

# STATISTICAL LEARNING: GAUSSIAN PROCESS REGRESSION

---

Jakob Söhl

Delft University of Technology

Introduction

Gaussian process regression

Advanced topics (not part of the exam)

Hyperparameter estimation by empirical Bayes

Gaussian process classification (RG section 8.3)

Gibbs sampling for GP classification (RG section 9.2)

# Introduction

---

## Some historical background

- Long history in geostatistics/spatial statistics: kriging.
- Popularised in machine learning 1995–
- Influential papers by Radford Neal.
- Well known book: *Gaussian Processes for Machine Learning*, Carl Edward Rasmussen and Chris Williams, the MIT Press, 2006.  
<https://gaussianprocess.org/gpml/>
- Builds upon theory on stochastic processes in continuous time.

## Remark on notation

- I try to avoid boldface, bars, etc.
- So a vector is written like  $y \in \mathbb{R}^n$ .
- If there is risk to confusion I use boldface to denote a vector.

Example:  $\boldsymbol{x}_i = \begin{bmatrix} 1 & x_i \end{bmatrix}^T$ .

# Gaussian process regression

---

# Revisiting the Bayesian linear model

The model is given by

$$\begin{aligned}y \mid \theta &\sim N(X\theta, \Sigma) \\ \theta &\sim N(0, \Sigma_0)\end{aligned}$$

Assume the noise covariance  $\Sigma$  is known.

**Claim:**  $y \sim N(0, \mathcal{K})$  with  $\mathcal{K} = \Sigma + X\Sigma_0X^T$ .

Checking that the mean and covariance are as stated:

1.  $\mathbb{E}[y] = \mathbb{E}\mathbb{E}[y \mid \theta] = \mathbb{E}[X\theta] = X\mathbb{E}\theta = 0$ .
2. Use the [law of total covariance](#):

$$\begin{aligned}\text{Cov } y &= \mathbb{E} \text{Cov}(y \mid \theta) + \text{Cov } \mathbb{E}[y \mid \theta] \\ &= \mathbb{E}\Sigma + \text{Cov}(X\theta) = \Sigma + X\Sigma_0X^T.\end{aligned}$$

# Implied covariance structure

Assume  $\Sigma = \sigma^2 I_n$ ,  $\Sigma_0 = \sigma_0^2 I_p$ .

$$\text{Cov } y = \Sigma + X\Sigma_0X^T = \sigma^2 I_n + \sigma_0^2 XX^T.$$

This implies

$$\text{Cov}(y_i, y_j) = \sigma^2 \delta_{ij} + \sigma_0^2 x_i^T x_j,$$

where

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

and  $x_i^T$  is the  $i$ -th row of  $X$ . Let

$$\mathcal{K}_{ij} = \sigma^2 \delta_{ij} + \sigma_0^2 x_i^T x_j.$$

If we define

$$K(x, \tilde{x}) = \sigma^2 \mathbf{1}_{\{x=\tilde{x}\}} + \sigma_0^2 x^T \tilde{x}$$

then  $\mathcal{K}_{ij} = K(x_i, x_j)$ .



## Example: linear regression with $M$ basis functions (1/2)

Assume observation pairs  $(u_i, y_i)$  and the model

$$y_i = \sum_{k=1}^p \theta_k \varphi_k(u_i) + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2).$$

Thus, **feature vector**

$$x_i = \begin{bmatrix} \varphi_1(u_i) & \varphi_2(u_i) & \dots & \varphi_p(u_i) \end{bmatrix}^T$$

and **parameter vector**  $\theta = \begin{bmatrix} \theta_1 & \dots & \theta_p \end{bmatrix}^T$ .

**Bayesian analysis:** suppose prior

$$\{\theta_k\} \stackrel{\text{ind}}{\sim} N(0, \lambda_k), \quad k = 1, \dots, p$$

so  $\theta \sim N_p(0, \Sigma_0)$  with  $\Sigma_0 = \text{diag}(\lambda_1, \dots, \lambda_p)$ .

## Example: linear regression with $M$ basis functions (2/2)

Recall

$$y \sim N(0, \mathcal{K}), \quad \mathcal{K} = X\Sigma_0 X^T + \Sigma.$$

This implies

$$\mathcal{K}_{ij} = \sum_{k=1}^p \lambda_k \varphi_k(u_i) \varphi_k(u_j) + \delta_{ij} \sigma^2.$$

Gaussian Process Regression (GPR):

- Consider  $p \rightarrow \infty$  (implies conditions on decay of  $\lambda_k$ , as  $k \rightarrow \infty$ ).
- Key idea: start from the kernel  $K$ .

# Gaussian processes in case of noiseless regression

Main idea: nonparametric regression by assuming

$$y_i = f(x_i) + \underbrace{\epsilon_i}_{\text{assume 0 (just for now)}} .$$

Denote  $f_i \equiv f(x_i)$ . So that  $y = If$  (with  $I$  the identity matrix), a linear model without noise.

Define the prior

$$\begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix} \sim N_n \left( \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_n) \\ \vdots & \vdots & \dots & \vdots \\ K(x_n, x_1) & K(x_n, x_2) & \dots & K(x_n, x_n) \end{bmatrix} \right),$$

where the **kernel function**  $K(x, y)$  needs to be such that the covariance matrix is symmetric and positive semi definite (psd).

# What is a Gaussian process?

**Definition.** Let  $x_i, i = 1, \dots, n$  be feature vectors, say in  $\mathbb{R}^p$ . Let

- $\mu : \mathbb{R}^p \rightarrow \mathbb{R}$  be the **mean function**
- $K : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$  be the **(covariance) kernel function**.

Then  $x \mapsto f_x$  is a **Gaussian process** with mean function  $\mu$  and kernel  $K$  if for *any* finite set  $\{x_1, \dots, x_n\}$ , with  $x_i \in \mathbb{R}^p$  for all  $i$ , we have

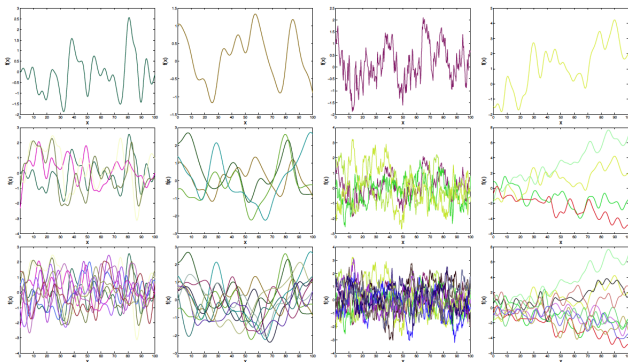
$$\begin{bmatrix} f_{x_1} \\ f_{x_2} \\ \vdots \\ f_{x_n} \end{bmatrix} \sim N_n \left( \begin{bmatrix} \mu(x_1) \\ \mu(x_2) \\ \vdots \\ \mu(x_n) \end{bmatrix}, \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \dots & K(x_1, x_n) \\ K(x_2, x_1) & K(x_2, x_2) & \dots & K(x_2, x_n) \\ \vdots & \vdots & \dots & \vdots \\ K(x_n, x_1) & K(x_n, x_2) & \dots & K(x_n, x_n) \end{bmatrix} \right).$$

We write  $f \sim GP(\mu, K)$ . The matrix with elements  $K_{ij} \equiv K(x_i, x_j)$  is called the **Gram**-matrix.

It can be shown that such processes exist, with realisations that are continuous functions (under certain conditions).

# Some realisations of Gaussian processes

Samples from GPs with different  $K(x, x')$



- $K$  needs to be such that the covariance matrix of  $(Z_{t_1}, \dots, Z_{t_n})$  is positive definite (advanced: Bochner's theorem gives conditions for stationary covariance functions).
- A covariance function is **stationary** if  $K(x, y)$  is a function of  $x - y$ . This yields invariance with respect to translations in feature (input) space.
- A covariance function is **isotropic** if  $K(x, y)$  is a function of  $\|x - y\|$ . This yields invariance with respect to rigid motions in input space.

- Popular choice **Square exponential kernel**

$$K(x, \tilde{x}) = v_0 \exp \left( -\frac{\|x - \tilde{x}\|^2}{2\ell^2} \right),$$

with

- $v_0$  signal variance
  - $\ell$  characteristic length scale
- **Polynomial kernels**

$$K(x, \tilde{x}) = \alpha(1 + x^T \tilde{x})^u.$$

[https:](https://distill.pub/2019/visual-exploration-gaussian-processes/)

[//distill.pub/2019/visual-exploration-gaussian-processes/](https://distill.pub/2019/visual-exploration-gaussian-processes/)

# Bayesian nonparametric model

Gaussian Process regression assumes the following model:

$$\begin{aligned}y_i \mid f &\stackrel{\text{ind}}{\sim} N(f(x_i), \sigma^2) \\ f &\sim GP(0, K)\end{aligned}$$

Note that  $x_i$  may be vector valued.

**Key result:** the posterior is a Gaussian process.

In the following we assume  $\sigma$  is known and fixed.



## Finding the posterior: case $\sigma = 0$

$\sigma = 0$  implies  $y_i = f_i$ , with  $f_i = f(x_i)$ . Set

$$\mathbf{f} = \begin{bmatrix} f_1 & \dots & f_n \end{bmatrix}^T.$$

Suppose we are interested in  $f_1^*, \dots, f_L^*$ , where  $f_k^* = f(x_k^*)$ . Set

$$\mathbf{f}^* = \begin{bmatrix} f_1^* & \dots & f_L^* \end{bmatrix}^T.$$

By definition of the Gaussian process:

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}^* \end{bmatrix} \sim N \left( \begin{bmatrix} \boldsymbol{\mu} \\ \boldsymbol{\mu}^* \end{bmatrix}, \begin{bmatrix} \mathcal{K} & R \\ R^T & \mathcal{K}^* \end{bmatrix} \right)$$

$$\mu_i = \mu(x_i) \quad \mu^*(x_k) = \mu(x_k)$$

$$\mathcal{K}_{ij} = K(x_i, x_j) \quad \mathcal{K}_{k\ell}^* = K(x_k^*, x_\ell^*) \quad R_{ik} = K(x_i, x_k^*).$$

Note that the matrix  $K$  is denoted by  $C$  in RG.

*Conditionals of the Gaussian distribution are Gaussian:*

$$\mathbf{f}^* \mid \mathbf{f} \sim N(\mu_{\text{post}}^*, \Sigma_{\text{post}}^*)$$

with

$$\mu_{\text{post}}^* = \mu^* + R^T \mathcal{K}^{-1}(\mathbf{f} - \mu) \quad \Sigma_{\text{post}}^* = \mathcal{K}^* - R^T \mathcal{K}^{-1} R$$

Note that the formulas in RG consider the special case where  $\mu = \mu^* = 0$ .

## Case $\sigma \neq 0$

Basic idea

$$\begin{aligned} p(\mathbf{f}^* | \mathbf{y}) &= \int p(\mathbf{f}^*, \mathbf{f} | \mathbf{y}) d\mathbf{f} \\ &= \int p(\mathbf{f}^* | \mathbf{f}) p(\mathbf{f} | \mathbf{y}) d\mathbf{f} \end{aligned}$$

combined with

$$p(\mathbf{f} | \mathbf{y}) \propto p(\mathbf{y} | \mathbf{f}) p(\mathbf{f}).$$

*All densities appearing are Gaussian.*

Lengthy calculation (Section 8.2.5 in RG) gives

$$\mathbf{f}^* | \mathbf{y} \sim N(\mu_{\text{post}}^*, \Sigma_{\text{post}}^*)$$

with

$$\mu_{\text{post}}^* = \mu^* + R^T (\mathcal{K} + \sigma^2 I)^{-1} (\mathbf{y} - \mu) \quad \Sigma_{\text{post}}^* = \mathcal{K}^* - R^T (\mathcal{K} + \sigma^2 I)^{-1} R$$

- Incorporating noise can be done directly in the kernel.
- Formulas look somewhat unwieldy. Need careful numerical linear algebra.
- Computationally intensive part: inversion of  $n \times n$  matrix  $\mathcal{K}$ . This is  $\mathcal{O}(n^3)$ .
- Parameters in the kernel can be dealt with by either
  - use empirical Bayes to estimate these;
  - employ a prior on these and do MCMC.

Want to do this in Python? See for example

[https://juanitorduz.github.io/gaussian\\_process\\_reg/](https://juanitorduz.github.io/gaussian_process_reg/) and  
<https://andrewcharlesjones.github.io/posts/2020/11/gaussian-processes/>

## **Advanced topics (not part of the exam)**

---

Map 1-dimensional variable  $x$  to  $u(x) := (\cos(x), \sin(x))$  and use the squared exponential kernel to get

$$\begin{aligned} K(x, \tilde{x}) &= v_0 \exp \left( -\frac{\|u(x) - u(\tilde{x})\|^2}{2\ell^2} \right) \\ &= v_0 \exp \left( \frac{2 \sin^2((x - \tilde{x})/2)}{\ell^2} \right). \end{aligned}$$

This follows from

$$\|u(x) - u(\tilde{x})\|^2 = 4 \sin^2 \left( \frac{x - \tilde{x}}{2} \right).$$

# Building kernels from simpler kernels

Slightly out of scope for this course, but

- sum of kernels  $k_1(x, \tilde{x}) + k_2(x, \tilde{x})$  gives a kernel
- product of kernels  $k_1(x, \tilde{x})k_2(x, \tilde{x})$  gives a kernel
- direct sums of kernels  $k(x, \tilde{x}) = k_1(x_1, \tilde{x}_1) + k_2(x_2, \tilde{x}_2)$  gives a kernel
- tensor product  $k(x, \tilde{x}) = k_1(x_1, \tilde{x}_1)k_2(x_2, \tilde{x}_2)$  gives a kernel

## A Special kernel: Automatic Relevance Detection (ARD)

We can adjust the square exponential kernel to the **non-isotropic** case.

Let  $M$  be PSD and set

$$K(x, \tilde{x}) = v_0 \exp \left( -\frac{(x - \tilde{x})^T M (x - \tilde{x})}{2\ell^2} \right).$$

In particular, if  $M = \text{diag}(\gamma_1, \dots, \gamma_p)$  then (taking  $\ell = 1$ )

$$K(x, \tilde{x}) = v_0 \exp \left( -\sum_{j=1}^p \gamma_j (x_j - \tilde{x}_j)^2 \right).$$

Different dimensions in  $x$  get different weights! Idea goes back to MacKay and Neal (1993).

Nice Idea: include  $\gamma_1, \dots, \gamma_p$  as parameter and estimate these.



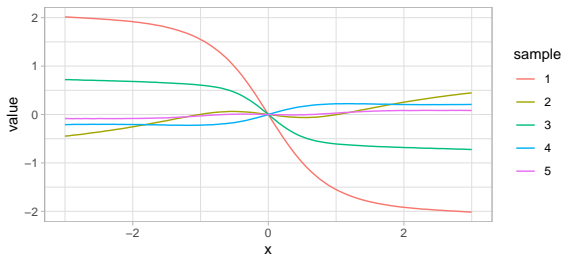
# A Special kernel: neuronal kernel

$$k(\mathbf{x}, \tilde{\mathbf{x}}) = \frac{2}{\pi} \arcsin \left( \frac{a(\mathbf{x}, \tilde{\mathbf{x}})}{\sqrt{(1 + a(\mathbf{x}, \mathbf{x}))(1 + a(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}))}} \right)$$

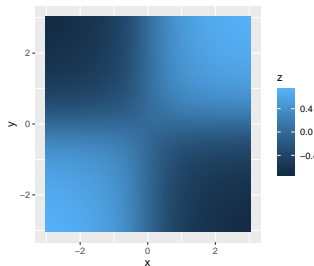
where

$$a(\mathbf{x}, \tilde{\mathbf{x}}) = 2\tilde{\mathbf{x}}^T \Sigma \mathbf{x}.$$

Draws neuronal(0.5) kernel



Cov struct neuronal(0.5) kernel



## One layer perceptron network: neuronal kernel. RG section 8.2.6.3.

Assume  $y \in \mathbb{R}$  and

$$y = f(\mathbf{x})$$

$$f(x) = \epsilon + \sum_{j=1}^N v_j h(\mathbf{x}; u_j), \quad h(x, u) = \operatorname{erf}(\mathbf{x}^T u)$$

$$u_1, \dots, u_N \stackrel{\text{ind}}{\sim} N(0, \Sigma),$$

where  $\operatorname{erf}(z) = 1/\sqrt{\pi} \int_{-z}^z e^{-t^2} dt$ .

- $v_j$ : hidden-to-output weights
- $h(\mathbf{x}, u_j)$ : hidden unit transfer functions

# One layer perceptron network

Assume  $\epsilon$ ,  $\{v_j\}$  are independent, zero mean, with variances  $\sigma_\epsilon^2$  and  $\sigma_v^2$  respectively.

$$\mathbb{E}f(\mathbf{x}) = \mathbb{E}[f(\mathbf{x}) \mid \{u_j\}] = 0.$$

$$\begin{aligned}\text{Cov}(f(\mathbf{x}), f(\tilde{\mathbf{x}})) &= \mathbb{E} [f(\mathbf{x})f(\tilde{\mathbf{x}})] \\ &= \mathbb{E} [\mathbb{E} [f(\mathbf{x})f(\tilde{\mathbf{x}}) \mid \{u_j\}_j]] \\ &= \sigma_\epsilon^2 + N\sigma_v^2 \mathbb{E}_u [h(\mathbf{x}, u)h(\tilde{\mathbf{x}}, u)] \\ &= \sigma_\epsilon^2 + \omega^2 \mathbb{E}_u [h(\mathbf{x}, u)h(\tilde{\mathbf{x}}, u)],\end{aligned}$$

if  $N\sigma_v^2 = \omega^2$ .

# Kernel for one layer perceptron network

Remains to compute

$$k(\mathbf{x}, \tilde{\mathbf{x}}) := \mathbb{E}_u [h(\mathbf{x}, u)h(\tilde{\mathbf{x}}, u)]$$

with  $h(\mathbf{x}, u) = \text{erf}(\mathbf{x}^T u)$  and  $u \sim N(0, \Sigma)$

Then

$$k(\mathbf{x}, \tilde{\mathbf{x}}) = \frac{2}{\pi} \arcsin \left( \frac{a(\mathbf{x}, \tilde{\mathbf{x}})}{\sqrt{(1 + a(\mathbf{x}, \mathbf{x}))(1 + a(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}))}} \right)$$

where

$$a(\mathbf{x}, \tilde{\mathbf{x}}) = 2\tilde{\mathbf{x}}^T \Sigma \mathbf{x}.$$

## Something cool, *advanced*

- A popular class of kernels are called Matérn covariance functions:

$$K(x, y) = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left( \sqrt{2\nu} \frac{\tau}{\ell} \right)^\nu \mathcal{K}_\nu \left( \frac{\tau}{\ell} \right),$$

where  $\mathcal{K}_\nu$  is the modified Bessel function,  $\tau = \|x - y\|$ , and  $\nu, \sigma, \ell$  are the smoothness, magnitude and length scale parameters.

- Realisations can be obtained also as solution to a Stochastic Differential Equation (SDE).
- For example, if  $\nu = 3/2$  then

$$df(t) = \begin{bmatrix} 0 & 1 \\ -\lambda^2 & -2\lambda \end{bmatrix} f(t) dt + \begin{bmatrix} 0 \\ 1 \end{bmatrix} dW_t,$$

with  $\lambda = \sqrt{2\nu}/\ell$ .

- Use Kalman smoothing to solve the interpolation problem: much lower computational cost! Not specific to Matérn covariance functions.

# Hyperparameter estimation by empirical Bayes

---

# Different ways to determine hyperparameters

1. Fully Bayesian: Put a hyperprior on the hyperparameters.
2. Empirical Bayes: Compute the probability of the model given the data.
3. Compute the generalisation error (the average error on unseen test examples).
4. Bound the generalisation error.

# Compute the probability of the model given the data

Assume  $\sigma$  fixed. As all is conditional on  $X$ , drop it from the notation.

$$\begin{aligned}y \mid f &\sim N(f, \sigma^2 I) \\ f &\sim N(0, \mathcal{K})\end{aligned}$$

Then  $y \sim N(0, \sigma^2 I + \mathcal{K})$ .

“Check”:

$$\begin{aligned}\mathbb{E} y &= \mathbb{E} \mathbb{E}[y \mid f] = \mathbb{E} f = 0 \\ \text{Cov } y &= \mathbb{E} \text{Cov}(y \mid f) + \text{Cov } \mathbb{E}[y \mid f] = \sigma^2 I + \mathcal{K}.\end{aligned}$$



Marginal likelihood:

$$p(y; \sigma^2, \eta) = (2\pi)^{-\frac{n}{2}} |\det(\mathcal{K}_y)|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} y^T \mathcal{K}_y^{-1} y\right),$$
$$\log p(y; \sigma^2, \eta) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\det(\mathcal{K}_y)| - \frac{1}{2} y^T \mathcal{K}_y^{-1} y$$

with  $\mathcal{K}_y = \sigma^2 I + \mathcal{K}$  and  $\eta$  the parameter vector of the kernel  $K$ .

Gradient methods can be used to optimise with respect to  $(\sigma^2, \eta)$ . This is an **empirical Bayes** method.

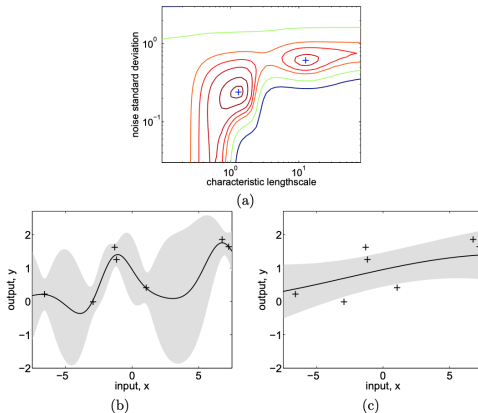


Figure 5.5: Panel (a) shows the marginal likelihood as a function of the hyperparameters  $\ell$  (length-scale) and  $\sigma_n^2$  (noise standard deviation), where  $\sigma_f^2 = 1$  (signal standard deviation) for a data set of 7 observations (seen in panels (b) and (c)). There are two local optima, indicated with '+': the global optimum has low noise and a short length-scale; the local optimum has a high noise and a long length-scale. In (b) and (c) the inferred underlying functions (and 95% confidence intervals) are shown for each of the two solutions. In fact, the data points were generated by a Gaussian process with  $(\ell, \sigma_f^2, \sigma_n^2) = (1, 1, 0.1)$  in eq. (5.1).

Figure taken Rasmussen & Williams, Gaussian Processes for Machine Learning, the MIT Press, 2006.

## **Gaussian process classification (RG section 8.3)**

---

# Gaussian process classification

- Let  $\psi$  denote the logistic function. Assume data  $\{(x_i, y_i), i = 1, \dots, n\}$  with  $y_i \in \{0, 1\}$ .
- Let  $X$  denote the design matrix.

Assume

$$\begin{aligned} y_1, \dots, y_n \mid f, X &\stackrel{\text{ind}}{\sim} \text{Ber}(\psi(f(x_i))) \\ f &\sim \text{GP}(0, K) \end{aligned}$$

(note that  $K$  is the kernel function).

# Prediction under Gaussian process classification

Suppose new data  $x_{\text{new}} = x^*$  and we wish to predict  $y_{\text{new}} = y^*$ .

Let  $\mathbf{f} = (f_1, \dots, f_n)$  with  $f_i = f(x_i)$ ,  $y = (y_1, \dots, y_n)$  and  $f^* = f(x^*)$ .

## Decomposition 1.

$$p(y^* | x^*, X, y) = \int p(y^* | f^*, x^*, X, y) p(f^* | x^*, X, y) df^*.$$

## Decomposition 2.

$$p(f^* | x^*, X, y) = \int p(f^* | \mathbf{f}, x^*, X, y) p(\mathbf{f} | x^*, X, y) d\mathbf{f}$$

(blue things can be “erased” by conditional independence). To sample from the predictive density:

- (A) Sample  $\mathbf{f}$  from the posterior  $p(\mathbf{f} | X, y)$ .
- (B) Sample  $f^*$  from  $p(f^* | \mathbf{f}, x^*)$ .
- (C) Sample  $y^* \sim \text{Ber}(\psi(f^*))$ .

## Step A: Sample $f$ from the posterior $p(f \mid X, y)$

The likelihood is

$$p(y \mid X, f) = \prod_{i=1}^n \psi(f_i)^{y_i} (1 - \psi(f_i))^{1-y_i}$$

The prior is

$$\mathbf{f} \sim N_n(0, \mathcal{K}) \quad \text{with} \quad \mathcal{K}_{ij} = K(x_i, x_j).$$

*The structure of the problem is exactly as in “ordinary” logistic regression; only the prior covariance is more sophisticated.*

In ordinary logistic regression we would have  $f_i = \theta^T x_i$  and prior  $\theta \sim N_p(0, \Sigma)$ .

## Step B: Sample $f^*$ from $p(f^* \mid \mathbf{f}, x^*)$

Note that (notation is somewhat sloppy here)

$$\begin{bmatrix} \mathbf{f} \\ f^* \end{bmatrix} \sim N \left( \begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} (K(x_i, x_j))_{i,j} & (K(x_i, x^*))_i \\ (K(x_i, x^*))_i & K(x^*, x^*) \end{bmatrix} \right).$$

This implies  $f^* \mid \mathbf{f}, x^*$  is normally distributed.

See Comment 8.2 in RG. *This is why the multivariate normal distribution is so convenient.*

## Step C: Sample $Y^* \sim \text{Ber}(\psi(f^*))$

This is really simple: once we have  $f^*$  we compute  $\psi(f^*)$  and draw from the Bernoulli distribution.



## Remarks on computational effort

- Step A cannot be done in closed form. For Bayesian logistic regression one option is to use the MH-algorithm.
- Note that  $\mathbf{f} \in \mathbb{R}^n$ , so we need the MH-algorithm in dimension  $n$ , rather than dimension  $p$  (if  $\theta \in \mathbb{R}^p$  for ordinary logistic regression).
- MCMC in high dimension is a very active topic of research.
- For step B, we need to invert the  $n \times n$  matrix  $(K(x_i, x_j))_{i,j}$ . The computational cost is  $\mathcal{O}(n^3)$ .

## **Gibbs sampling for GP classification (RG section 9.2)**

---

## Gibbs sampling for GP classification (RG section 9.2)

Assume the **probit-regression** model:

$$\begin{aligned} y_1, \dots, y_n \mid f, X &\stackrel{\text{ind}}{\sim} \text{Ber}(\Phi(f(x_i))) \\ f &\sim \text{GP}(0, K) \end{aligned}$$

with

$$\Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx.$$

With  $\Phi$  as *inverse-link* function we can avoid tuning MH-steps needed in step A. Instead: use Gibbs sampling.

## Turning the model in parametric form

Rewrite the model, using  $f_i := f(x_i)$ :

$$\begin{aligned} y_1, \dots, y_n \mid f, X &\stackrel{\text{ind}}{\sim} \text{Ber}(\Phi(f_i)) \\ f &\sim N(0, \mathcal{K}) \end{aligned}$$

with  $\mathcal{K}_{i,j} = K(x_i, x_j)$ .

First objective is to sample from  $f \mid y$ , where  $y = (y_1, \dots, y_n)$ .

## Introducing auxiliary variables

Trick: introduce auxiliary variables  $z_1, \dots, z_n$  such that

$$\begin{aligned} y_i &= \mathbf{1}\{z_i \geq 0\} \\ z_i \mid f_i &\stackrel{\text{ind}}{\sim} N(f_i, 1). \end{aligned}$$

This gives the correct Bernoulli distribution for  $y_i \mid f_i$  since

$$P(y_i = 1 \mid f_i) = P(z_i \geq 0 \mid f_i) = P\left(\underbrace{z_i - f_i}_{N(0,1)} \geq -f_i \mid f_i\right) = \Phi(f_i).$$

[Use that for  $Z \sim N(0, 1)$  we have  $P(Z \geq -z) = P(Z \leq z) = \Phi(z)$ .]

In the following let  $z = (z_1, \dots, z_n)$  and let  $\varphi_n(x; \mu, \Sigma)$  denote the density of the  $N(\mu, \Sigma)$ -distribution, evaluated at  $x$ .

# Gibbs sampler (1/2)

Structure:

$$f_i \longrightarrow z_i \longrightarrow y_i.$$

Gibbs sampling for

$$\begin{aligned} p(z, f \mid y) &\propto p(y, z, f) = p(y \mid z)p(z \mid f)p(f) \\ &= p(f)p(z \mid f) \prod_{i=1}^n p(y_i \mid z_i). \end{aligned}$$

- $p(f) = \varphi_n(f; 0, \mathcal{K})$ , where  $\mathcal{K}_{ij} = K(x_i, x_j)$   
( $K$  is kernel from Gaussian process),
- $p(z \mid f) = \prod_{i=1}^n \varphi(z_i; f_i, 1) = \varphi_n(z; f, I_n)$ ,
- $p(y_i \mid z_i) = \alpha_i^{y_i} (1 - \alpha_i)^{1-y_i}$  with  $\alpha_i = \mathbf{1}_{\{z_i \geq 0\}}$   
(because  $y_i \sim \text{Ber}(\alpha_i)$ ).

## Gibbs sampler (2/2)

- Updating  $f$ :

$$p(f \mid y, z) \propto \varphi_n(f; 0, \mathcal{K}) \varphi_n(z; f, I_n).$$

Conjugate step: amounts to drawing from the multivariate normal distribution.<sup>1</sup>

- Updating  $z_i$  can be done independently:

$$p(z_i \mid y, f) \propto \alpha_i^{y_i} (1 - \alpha_i)^{1-y_i} \varphi(z_i; f_i, 1)$$

This means

- If  $y_i = 1$ , draw  $z_i$  from  $N(f_i, 1)$ -distribution, conditioned to be positive.
- If  $y_i = 0$ , draw  $z_i$  from  $N(f_i, 1)$ -distribution, conditioned to be negative.

---

<sup>1</sup>In fact, conditional on  $z$ ,  $y$  plays no role and  $f \mid y, z \sim N(\mu_f, \Sigma_f)$ , where  $\mu_f = \Sigma_f z$  and  $\Sigma_f = (\mathcal{K}^{-1} + I)^{-1}$ .