

305C

~~STATS271/371~~: Applied ~~Bayesian~~ Statistics III

Gaussian processes

Scott Linderman

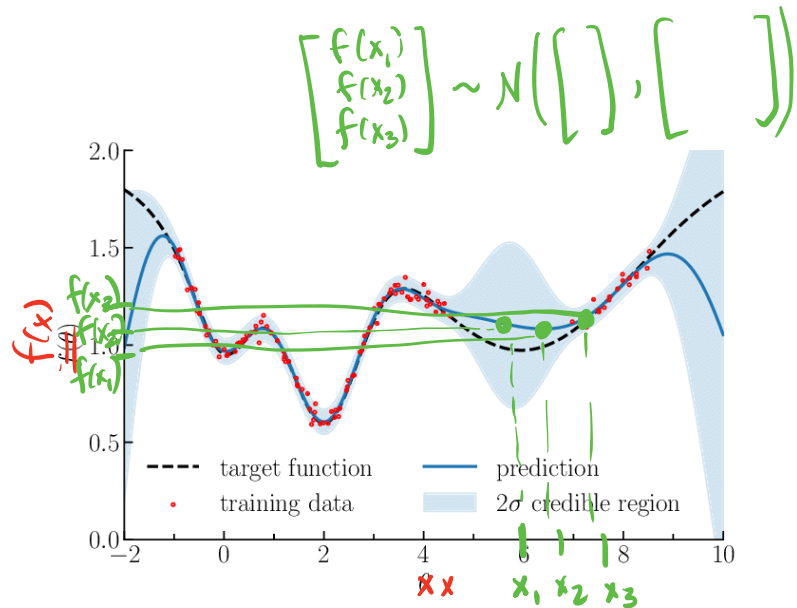
May 16, 2022

~~Lap 8~~^{Week 8}: Gaussian processes

- ▶ Gaussian processes and kernels
- ▶ Gaussian process regression
- ▶ Gaussian process classification via augmentation

Gaussian processes

- ▶ Gaussian processes are distributions on functions $f : \mathbb{R}^D \rightarrow \mathbb{R}$. (We can generalize to other domains as well.)
- ▶ Equivalently, a GP is a continuous set of r.v.'s $\{f(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^D\}$; i.e. a stochastic process.



https://www.researchgate.net/figure/Illustration-of-Gaussian-process-regression-in-one-dimension-fig1_327613136

Gaussian processes II

We say $f \sim \text{GP}(\mu(\cdot), K(\cdot, \cdot))$ if

$$\begin{bmatrix} f(\mathbf{x}_1) \\ \vdots \\ f(\mathbf{x}_N) \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu(\mathbf{x}_1) \\ \vdots \\ \mu(\mathbf{x}_N) \end{bmatrix}, \begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \cdots & K(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & & \vdots \\ K(\mathbf{x}_N, \mathbf{x}_1) & \cdots & K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix} \right) \quad (1)$$

for all finite subsets of points $\{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subset \mathbb{R}^D$.

Here, $\mu : \mathbb{R}^D \rightarrow \mathbb{R}$ is the **mean function** and $K : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$ is the **covariance function**, or **kernel**.

The covariance matrix obtained by applying the covariance function to each pair of data points above is called the **Gram matrix**.

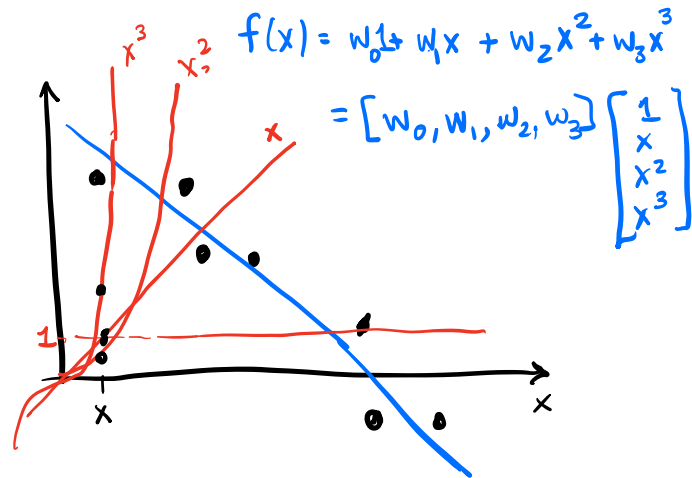
The covariance function must be positive definite; i.e. the Gram matrix must be positive definite for any subset of points.

From linear regression to GPs

- Think back to HW1 on Bayesian linear regression.
- Our motivating example was approximating a $D = 1$ dimensional function $f(x) : \mathbb{R} \rightarrow \mathbb{R}$ given noisy observations $y_n \sim \mathcal{N}(f(x_n), \sigma^2)$.
- We can cast polynomial regression as linear regression by encoding the inputs x_n with feature vectors,

$$\phi(x_n) = (x_n^0, x_n^1, \dots, x_n^{p-1}) \in \mathbb{R}^p.$$

<Draw the polynomial regression setup>



From linear regression to GPs II

More generally, let $\boldsymbol{\phi}(x) = (\phi_1(x), \dots, \phi_P(x)) \in \mathbb{R}^P$ be a function that encodes x in a P -dimensional feature space.

With these features, our linear model was,

$$\mathbb{E}[y_n | x_n] = \sum_{p=1}^P w_p \phi_p(x_n) = \boldsymbol{\phi}(x_n)^\top \mathbf{w} \triangleq f(x_n). \quad (2)$$

Now assume a Gaussian prior, $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \frac{\lambda}{P} \mathbf{I})$. Then for any $\{x_1, \dots, x_N\} \subset \mathbb{R}$,

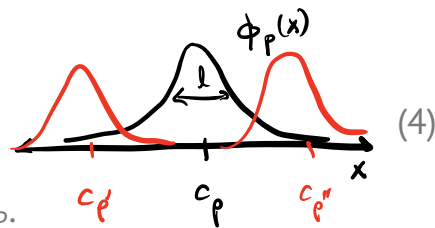
$$\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_N) \end{bmatrix} = \begin{bmatrix} \phi_1(x_1) & \cdots & \phi_P(x_1) \\ \vdots & & \vdots \\ \phi_1(x_N) & \cdots & \phi_P(x_N) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_P \end{bmatrix} = \boldsymbol{\Phi} \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \frac{\lambda}{P} \boldsymbol{\Phi} \boldsymbol{\Phi}^\top) \quad (3)$$

The function $f(\cdot)$ is a GP! It has kernel $K(x_i, x_j) = \frac{\lambda}{P} \boldsymbol{\phi}(x_i)^\top \boldsymbol{\phi}(x_j)$

Example: radial basis functions

Instead of polynomial features, let's use **radial basis functions**,

$$\phi_p(x) = e^{-\frac{1}{2\ell^2}(x-c_p)^2}.$$



These are un-normalized Gaussian bumps of width ℓ located at c_1, \dots, c_p .

With this feature encoding, the kernel is,

$$K(x_i, x_j) = \frac{\lambda}{P} \boldsymbol{\phi}(x_i)^\top \boldsymbol{\phi}(x_j) = \frac{\lambda}{P} \sum_{p=1}^P e^{-\frac{1}{2\ell^2}(x_i-c_p)^2} e^{-\frac{1}{2\ell^2}(x_j-c_p)^2} \quad (5)$$

Now take the limit of having $P \rightarrow \infty$ equally spaced centers. Then,

$$\lim_{P \rightarrow \infty} K(x_i, x_j) = \lambda \int_{-\infty}^{\infty} e^{-\frac{1}{2\ell^2}(x_i-c)^2} e^{-\frac{1}{2\ell^2}(x_j-c)^2} dc = \sqrt{\pi} \ell \lambda e^{-\frac{1}{2(\sqrt{2}\ell)^2}(x_i-x_j)^2} \quad (6)$$

This is called a **squared exponential kernel** with length scale $\sqrt{2}\ell$ and variance $\sqrt{\pi}\ell\lambda$.

Demo: sampling GPs with squared exponential kernels

https://colab.research.google.com/github/slinderman/stats305c/blob/master/notebooks/lecture15_gps.ipynb

The Matérn family of kernels

The Matérn family of kernels is defined by

$$K(x_i, x_j) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}\Delta_{ij}}{\ell} \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}\Delta_{ij}}{\ell} \right), \quad (7)$$

where

- ▶ $\Delta_{ij} = x_i - x_j$ is the difference of the points,
- ▶ ℓ is a positive length scale,
- ▶ the positive parameter ν controls the smoothness of the function, and
- ▶ and K_ν (in a potentially confusing overloading of notation) denotes the modified Bessel function.

When $\nu \rightarrow \infty$, the Matérn kernel converges to the squared exponential.

When $\nu = \frac{1}{2}$ the kernel is $K(x_i, x_j) = e^{-\frac{\Delta_{ij}}{\ell}}$, the covariance function of the **Ornstein-Uhlenbeck (OU) process**, a continuous-time AR(1) process.

When $\nu = p - \frac{1}{2}$, the kernel is the covariance of a continuous-time AR(p) process.

The Matérn family of kernels II

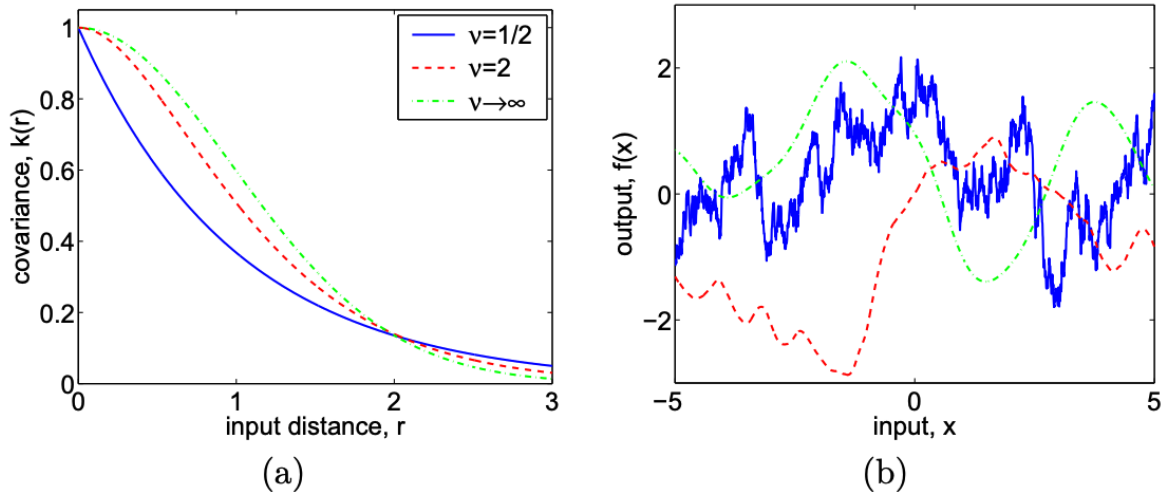


Figure 4.1: Panel (a): covariance functions, and (b): random functions drawn from Gaussian processes with Matérn covariance functions, eq. (4.14), for different values of ν , with $\ell = 1$. The sample functions on the right were obtained using a discretization of the x -axis of 2000 equally-spaced points.

Stationary kernels and Bochner's theorem

Note that both the squared exponential and the Matérn kernels are **stationary** in that $K(\mathbf{x}_i, \mathbf{x}_j)$ only depends on $\Delta_{ij} = \mathbf{x}_i - \mathbf{x}_j$.

Stationary kernels are particularly interesting because they can be defined by their power spectrum.

Theorem (Bochner's Theorem)

A complex-valued function $K(\Delta)$ on \mathbb{R}^D is the covariance function of a weakly stationary mean square continuous complex valued random process on \mathbb{R}^D if and only if it can be represented as

$$K(\Delta) = \int_{\mathbb{R}^D} e^{2\pi i \mathbf{s} \cdot \Delta} d\mu(\mathbf{s}) \quad (8)$$

where μ is a positive finite measure.

If μ has a density $S(\mathbf{s})$ then S is known as the **spectral density** or **power spectrum** corresponding to $K(\Delta)$. (This is Theorem 4.1 of Williams and Rasmussen [1996].)

Power spectrum intuition

Think of $\mu(\mathbf{s})$ as the power put into frequency \mathbf{s} . The constraint that it be positive is akin to requiring that the covariance be positive definite.

The **Wiener-Khinchine theorem** says that covariance function and the spectral density (assuming it exists) are Fourier duals of one another,

$$K(\Delta) = \int S(\mathbf{s}) e^{2\pi i \mathbf{s} \cdot \Delta} d\mathbf{s}, \quad (9)$$

$$S(\mathbf{s}) = \int K(\Delta) e^{-2\pi i \mathbf{s} \cdot \Delta} d\Delta \quad (10)$$

The variance of the process is $K(\mathbf{0}) = \int S(\mathbf{s}) d\mathbf{s}$ so the spectral density must be integrable (it must decay sufficiently fast as $|\mathbf{s}| \rightarrow \infty$) to define a valid process.

One nice consequence: the Gram matrix of a stationary kernel evaluated at evenly spaced points x_1, \dots, x_N in 1D is a **Toeplitz matrix**, which can be inverted in only $O(N \log N)$ time using the Fourier transform [Storkey, 1999, Cunningham et al., 2008]

Common kernels

covariance function	expression	S	ND
constant	σ_0^2	✓	
linear	$\sum_{d=1}^D \sigma_d^2 x_d x'_d$		
polynomial	$(\mathbf{x} \cdot \mathbf{x}' + \sigma_0^2)^p$		
squared exponential	$\exp(-\frac{r^2}{2\ell^2})$	✓	✓
Matérn	$\frac{1}{2^{\nu-1}\Gamma(\nu)} \left(\frac{\sqrt{2\nu}}{\ell}r\right)^\nu K_\nu\left(\frac{\sqrt{2\nu}}{\ell}r\right)$	✓	✓
exponential	$\exp(-\frac{r}{\ell})$	✓	✓
γ -exponential	$\exp\left(-\left(\frac{r}{\ell}\right)^\gamma\right)$	✓	✓
rational quadratic	$(1 + \frac{r^2}{2\alpha\ell^2})^{-\alpha}$	✓	✓
neural network	$\sin^{-1}\left(\frac{2\tilde{\mathbf{x}}^\top \Sigma \tilde{\mathbf{x}}'}{\sqrt{(1+2\tilde{\mathbf{x}}^\top \Sigma \tilde{\mathbf{x}})(1+2\tilde{\mathbf{x}}'^\top \Sigma \tilde{\mathbf{x}}')}}\right)$		✓

Table 4.1: Summary of several commonly-used covariance functions. The covariances are written either as a function of \mathbf{x} and \mathbf{x}' , or as a function of $r = |\mathbf{x} - \mathbf{x}'|$. Two columns marked ‘S’ and ‘ND’ indicate whether the covariance functions are stationary and nondegenerate respectively. Degenerate covariance functions have finite rank, see section 4.3 for more discussion of this issue.

The Kernel Cookbook: **Advice on Covariance functions** by **David Duvenaud**

<https://www.cs.toronto.edu/~duvenaud/cookbook/>

The “Automated Statistician”

Duvenaud et al. [2013] proposed a method to search over compositions of kernels to best fit the data. The idea is very cool, but unfortunately they came up with the *worst name ever*.

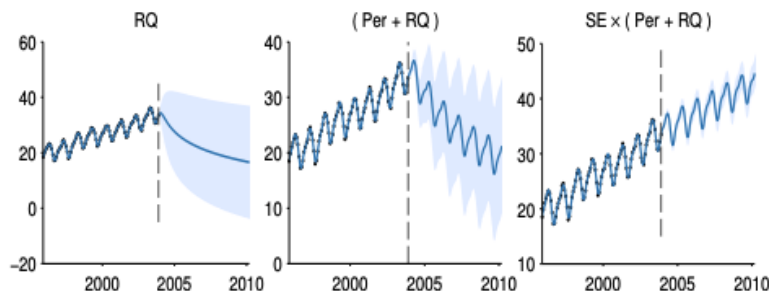


Figure 3. Posterior mean and variance for different depths of kernel search. The dashed line marks the extent of the dataset. In the first column, the function is only modeled as a locally smooth function, and the extrapolation is poor. Next, a periodic component is added, and the extrapolation improves. At depth 3, the kernel can capture most of the relevant structure, and is able to extrapolate reasonably.

Lap 8: Gaussian processes

- ▶ Gaussian processes and kernels
- ▶ **Gaussian process regression**
- ▶ Gaussian process classification via augmentation

Gaussian process regression

Now that we have a prior distribution on functions, suppose we observe (y_n, \mathbf{x}_n) pairs. Assume,

$$f \sim \text{GP}(\mu(\cdot), K(\cdot, \cdot)) \quad y_n \sim \mathcal{N}(f(\mathbf{x}_n), \sigma^2) \quad (11)$$

independently for $n = 1, \dots, N$.

Let $\mathbf{y} = [y_1, \dots, y_N]^\top$, $\boldsymbol{\mu} = [\mu(\mathbf{x}_1), \dots, \mu(\mathbf{x}_N)]^\top$, $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^\top$, and \mathbf{K} denote the Gram matrix of $\mathbf{x}_1, \dots, \mathbf{x}_N$. Under the GP prior,

$$p(\mathbf{f} \mid \{\mathbf{x}_n, y_n\}_{n=1}^N) \propto \left[\prod_{n=1}^N p(y_n \mid f(\mathbf{x}_n)) \right] p(\mathbf{f}) d\mathbf{f} \quad (12)$$

$$= \mathcal{N}(\mathbf{y} \mid \mathbf{f}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{f} \mid \boldsymbol{\mu}, \mathbf{K}) d\mathbf{f} \quad (13)$$

$$\propto \mathcal{N}(\mathbf{f} \mid \boldsymbol{\mu}', \mathbf{K}'), \quad (14)$$

where

$$\mathbf{K}' = (\mathbf{K}^{-1} + \sigma^{-2} \mathbf{I})^{-1} \quad \boldsymbol{\mu}' = \mathbf{K}'(\mathbf{K}^{-1} \boldsymbol{\mu} + \sigma^{-2} \mathbf{y}) \quad (15)$$

Marginal distribution

Likewise, we can integrate over the random function to obtain the marginal distribution,

$$p(\mathbf{y}) = \int p(\mathbf{y} | \mathbf{f}) p(\mathbf{f}) d\mathbf{f} \quad (16)$$

$$= \int \mathcal{N}(\mathbf{y} | \mathbf{f}, \sigma^2 I) \mathcal{N}(\mathbf{f} | \boldsymbol{\mu}, \mathbf{K}) d\mathbf{f} \quad (17)$$

$$= \mathcal{N}(\mathbf{y} | \boldsymbol{\mu}, \mathbf{K} + \sigma^2 I). \quad (18)$$

This allows us to compute the marginal likelihood of the data exactly.

Posterior predictive distribution

What is the posterior predictive distribution $p(y_{N+1} \mid \mathbf{x}_{N+1}, \{\mathbf{x}_n, y_n\}_{n=1}^N)$?

It's one big Gaussian model,

$$\begin{bmatrix} y_1 \\ \vdots \\ y_N \\ y_{N+1} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu(\mathbf{x}_1) \\ \vdots \\ \mu(\mathbf{x}_N) \\ \mu(\mathbf{x}_{N+1}) \end{bmatrix}, \begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \cdots & K(\mathbf{x}_1, \mathbf{x}_N) & K(\mathbf{x}_1, \mathbf{x}_{N+1}) \\ \vdots & & \vdots & \\ K(\mathbf{x}_N, \mathbf{x}_1) & \cdots & K(\mathbf{x}_N, \mathbf{x}_N) & K(\mathbf{x}_N, \mathbf{x}_{N+1}) \\ K(\mathbf{x}_{N+1}, \mathbf{x}_1) & \cdots & K(\mathbf{x}_{N+1}, \mathbf{x}_N) & K(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) \end{bmatrix} + \sigma^2 \mathbf{I} \right) \quad (19)$$

We obtain the predictive likelihood via a **Schur complement**,

$$y_{N+1} \mid \mathbf{x}_{N+1}, \{\mathbf{x}_n, y_n\}_{n=1}^N \sim \mathcal{N}(m_{N+1}, v_{N+1}) \quad (20)$$

$$m_{N+1} = \mu(\mathbf{x}_{N+1}) + \mathbf{k}^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} (\mathbf{y} - \boldsymbol{\mu}) \quad (21)$$

$$v_{N+1} = K(\mathbf{x}_{N+1}, \mathbf{x}_{N+1}) - \mathbf{k}^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k} + \sigma^2. \quad (22)$$

where $\mathbf{k} = [K(\mathbf{x}_1, \mathbf{x}_{N+1}), \dots, K(\mathbf{x}_N, \mathbf{x}_{N+1})]^\top \in \mathbb{R}^N$.

GPs are linear predictors

Assume that $\mu(\mathbf{x}) = 0$ everywhere. Then we can write the predictive mean as,

$$m_{N+1} = \sum_{n=1}^N \alpha_n y_n, \quad (23)$$

where

$$\alpha_n = [K^{-1} \mathbf{k}]_n \quad (24)$$

Handwritten red text: $K \in \mathbb{R}^{N \times N}$

Question: what is the complexity of computing the predictive mean? *$O(N^3)$*

Recall the posterior predictive distribution in Bayesian linear regression

In Bayesian linear regression, we first compute the posterior mean of the weights.

In our notation from Slide 6, the posterior mean (under an uninformative prior) was,

$$\mathbb{E}[\mathbf{w} \mid \{\mathbf{x}_n, y_n\}_{n=1}^N] = \underbrace{(\Phi^\top \Phi + I)}_{\Phi^\top \Phi \in \mathbb{R}^{p \times p}}^{-1} (\Phi^\top \mathbf{y}) \quad (25)$$

Handwritten notes: $\Phi \in \mathbb{R}^{N \times p}$ (with an arrow pointing to Φ), $\Phi^\top \Phi \in \mathbb{R}^{p \times p}$ (under the bracket).

so the prediction would be

$$\mathbb{E}[y_{N+1} \mid \mathbf{x}_{N+1}, \{\mathbf{x}_n, y_n\}_{n=1}^N] = \phi(\mathbf{x}_{N+1})^\top (\Phi^\top \Phi + I)^{-1} (\Phi^\top \mathbf{y}) \quad (26)$$

Question: what is the complexity of computing the predictive mean in Bayesian linear regression?

Handwritten note: $\sim O(p^3)$

Question: When is it more efficient to work with kernels rather than weights?

Handwritten notes:
 $p \ll N \Rightarrow$ weight space
else \Rightarrow func space

Lap 8: Gaussian processes

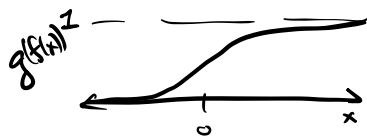
- ▶ Gaussian processes and kernels
- ▶ Gaussian process regression
- ▶ **Gaussian process classification via augmentation**

GP Classification

We've seen how Gaussian processes can be used for regression problems with $y_n \in \mathbb{R}$. What if we have binary observations under the following model,

$$f \sim \text{GP}(\mu(\cdot), K(\cdot, \cdot)) \qquad y_n \sim \text{Bern}(g(f(\mathbf{x}_n))) \qquad (27)$$

where $g : \mathbb{R} \rightarrow [0, 1]$ is, for example, the logistic or probit function.



This is the GP analog of **logistic regression** or **probit regression**.

As in standard logistic regression, the posterior is *not* available in closed form. All we know is,

$$p(\mathbf{f} \mid \{\mathbf{x}_n, y_n\}_{n=1}^N) \propto \left[\prod_{n=1}^N \text{Bern}(y_n \mid g(f(\mathbf{x}_n))) \right] \mathcal{N}(\mathbf{f} \mid \boldsymbol{\mu}, \mathbf{K}) \qquad (28)$$

We will show how one can perform Bayesian inference in these models via **augmentation**.

GP Classification with a probit mean function

First, consider the **probit function**, a.k.a. the normal c.d.f.

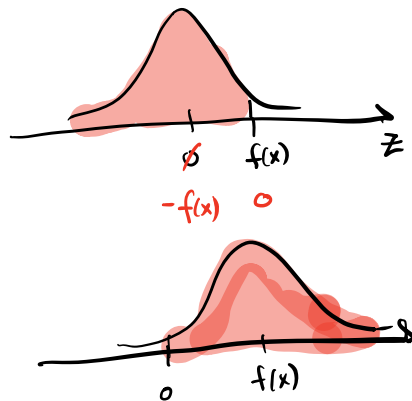
$$g(f(\mathbf{x})) = \Pr(z \leq f(\mathbf{x})) \quad \text{where} \quad z \sim \mathcal{N}(0, 1). \quad (29)$$

Exercise: Show that $g(f(\mathbf{x})) = \int_0^\infty \mathcal{N}(z | f(\mathbf{x}), 1) dz$.

$$= \int_{-\infty}^{f(\mathbf{x})} \mathcal{N}(z | 0, 1) dz$$

$$= \int_{-\infty}^0 \mathcal{N}(z | -f(\mathbf{x}), 1) dz$$

$$= \int_0^\infty \mathcal{N}(z | f(\mathbf{x}), 1) dz$$



GP Classification with a probit mean function II

Now plug this identity into the likelihood from (28),

$$p(y_n | \mathbf{f}, \mathbf{x}_n) \propto \text{Bern}(y_n | g(f(\mathbf{x}_n))) \quad (30)$$

$$= g(f(\mathbf{x}_n))^{\mathbb{I}[y_n=1]} (1 - g(f(\mathbf{x}_n)))^{\mathbb{I}[y_n=0]} \quad (31)$$

$$= \left(\int_0^\infty \mathcal{N}(z_n | f(\mathbf{x}_n), 1) \mathrm{d}z_n \right)^{\mathbb{I}[y_n=1]} \left(\int_{-\infty}^0 \mathcal{N}(z_n | f(\mathbf{x}_n), 1) \mathrm{d}z_n \right)^{\mathbb{I}[y_n=0]} \quad (32)$$

$$= \int_{-\infty}^\infty p(y_n | z_n) \mathcal{N}(z_n | f(\mathbf{x}_n), 1) \mathrm{d}z_n \quad (33)$$

$$= \int_{-\infty}^\infty p(y_n, z_n | \mathbf{f}, \mathbf{x}_n) \mathrm{d}z_n \quad (34)$$

where

$$p(y_n | z_n) = \begin{cases} \mathbb{I}[z_n > 0] & \text{if } y_n = 1 \\ \mathbb{I}[z_n \leq 0] & \text{if } y_n = 0 \end{cases} \quad (35)$$

GP Classification with a probit mean function III

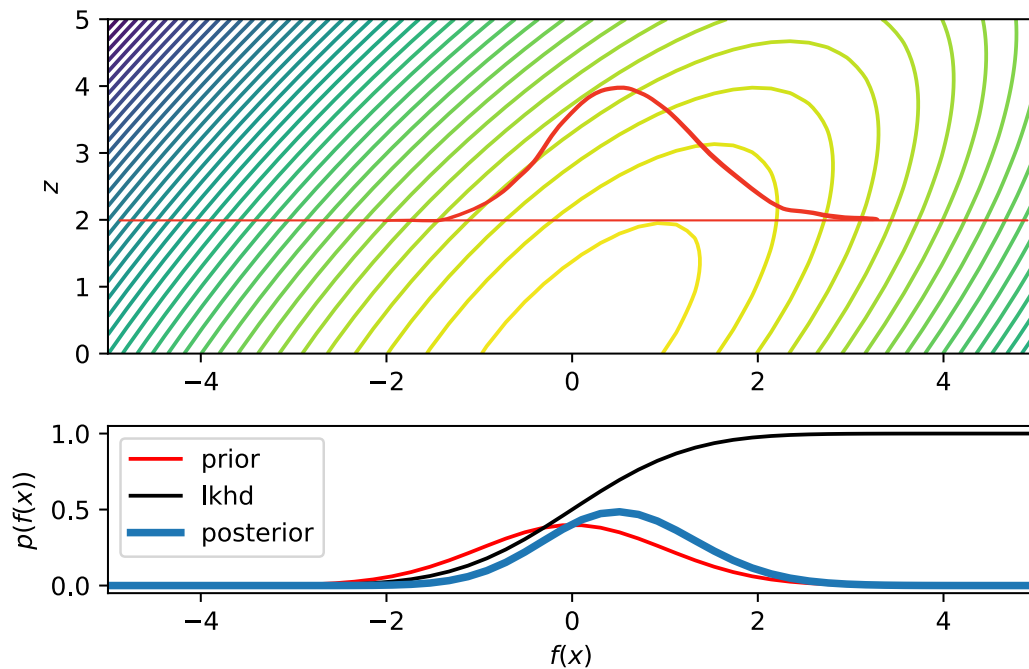


Figure: Augmentation scheme for GP Probit Classification [Albert and Chib, 1993, Girolami and Rogers, 2006]

GP Classification with a probit mean function IV

This suggests the following **augmentation scheme** [Albert and Chib, 1993, Girolami and Rogers, 2006].
Design a Gibbs sampler to target the augmented posterior distribution,

$$p(\mathbf{f}, \{z_n\}_{n=1}^N \mid \{\mathbf{x}_n, y_n\}_{n=1}^N) \propto \prod_{n=1}^N \left[p(y_n \mid z_n) \mathcal{N}(z_n \mid f(\mathbf{x}_n), 1) \right] \mathcal{N}(\mathbf{f} \mid \boldsymbol{\mu}, \mathbf{K}). \quad (36)$$

Exercise: Derive the conditional distributions,

$$p(z_n \mid \mathbf{f}, \{\mathbf{x}_n, y_n\}_{n=1}^N) = \text{truncated normal} \quad (37)$$

and

$$p(\mathbf{f} \mid \{\mathbf{x}_n, y_n, z_n\}_{n=1}^N) = \text{gaussian} \quad (38)$$

Extra Slides

- ▶ Elliptical slice sampling
- ▶ Sparse GPs and inducing points
- ▶ Stochastic variational inference for sparse GPs

Slice sampling intuition

Slice sampling

Slice sampling [Neal, 2003] is a MCMC algorithm for distributions with intractable normalization constants.

Suppose we want to sample a posterior distribution,

$$p(\boldsymbol{\theta} \mid \mathbf{y}) = \frac{p(\boldsymbol{\theta}, \mathbf{y})}{p(\mathbf{y})}. \quad (39)$$

Assume we can evaluate the joint probability (numerator) but not the marginal (denominator).

Slice sampling works by introducing an **auxiliary variable** $u \in \mathbb{R}$ so that,

$$p(u, \boldsymbol{\theta} \mid \mathbf{y}) = \begin{cases} \frac{1}{p(\mathbf{y})} & \text{if } 0 \leq u \leq p(\boldsymbol{\theta}, \mathbf{y}) \\ 0 & \text{o.w.} \end{cases} \quad (40)$$

Note that marginalizing over u recovers the posterior,

$$\int p(u, \boldsymbol{\theta} \mid \mathbf{y}) \, du = \int_0^{p(\boldsymbol{\theta}, \mathbf{y})} \frac{1}{p(\mathbf{y})} \, du = \frac{p(\boldsymbol{\theta}, \mathbf{y})}{p(\mathbf{y})} = p(\boldsymbol{\theta} \mid \mathbf{y}). \quad (41)$$

If we can sample $p(u, \boldsymbol{\theta} \mid \mathbf{y})$, we can simply throw away the u 's to get samples from $p(\boldsymbol{\theta} \mid \mathbf{y})$.

Slice sampling II

Gibbs sampling is often easy in the augmented space of $(u, \boldsymbol{\theta})$:

Given $\boldsymbol{\theta}$, the auxiliary variable u follows a uniform distribution,

$$p(u \mid \boldsymbol{\theta}, \mathbf{y}) \propto p(u, \boldsymbol{\theta} \mid \mathbf{y}) = \frac{1}{p(\mathbf{y})} \mathbb{I}[0 \leq u \leq p(\boldsymbol{\theta}, \mathbf{y})] \propto \text{Unif}(u; [0, p(\boldsymbol{\theta}, \mathbf{y})]). \quad (42)$$

Given u , the variable of interest $\boldsymbol{\theta}$ is uniformly distributed as well,

$$p(\boldsymbol{\theta} \mid u, \mathbf{y}) \propto p(u, \boldsymbol{\theta} \mid \mathbf{y}) = \frac{1}{p(\mathbf{y})} \mathbb{I}[p(\boldsymbol{\theta}, \mathbf{y}) \geq u] \propto \text{Unif}(\boldsymbol{\theta}; \Theta(u)), \quad (43)$$

where $\Theta(u) = \{\boldsymbol{\theta}' : p(\boldsymbol{\theta}', \mathbf{y}) \geq u\}$ is a **slice** of parameter space with joint probability at least u .

If we can find that slice and sample uniformly from it, we're in business!

Slice sampling III

$\Theta(u)$ is generally a complex subset that is hard to sample uniformly.

Instead, we typically sampling one coordinate of θ at a time, holding the rest of fixed. (This is also a valid Gibbs update.)

Still, we need some way of finding the slice. Neal [2003] proposes a “stepping out” procedure to compute a 1D slice.

Slice sampling for GPs

Consider a GP classification with 2 data points $(\mathbf{x}_1, \mathbf{x}_2)$. Assume they are close relative to the length-scale of a squared exponential kernel.

Accounting for the correlated Gaussian prior

Single variable slice sampling, like all Gibbs sampling methods, suffers when the coordinates are correlated.

For Gaussian processes, the prior induces correlations between the function values, and this can seriously impair such coordinate-wise update algorithms.

Can we reparameterize the problem so that our 1D slice is better aligned with the prior?

A strange way to sample a Gaussian...

Suppose $\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \Sigma)$. A strange but interesting way to sample \mathbf{f} is via the following program:

$$\mathbf{v}_0 \sim \mathcal{N}(\mathbf{0}, \Sigma) \quad (44)$$

$$\mathbf{v}_1 \sim \mathcal{N}(\mathbf{0}, \Sigma) \quad (45)$$

$$\theta \sim \text{Unif}([0, 2\pi)) \quad (46)$$

$$\mathbf{f} = \mathbf{v}_0 \sin \theta + \mathbf{v}_1 \cos \theta \quad (47)$$

To see this, define $\mathbf{v}' = [\mathbf{v}_0 \quad \mathbf{v}_1]^\top$ and note that for any $\theta \in [0, 2\pi)$

$$\mathbf{v}' \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Sigma & \mathbf{0} \\ \mathbf{0} & \Sigma \end{bmatrix}\right), \quad (48)$$

$$\mathbf{f} = [\sin \theta I \quad \cos \theta I] \mathbf{v}' \sim \mathcal{N}\left(\mathbf{0}, [\sin \theta I \quad \cos \theta I] \begin{bmatrix} \Sigma & \mathbf{0} \\ \mathbf{0} & \Sigma \end{bmatrix} \begin{bmatrix} \sin \theta I \\ \cos \theta I \end{bmatrix}\right) \quad (49)$$

$$\sim \mathcal{N}(\mathbf{0}, (\sin^2 \theta + \cos^2 \theta) \Sigma) = \mathcal{N}(\mathbf{0}, \Sigma). \quad (50)$$

Elliptical slice sampling

With this sampling procedure in mind, Murray et al. [2010] proposed **elliptical slice sampling**.

Reparameterize \mathbf{f} in terms of $\mathbf{v}_0, \mathbf{v}_1, \theta$ and sample the posterior,

$$p(\mathbf{v}_0, \mathbf{v}_1, \theta \mid \mathbf{y}) \propto \mathcal{N}(\mathbf{v}_0 \mid \mathbf{0}, \Sigma) \mathcal{N}(\mathbf{v}_1 \mid \mathbf{0}, \Sigma) \text{Unif}(\theta \mid [0, 2\pi)) p(\mathbf{y} \mid \mathbf{f}(\mathbf{v}_0, \mathbf{v}_1, \theta)) \quad (51)$$

Their sampling algorithm consists of two steps:

1. Sample $\mathbf{v}_0, \mathbf{v}_1, \theta \mid \mathbf{y}, \mathbf{f} = \mathbf{v}_0 \sin \theta + \mathbf{v}_1 \cos \theta$
 - 1.1 Sample $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ and $\theta \sim \text{Unif}([0, 2\pi))$.
 - 1.2 Set $\mathbf{v}_0 = \mathbf{f} \sin \theta + \mathbf{v} \cos \theta$ and $\mathbf{v}_1 = \mathbf{f} \cos \theta - \mathbf{v}' \sin \theta$.
2. Sample $\theta \mid \mathbf{v}_0, \mathbf{v}_1, \mathbf{y}$ by slice sampling its conditional $p(\theta \mid \mathbf{v}_0, \mathbf{v}_1, \mathbf{y}) \propto p(\mathbf{y} \mid \mathbf{v}_0 \cos \theta + \mathbf{v}_1 \sin \theta)$.

We can think of step 2 as slice sampling on an elliptical slice, hence the name.

Elliptical slice sampling II

This formulation is a little complicated. There's no need to actually instantiate ν_0 and ν_1 . Instead, we can just work with ν and θ .

Input: current state \mathbf{f} , a routine that samples from $\mathcal{N}(0, \Sigma)$, log-likelihood function $\log L$.

Output: a new state \mathbf{f}' . When \mathbf{f} is drawn from $p^*(\mathbf{f}) \propto \mathcal{N}(\mathbf{f}; 0, \Sigma) L(\mathbf{f})$, the marginal distribution of \mathbf{f}' is also p^* .

1. Choose ellipse: $\nu \sim \mathcal{N}(0, \Sigma)$

2. Log-likelihood threshold:

$$u \sim \text{Uniform}[0, 1]$$

$$\log y \leftarrow \log L(\mathbf{f}) + \log u$$

3. Draw an initial proposal, also defining a bracket:

$$\theta \sim \text{Uniform}[0, 2\pi]$$

$$[\theta_{\min}, \theta_{\max}] \leftarrow [\theta - 2\pi, \theta]$$

4. $\mathbf{f}' \leftarrow \mathbf{f} \cos \theta + \nu \sin \theta$

5. **if** $\log L(\mathbf{f}') > \log y$ **then:**

6. Accept: **return** \mathbf{f}'

7. **else:**

 Shrink the bracket and try a new point:

8. **if** $\theta < 0$ **then:** $\theta_{\min} \leftarrow \theta$ **else:** $\theta_{\max} \leftarrow \theta$

9. $\theta \sim \text{Uniform}[\theta_{\min}, \theta_{\max}]$

10. **GoTo** 4.

Elliptical slice sampling III

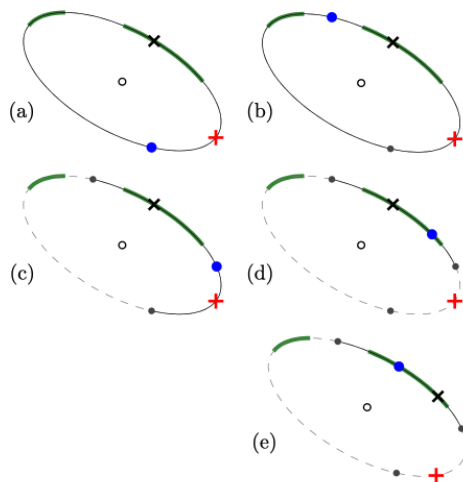


Figure 3: **(a)** The algorithm receives $\mathbf{f}=\mathbf{x}$ as input. Step 1 draws auxiliary variate $\nu=+$, defining an ellipse centred at the origin (o). Step 2: a likelihood threshold defines the ‘slice’ (—). Step 3: an initial proposal \bullet is drawn, in this case not on the slice. **(b)** The first proposal defined both edges of the $[\theta_{\min}, \theta_{\max}]$ bracket; the second proposal (\bullet) is also drawn from the whole range. **(c)** One edge of the bracket (—) is moved to the last rejected point such that \mathbf{x} is still included. Proposals are made with this shrinking rule until one lands on the slice. **(d)** The proposal here (\bullet) is on the slice and is returned as \mathbf{f}' . **(e)** Shows the reverse configuration discussed in Section 2.3: \mathbf{x} is the input \mathbf{f}' , which with auxiliary $\nu'=+$ defines the same ellipse. The brackets and first three proposals (\bullet) are the same. The final proposal (\bullet) is accepted, a move back to \mathbf{f} .

Computational complexity

Question: What is the computational complexity of one elliptical slice sampling update? Assume conditionally independent Bernoulli observations, as in GP classification.

Lap 8: Gaussian processes, elliptical slice sampling, and Bayesian optimization

- ▶ **Model:** Gaussian processes
- ▶ **Algorithm:** Elliptical slice sampling
- ▶ **Algorithm:** Variational Inference with Sparse GPs
- ▶ **Application:** Bayesian optimization

Inducing point methods

The key limitation of GPs is the cubic complexity of inference (with quadratic memory complexity).

One way of circumventing this complexity is via **inducing points** [???]. These form the basis of an approximate model called **Sparse GPs**.

We follow the presentation by Hensman et al. [2013].

Sparse GPs and Inducing Points

Consider a Gaussian process regression with inputs $\{\mathbf{x}_n\}_{n=1}^N \subset \mathbb{R}^D$ and observations $\{y_n\}_{n=1}^N \subset \mathbb{R}$. As above, let $\mathbf{y} = [y_1, \dots, y_N]^\top$ and $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^\top$

Introduce a set of **inducing points** $\{\mathbf{z}_m\}_{m=1}^M \subset \mathbb{R}^D$ with corresponding function values $u_m = f(\mathbf{z}_m)$.

Using the GP predictive distribution, we can write,

$$p(\mathbf{y} | \mathbf{f}) = \mathcal{N}(\mathbf{y} | \mathbf{f}, \sigma^2 \mathbf{I}) \quad (52)$$

$$p(\mathbf{f} | \mathbf{u}) = \mathcal{N}(\mathbf{f} | \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{u}, \tilde{\mathbf{K}}) \quad (53)$$

where \mathbf{K}_{nm} is the Gram matrix of kernel evaluations for each $(\mathbf{x}_n, \mathbf{z}_m)$ pair, \mathbf{K}_{mm} is the Gram matrix for each $(\mathbf{z}_m, \mathbf{z}_m)$ pair, and

$$\tilde{\mathbf{K}} = \mathbf{K}_{nn} - \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn}. \quad (54)$$

Note that $\tilde{\mathbf{K}} \in \mathbb{R}^{N \times N}$, so we're still looking at $O(N^3)$ complexity.

Variational inference for Sparse GPs

Now we will lower bound the marginal log likelihood using Jensen's inequality,

$$\log p(\mathbf{y} \mid \mathbf{u}) = \log \int p(\mathbf{y}, \mathbf{f} \mid \mathbf{u}) d\mathbf{f} \quad (55)$$

$$= \log \mathbb{E}_{p(\mathbf{f} \mid \mathbf{u})} [p(\mathbf{y} \mid \mathbf{f})] \quad (56)$$

$$\geq \mathbb{E}_{p(\mathbf{f} \mid \mathbf{u})} [\log p(\mathbf{y} \mid \mathbf{f})] \quad (57)$$

$$\triangleq \mathcal{L}_1. \quad (58)$$

Assume $p(\mathbf{y} \mid \mathbf{f}) = \prod_{n=1}^N p(y_n \mid f(\mathbf{x}_n))$; i.e. the observations are conditionally independent.

Question: what is the complexity of evaluating the bound in eq. (57)?

Question: when is the bound maximized?

Variational inference for Sparse GPs II

Using the bound from above, we can obtain a lower bound on the marginal likelihood of the data, integrating out the inducing variables [Titsias, 2009],

$$\log p(\mathbf{y} \mid \mathbf{X}) = \log \int p(\mathbf{y} \mid \mathbf{u}) p(\mathbf{u}) d\mathbf{u} \quad (59)$$

$$\geq \log \int \exp\{\mathcal{L}_1\} p(\mathbf{u}) d\mathbf{u} \quad (60)$$

$$\triangleq \mathcal{L}_2. \quad (61)$$

After some algebra, this bound equals,

$$\mathcal{L}_2 = \log \mathcal{N}(\mathbf{y} \mid \mathbf{0}, \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn} + \sigma^2 \mathbf{I}) - \frac{1}{2\sigma^2} \text{Tr}(\tilde{\mathbf{K}}). \quad (62)$$

We can view this bound as arising from a variational approximation $q(\mathbf{u}) = \mathcal{N}(\hat{\mathbf{u}}, \mathbf{\Lambda}^{-1})$ where

$$\mathbf{\Lambda} = \sigma^{-2} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn} \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} + \mathbf{K}_{mm}^{-1} \quad (63)$$

$$\hat{\mathbf{u}} = \sigma^{-2} \mathbf{\Lambda}^{-1} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn} \mathbf{y}. \quad (64)$$

Question: what is the complexity of computing \mathcal{L}_2 ?

Stochastic variational inference for Sparse GPs

Marginalizing over \mathbf{u} coupled the observations \mathbf{y} so that the likelihood was a multivariate Gaussian (albeit with a low rank plus diagonal covariance).

As we briefly discussed in Lap 5: Mixed Membership Models, if the likelihood factors into a product over data points we can use **stochastic variational inference**, operating on mini-batches of data instead.

Can we do something similar here?

Stochastic variational inference for Sparse GPs II

Hensman et al. [2013] showed that if we explicitly represent $q(\mathbf{u})$, the likelihood $p(\mathbf{y} \mid \mathbf{u})$ factors into a product over the N data points.

Define a new lower bound,

$$\log p(\mathbf{y} \mid \mathbf{X}) \geq \mathbb{E}_{q(\mathbf{u})} [\log p(\mathbf{y} \mid \mathbf{u}) + \log p(\mathbf{u}) - \log q(\mathbf{u})] \quad (65)$$

$$\geq \mathbb{E}_{q(\mathbf{u})} [\mathcal{L}_1 + \log p(\mathbf{u}) - \log q(\mathbf{u})] \quad (66)$$

$$\triangleq \mathcal{L}_3, \quad (67)$$

and assume $q(\mathbf{u}) = \mathcal{N}(\mathbf{m}, \mathbf{S})$. Then the bound simplifies to,

$$\mathcal{L}_3 = \sum_{n=1}^N \left[\log \mathcal{N}(y_n \mid \mathbf{k}_n^\top \mathbf{K}_{mm}^{-1} \mathbf{m}, \sigma^2) - \frac{1}{2\sigma^2} \tilde{K}_{n,n} - \frac{1}{2} \text{Tr}(\mathbf{S} \mathbf{\Lambda}_n) \right] - D_{\text{KL}}(q(\mathbf{u}) \parallel p(\mathbf{u})) \quad (68)$$

where \mathbf{k}_n is the k -th column of \mathbf{K}_{mm} and $\mathbf{\Lambda}_n = \sigma^{-1} \mathbf{K}_{mm}^{-1} \mathbf{k}_n \mathbf{k}_n^\top \mathbf{K}_{mm}^{-1}$.

This lower bound and its gradients wrt \mathbf{m} and \mathbf{S} can be approximated with Monte Carlo using mini-batches. Hensman et al. [2013] used this trick to fit sparse GPs to millions of data points.

Choosing the inducing points

Question: how should we obtain the inducing points?

References I

Christopher KI Williams and Carl Edward Rasmussen. *Gaussian processes for regression*. MIT Press, 1996.

Amos J Storkey. Truncated covariance matrices and toeplitz methods in gaussian processes. In *Artificial Neural Networks, 1999. ICANN 99. Ninth International Conference on (Conf. Publ. No. 470)*, volume 1, pages 55–60 vol.1. unknown, February 1999.

John P Cunningham, Krishna V Shenoy, and Maneesh Sahani. Fast Gaussian process methods for point process intensity estimation. In *Proceedings of the 25th International Conference on Machine Learning*, pages 192–199, New York, NY, USA, July 2008. Association for Computing Machinery.

David Duvenaud, James Lloyd, Roger Grosse, Joshua Tenenbaum, and Ghahramani Zoubin. Structure discovery in nonparametric regression through compositional kernel search. In *International Conference on Machine Learning*, pages 1166–1174. PMLR, 2013.

James H Albert and Siddhartha Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American statistical Association*, 88(422):669–679, 1993.

References II

- Mark Girolami and Simon Rogers. Variational Bayesian multinomial probit regression with Gaussian process priors. *Neural Computation*, 18(8):1790–1817, 2006.
- Radford M Neal. Slice sampling. *Annals of Statistics*, 31(3):705–767, June 2003.
- Iain Murray, Ryan Adams, and David MacKay. Elliptical slice sampling. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 541–548. jmlr.org, March 2010.
- James Hensman, Nicolò Fusi, and Neil D Lawrence. Gaussian processes for big data. In *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence*, UAI’13, pages 282–290, Arlington, Virginia, USA, August 2013. AUAI Press.
- Michalis Titsias. Variational learning of inducing variables in sparse gaussian processes. In David van Dyk and Max Welling, editors, *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics*, volume 5 of *Proceedings of Machine Learning Research*, pages 567–574, Hilton Clearwater Beach Resort, Clearwater Beach, Florida USA, 2009. PMLR.