr}(K_y^{-1} \frac{\partial K_y}{\partial \theta_j}) \\ &= -\frac{1}{2} \text{tr} \left((\alpha \alpha^T - K_y^{-1}) \frac{\partial K_y}{\partial \theta_j} \right) \end{aligned}$$

with $\alpha = K_y^{-1} y$

Cholesky decomposition

How to solve $Ax = b$?

Matrix inversion: $x = A^{-1}b$ is problematic because

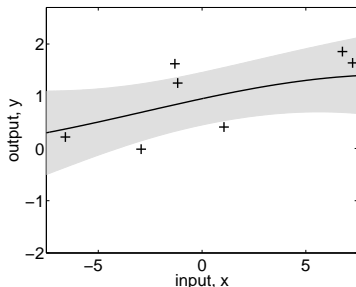
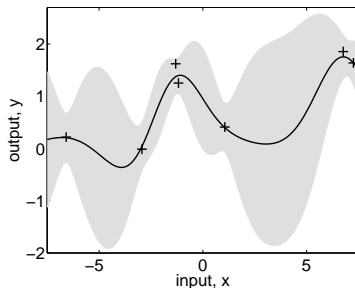
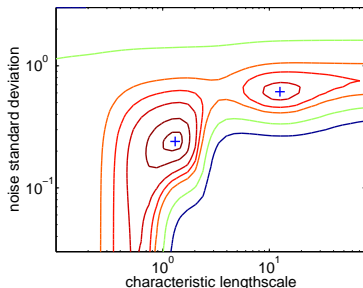
- ▶ Numerically unstable
- ▶ A^{-1} is dense, even if A is sparse

Better: LU -decomposition

$$Ax = b \xrightarrow{A=LU} \begin{cases} Lz = b \\ Ux = z \end{cases}$$

- ▶ L and U lower/upper triangular
- ▶ if A symmetric pos.-definite, then (L, U) can be chosen s.t. $U = L^T$ (Cholesky-decomposition)

Local Minima for Kernel Parameters



- ▶ top: $(l, \sigma_y) \approx (10, 0.8)$
- ▶ left: $(l, \sigma_y) \approx (1, 0.1)$
- ▶ in both cases $\sigma_f = 1$
(fixed)

[Mur12, fig. 15.5]

Semi-parametric GPs

$$f(x) = \beta^T \phi(x) + r(x)$$

$$r(X) \sim \text{GP}(r \mid 0, k(X, X))$$

Assuming

$$\beta \sim \mathcal{N}(\beta \mid b, B), \quad \text{e.g., } b := 0, B := \sigma_\beta^2 I$$

yields just another GP

$$f(X) \sim \text{GP}(\phi(X)^T b, k(X, X) + \phi(X) B \phi(X)^T)$$

where

$$\phi(X) := (\phi(x_1), \dots, \phi(x_N))^T$$

Outline

1. GPs for Regression

2. GPs for Classification

Model

$$p(y | x) := s(y f(x)), \quad y \in \{+1, -1\}, \quad s = \text{logistic}$$
$$f \sim \text{GP}(0, K(X, X))$$

► f : **latent score**

Inference

Two-step inference (given training data-set $\mathcal{D} = (X, y)$)

- infer latent score variable:

$$p(f_* | X, y, x_*) = \int p(f_* | X, x_*, f) p(f | X, y) df$$

with $p(f | X, y) = p(y | f)p(f | X)/p(y | X)$ (Bayes thm.)

- infer target:

$$\pi_* := p(y_* = +1 | X, y, x_*) = \int s(f_*) p(f_* | X, y, x_*) df_*$$

Non-Gaussians are analytically intractable.

↪ Gaussian approximation (**Laplace approximation**)

↪ **Expectation Propagation (EP)**

↪ further methods

Laplace Approximation

If h has a unique global maximum in x_0 then

$$\int_{-\infty}^{+\infty} e^{h(x)} dx \approx \sqrt{\frac{2\pi}{-h''(x_0)}} e^{h(x_0)}$$

Proof: Via Taylor $h(x) \approx h(x_0) + \frac{1}{2}h''(x_0)(x - x_0)^2$

Apply on the marginal likelihood

$$p(y | X) = \int p(y | f)p(f | x)df = \int e^{\ell(f)} \approx e^{\ell(\hat{f})} \frac{\sqrt{2\pi}^n}{\sqrt{|\nabla^2 \ell(\hat{f})|}}$$

$$\implies \log p(y | X) \approx \log p(y | \hat{f}) + \log p(\hat{f} | x) - \frac{1}{2} \log |\nabla^2 \ell(\hat{f})| + \frac{n}{2} \log(2\pi)$$

How to find \hat{f} and $\nabla^2 \ell(\hat{f})$?

Posterior

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

$$\begin{aligned} \ell(f) &= \log p(y | f) + \log p(f | X) \\ &= \log p(y | f) - \frac{1}{2} f^T K^{-1} f - \frac{1}{2} \log |K| - \frac{N}{2} \log 2\pi \\ \implies \nabla \ell(f) &= \nabla \log p(y | f) - K^{-1} f \\ \implies \nabla^2 \ell(f) &= \nabla^2 \log p(y | f) - K^{-1} \end{aligned}$$

for logistic ($p(y | f) = s(yf)$)

$$\begin{aligned} \nabla \log p(y | f) &= y - \pi \\ \nabla^2 \log p(y | f) &= \text{diag}(-\pi \circ (1 - \pi)) =: -W \end{aligned}$$

At maximum: $\nabla \ell(f) = 0 \implies f = K \nabla \log p(y | f)$

Posterior

at maximum:

$$\nabla \ell(f) = 0 \quad \implies \quad f = K \nabla \log p(y | f)$$

Use Newton to find a maximum:

$$\begin{aligned} f^{(t+1)} &:= f^{(t)} - (\nabla^2 \ell)^{-1} \nabla \ell \\ &= f^{(t)} + (K^{-1} + W^{(t)})^{-1} (\nabla \log p(y | f) - K^{-1} f^{(t)}) \\ &= (K^{-1} + W^{(t)})^{-1} (W^{(t)} f^{(t)} + \nabla \log p(y | f)) \end{aligned}$$

eventually yielding the maximum posterior \hat{f} at convergence. Then:

$$p(f | X, y) \approx q(f | X, y) := \mathcal{N}(f | \hat{f}, (K^{-1} + W)^{-1})$$

(Gaussian Approximation)

Predictions

$$\text{Recall: } \mathbb{E}_q[f_* | X, y, x_*] = K_*^T K_y^{-1} \hat{f} = K_*^T \nabla \log p(y | \hat{f})$$

$$\begin{aligned} \implies \mathbb{E}_p[f_* | X, y, x_*] &= \int \mathbb{E}[f_* | f, X, x_*] p(f | X, y) df \\ &= K_*^T K_y^{-1} \int f p(f | X, y) df \\ &= k(x_*)^T K^{-1} \mathbb{E}[f | X, y] \end{aligned}$$

approximated mean:

$$\mathbb{E}_q(f_* | X, y, x_*) = k(x_*)^T K^{-1} \hat{f}$$

variance:

$$\text{Var}_q(f_* | X, y, x_*) = k(x_*, x_*) - k_*^T (K + W^{-1})^{-1} k_*$$

predictions:

$$\bar{\pi}_* := \mathbb{E}_q(\pi_* | X, y, x_*) = \int s(f_*) q(f_* | X, y, x_*) df_*$$

Solve integral via MCMC or probit approximation (Murphy 8.4.4.2)

Algorithm (Step 1)

input:	K (covariance matrix), \mathbf{y} (± 1 targets), $p(\mathbf{y} \mathbf{f})$ (likelihood function)	
2: $\mathbf{f} := \mathbf{0}$		initialization
repeat		Newton iteration
4: $W := -\nabla\nabla \log p(\mathbf{y} \mathbf{f})$		eval. W e.g. using eq. (3.15) or (3.16)
	$L := \text{cholesky}(I + W^{\frac{1}{2}}KW^{\frac{1}{2}})$	$B = I + W^{\frac{1}{2}}KW^{\frac{1}{2}}$
6: $\mathbf{b} := W\mathbf{f} + \nabla \log p(\mathbf{y} \mathbf{f})$		} eq. (3.18) using eq. (3.27)
	$\mathbf{a} := \mathbf{b} - W^{\frac{1}{2}}L^{\top} \setminus (L \setminus (W^{\frac{1}{2}}K\mathbf{b}))$	
8: $\mathbf{f} := K\mathbf{a}$		
	until convergence	objective: $-\frac{1}{2}\mathbf{a}^{\top}\mathbf{f} + \log p(\mathbf{y} \mathbf{f})$
10: $\log q(\mathbf{y} X, \theta) := -\frac{1}{2}\mathbf{a}^{\top}\mathbf{f} + \log p(\mathbf{y} \mathbf{f}) - \sum_i \log L_{ii}$		eq. (3.32)
	return: $\mathbf{f} := \mathbf{f}$ (post. mode), $\log q(\mathbf{y} X, \theta)$ (approx. log marg. likelihood)	

Algorithm 3.1: Mode-finding for binary Laplace GPC. Commonly used convergence

Algorithm (Step 2)

input: $\hat{\mathbf{f}}$ (mode), X (inputs), \mathbf{y} (± 1 targets), k (covariance function),
 $p(\mathbf{y}|\mathbf{f})$ (likelihood function), \mathbf{x}_* test input

$$2: W := -\nabla\nabla \log p(\mathbf{y}|\hat{\mathbf{f}})$$

$$L := \text{cholesky}(I + W^{\frac{1}{2}}KW^{\frac{1}{2}})$$

$$B = I + W^{\frac{1}{2}}KW^{\frac{1}{2}}$$

$$4: \bar{f}_* := \mathbf{k}(\mathbf{x}_*)^T \nabla \log p(\mathbf{y}|\hat{\mathbf{f}})$$

$$\text{eq. (3.21)}$$

$$\mathbf{v} := L \setminus (W^{\frac{1}{2}}\mathbf{k}(\mathbf{x}_*))$$

$$\left. \begin{array}{l} \\ \end{array} \right\} \text{eq. (3.24) using eq. (3.29)}$$

$$6: \mathbb{V}[f_*] := k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^T \mathbf{v}$$

$$\bar{\pi}_* := \int \sigma(z) \mathcal{N}(z|\bar{f}_*, \mathbb{V}[f_*]) dz$$

$$\text{eq. (3.25)}$$

8: **return:** $\bar{\pi}_*$ (predictive class probability (for class 1))

Algorithm 3.2: Predictions for binary Laplace GPC. The posterior mode $\hat{\mathbf{f}}$ (which can be computed using Algorithm 3.1) is input. For multiple test inputs lines 4–7 are applied to each test input. Computational complexity is $n^3/6$ operations once (line 3) plus n^2 operations per test case (line 5). The one-dimensional integral in line 7 can be done analytically for cumulative Gaussian likelihood, otherwise it is computed using an approximation or numerical quadrature.

MCMC

How to compute integrals of the form

$$\int_a^b h(x)p(x)dx$$

where p is a probability density on $[a, b]$. LOTUS implies

$$\int_a^b h(x)p(x)dx = \mathbb{E}_p[h] \approx \frac{1}{N} \sum_{i=1}^N h(x_i) \quad (2)$$

when x_i are sampled iid from p . (**Monte-Carlo**-integration)

Markov-Chain-Monte-Carlo: Clever sampling strategy of x_i

Approximation Methods for Large Datasets

See recent literature:

- ▶ Filippone, M. and Engler, R. 2015.
Enabling scalable stochastic gradient-based inference for Gaussian processes by employing the Unbiased Linear System SolvEr (ULISSE), arXiv preprint arXiv:1501.05427. (2015).
- ▶ Dai, B., Xie, B., He, N., Liang, Y., Raj, A., Balcan, M.-F. and Song, L. 2014.
Scalable Kernel Methods via Doubly Stochastic Gradients. arXiv:1407.5599 [cs, stat]. (Jul. 2014).
- ▶ Hensman, J., Fusi, N. and Lawrence, N.D. 2013.
Gaussian processes for big data. arXiv preprint arXiv:1309.6835. (2013).

Summary

- ▶ **Gaussian processes** model continuous targets as jointly normally distributed.
 - ▶ correlated by covariance matrix depending on the predictors (**kernel**)
- ▶ **The squared exponential kernel** often is used as kernel.
 - ▶ having 2 kernel parameters: **horizontal length scale** and **vertical variation**
- ▶ **Noise variation** has to be added to the model
 - otherwise Gaussian processes interpolate the observed data.
- ▶ Kernel parameters can be learnt through gradient descent.
 - ▶ the objective is not convex, local minima need to be treated

Summary (2/2)

- ▶ For classification, Gaussian processes can be used to model
 - ▶ a **score function** f
 - ▶ that is mapped through the logistic function to probabilities π of target labels.
- ▶ The posterior is not Gaussian, but can be approximated by a Gaussian (**Laplace approximation**).
- ▶ Also the posterior predictive $E(\pi_* | x_*, X, y)$ cannot be computed analytically.
 - ▶ but it can be approximated by an integral over the (approximately) normally distributed predictive score f_*
 - ▶ and thus be computed by MCMC.

Further Readings

- ▶ Rasmussen & Williams: Gaussian Processes for Machine Learning
(free ebook!)
- ▶ See also [Mur12, chapter 15].
- ▶ Conditioning Gaussians: [Mur12, section 4.3].
- ▶ Derivatives of inverse of a matrix etc., see, e.g., *The Matrix Cookbook*, http://www.mit.edu/~wingated/stuff_i_use/matrix_cookbook.pdf

Some Matrix Derivatives

$$\begin{aligned}\partial(X^{-1}) &= -X^{-1}(\partial X)X^{-1} \\ \partial(\log(|X|)) &= \text{tr}(X^{-1}\partial X)\end{aligned}$$

Computing with traces:

$$\text{tr}(aa^T B) = a^T B a$$

References



Kevin P. Murphy.

Machine learning: a probabilistic perspective.

The MIT Press, 2012.