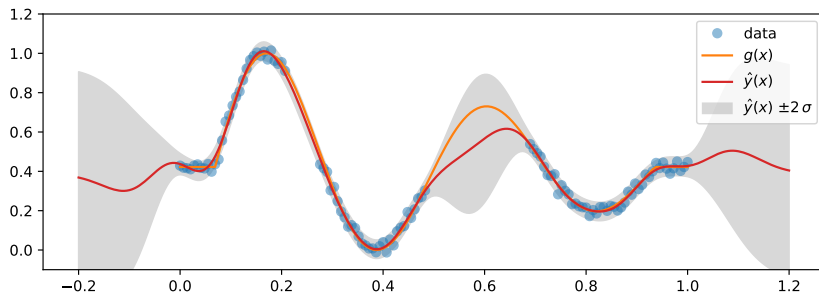


Introduction to Gaussian processes

Steve Schmerler

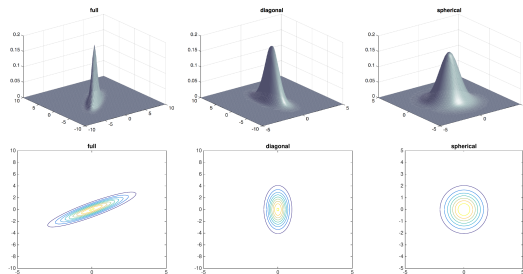
Helmholtz AI @HZDR

Motivation: Why GPs?



- ▶ interpolation or regression for low-dimensional problems ("smoothing device")
- ▶ **predictive uncertainty**
- ▶ building block for Bayesian optimization
- ▶ Bayesian stats and Gaussian process (GP) theory: understand uncertainty quantification (UQ) methods for neural networks (NNs)
- ▶ infinite width limits of NNs: neural network Gaussian process (NNGP) and the neural tangent kernel (NTK)
- ▶ two derivations: weight space, function space

Preliminaries: multivariate normal distribution



$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \text{cov}[x_1, x_2] \\ \text{cov}[x_1, x_2] & \sigma_2^2 \end{bmatrix}$$

Preliminaries: linear models

Linear model (parametric: $\dim \mathbf{w} = D \neq N$, data set content "compressed" into \mathbf{w})

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} = w_1 x_1 + w_2 x_2 + \cdots$$

$$f(\mathbf{x}) = \mathbf{w}^\top [1, \mathbf{x}] = w_0 + w_1 x_1 + w_2 x_2 + \cdots$$

Only regression models of the form

$$f : \mathbb{R}^D \rightarrow \mathbb{R}$$

Preliminaries: linear models

Linear model (parametric: $\dim \mathbf{w} = D \neq N$, data set content "compressed" into \mathbf{w})

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} = w_1 x_1 + w_2 x_2 + \dots$$

$$f(\mathbf{x}) = \mathbf{w}^\top [1, \mathbf{x}] = w_0 + w_1 x_1 + w_2 x_2 + \dots$$

Only regression models of the form

$$f: \mathbb{R}^D \rightarrow \mathbb{R}$$

Data set

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N = (\mathbf{X}, \mathbf{y})$$

$$\mathbf{x}_i \in \mathcal{X} = \mathbb{R}^D$$

$$y_i \in \mathcal{Y} = \mathbb{R}$$

$$\mathbf{X} \in \mathbb{R}^{N \times D}$$

Design matrix

$$\mathbf{X} = \overbrace{\begin{bmatrix} - & \mathbf{x}_1^\top & - \\ - & \mathbf{x}_2^\top & - \\ & \vdots & \\ - & \mathbf{x}_N^\top & - \end{bmatrix}}^D \in \mathbb{R}^{N \times D}$$

Preliminaries: linear models

Linear model (parametric: $\dim \mathbf{w} = D \neq N$, data set content "compressed" into \mathbf{w})

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} = w_1 x_1 + w_2 x_2 + \dots$$

$$f(\mathbf{x}) = \mathbf{w}^\top [1, \mathbf{x}] = w_0 + w_1 x_1 + w_2 x_2 + \dots$$

Only regression models of the form

$$f : \mathbb{R}^D \rightarrow \mathbb{R}$$

Notation

(noisy) data/target/label y

model output (train) $\mathbf{f} = \mathbf{w}^\top \mathbf{x}, \mathbf{f} = \mathbf{X} \mathbf{w}$

model output (test) $\mathbf{f}_* = \mathbf{w}^\top \mathbf{x}_*, \mathbf{f}_* = \mathbf{X}_* \mathbf{w}$

Data set

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N = (\mathbf{X}, \mathbf{y})$$

$$\mathbf{x}_i \in \mathcal{X} = \mathbb{R}^D$$

$$y_i \in \mathcal{Y} = \mathbb{R}$$

$$\mathbf{X} \in \mathbb{R}^{N \times D}$$

Design matrix

$$\mathbf{X} = \overbrace{\begin{bmatrix} - & \mathbf{x}_1^\top & - \\ - & \mathbf{x}_2^\top & - \\ & \vdots & \\ - & \mathbf{x}_N^\top & - \end{bmatrix}}^D \in \mathbb{R}^{N \times D}$$

Basis functions

Feature space mapping

$$\phi : \mathcal{X} \rightarrow \mathcal{F}$$

$$f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$$

$f(\mathbf{x})$ is nonlinear in \mathbf{x} but still linear in \mathbf{w}

$$\mathbf{X} = \begin{bmatrix} - & \mathbf{x}_1^\top & - \\ - & \mathbf{x}_2^\top & - \\ & \vdots & \\ - & \mathbf{x}_N^\top & - \end{bmatrix} \rightarrow \Phi = \begin{bmatrix} - & \phi(\mathbf{x}_1)^\top & - \\ - & \phi(\mathbf{x}_2)^\top & - \\ & \vdots & \\ - & \phi(\mathbf{x}_N)^\top & - \end{bmatrix}$$

Basis functions

Feature space mapping

$$\phi : \mathcal{X} \rightarrow \mathcal{F}$$

$$f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$$

$f(\mathbf{x})$ is nonlinear in \mathbf{x} but still linear in \mathbf{w}

$$\mathbf{X} = \begin{bmatrix} - & \mathbf{x}_1^\top & - \\ - & \mathbf{x}_2^\top & - \\ & \vdots & \\ - & \mathbf{x}_N^\top & - \end{bmatrix} \rightarrow \Phi = \begin{bmatrix} - & \phi(\mathbf{x}_1)^\top & - \\ - & \phi(\mathbf{x}_2)^\top & - \\ & \vdots & \\ - & \phi(\mathbf{x}_N)^\top & - \end{bmatrix}$$

Example: polynomial basis: $\mathbf{x} \in \mathbb{R}^2$, $\mathcal{F} = \mathbb{R}^5$, $\mathbf{w}, \phi(\mathbf{x}) \in \mathbb{R}^5$

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1^2, x_2^2]$$

$$f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2$$

`sklearn.preprocessing.PolynomialFeatures`

Kernels

Kernel function $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ as similarity measure

- ▶ symmetric: $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \kappa(\mathbf{x}_j, \mathbf{x}_i)$
- ▶ positive: $\kappa(\mathbf{x}_i, \mathbf{x}_j) \geq 0$

Kernels

Kernel function $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ as similarity measure

► symmetric: $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \kappa(\mathbf{x}_j, \mathbf{x}_i)$

► positive: $\kappa(\mathbf{x}_i, \mathbf{x}_j) \geq 0$

Gram matrix $\mathbf{K} \in \mathbb{R}^{N \times N}$

$$K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

$$\mathbf{K} := \kappa(\mathbf{X}, \mathbf{X})$$

Kernels

Kernel function $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ as similarity measure

► symmetric: $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \kappa(\mathbf{x}_j, \mathbf{x}_i)$

► positive: $\kappa(\mathbf{x}_i, \mathbf{x}_j) \geq 0$

Gram matrix $\mathbf{K} \in \mathbb{R}^{N \times N}$

"Kernel trick"

$$K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \equiv \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$

$$\mathbf{K} := \kappa(\mathbf{X}, \mathbf{X})$$

Rich theory (Reproducing kernel Hilbert space, Mercer's theorem, ...): no need to define ϕ explicitly, sufficient to define $\kappa(\cdot, \cdot)$, for certain κ we have $f(\mathbf{x}) = \sum_{i=1}^{\infty} w_i \phi_i(\mathbf{x})$

Kernels

Kernel function $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ as similarity measure

► symmetric: $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \kappa(\mathbf{x}_j, \mathbf{x}_i)$

► positive: $\kappa(\mathbf{x}_i, \mathbf{x}_j) \geq 0$

Gram matrix $\mathbf{K} \in \mathbb{R}^{N \times N}$ "Kernel trick"

$$K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \equiv \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$

$$\mathbf{K} := \kappa(\mathbf{X}, \mathbf{X})$$

Rich theory (Reproducing kernel Hilbert space, Mercer's theorem, ...): no need to define ϕ explicitly, sufficient to define $\kappa(\cdot, \cdot)$, for certain κ we have $f(\mathbf{x}) = \sum_{i=1}^{\infty} w_i \phi_i(\mathbf{x})$

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_j$$

Linear / dot product kernel

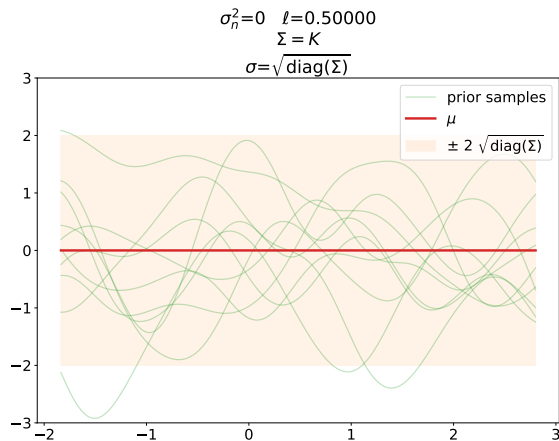
$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\ell^2}\right) = \begin{cases} 1 & \mathbf{x}_i = \mathbf{x}_j \\ < 1 & \text{else} \end{cases}$$

Gaussian/RBF/"squared exponential"
kernel, characteristic length scale ℓ

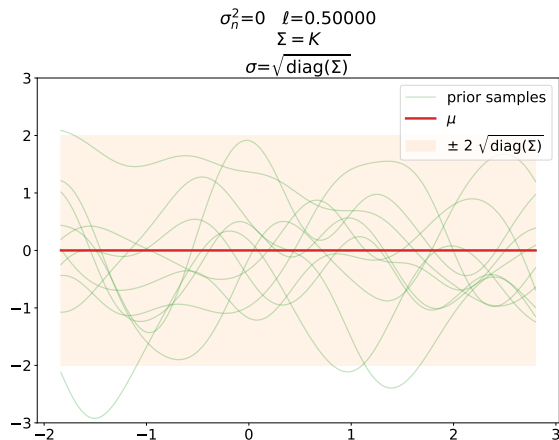
The GP prior for the RBF kernel, fixed ℓ

1D example where

$$\mathbf{x} = x \in \mathbb{R}, \mathbf{X} \in \mathbb{R}^{N \times 1}$$



The GP prior for the RBF kernel, fixed ℓ



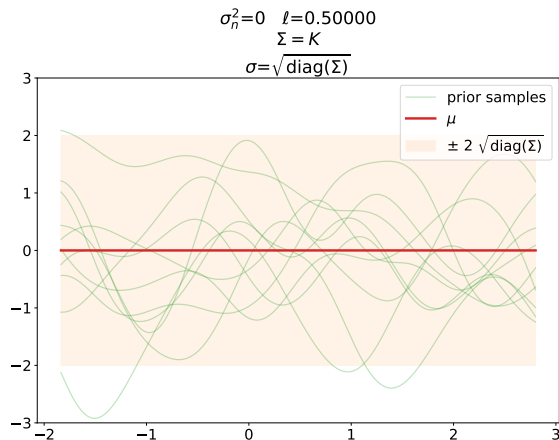
1D example where

$$\mathbf{x} = x \in \mathbb{R}, \mathbf{X} \in \mathbb{R}^{N \times 1}$$

GP prior for model $f = f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}}) \quad \text{weight prior}$$

The GP prior for the RBF kernel, fixed ℓ



1D example where

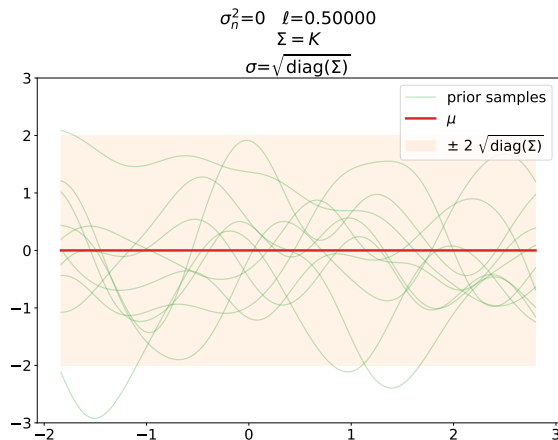
$$\mathbf{x} = x \in \mathbb{R}, \mathbf{X} \in \mathbb{R}^{N \times 1}$$

GP prior for model $f = f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}}) \quad \text{weight prior}$$

$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$$

The GP prior for the RBF kernel, fixed ℓ



1D example where

$$\mathbf{x} = x \in \mathbb{R}, \mathbf{X} \in \mathbb{R}^{N \times 1}$$

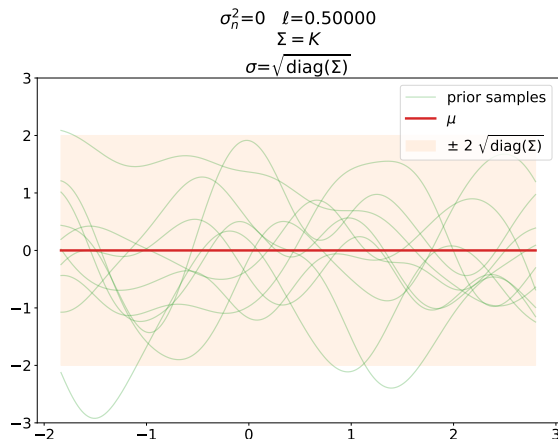
GP prior for model $f = f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}}) \quad \text{weight prior}$$

$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$$

$$\mathbb{E}[\mathbf{f}] = \mathbb{E}[\Phi \mathbf{w}] = \Phi \mathbb{E}[\mathbf{w}] = \mathbf{0}$$

The GP prior for the RBF kernel, fixed ℓ



1D example where

$$\mathbf{x} = x \in \mathbb{R}, \mathbf{X} \in \mathbb{R}^{N \times 1}$$

GP prior for model $\mathbf{f} = f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$

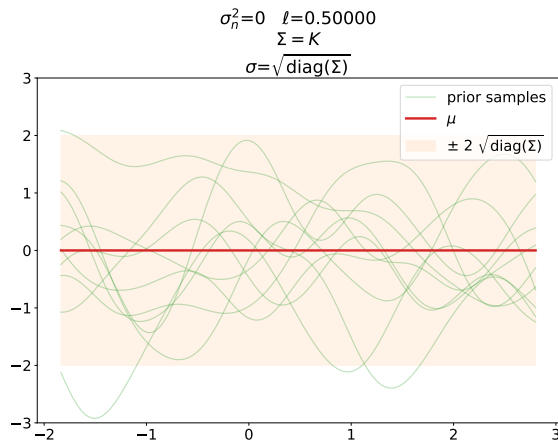
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}}) \quad \text{weight prior}$$

$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$$

$$\mathbb{E}[\mathbf{f}] = \mathbb{E}[\Phi \mathbf{w}] = \Phi \mathbb{E}[\mathbf{w}] = \mathbf{0}$$

$$\begin{aligned} \text{cov}[\mathbf{f}] &= \mathbb{E}[(\mathbf{f} - \mathbb{E}[\mathbf{f}]) (\mathbf{f} - \mathbb{E}[\mathbf{f}])^\top] \\ &= \Phi \Sigma_{\mathbf{w}} \Phi^\top =: \mathbf{K} \end{aligned}$$

The GP prior for the RBF kernel, fixed ℓ



1D example where

$$\mathbf{x} = x \in \mathbb{R}, \mathbf{X} \in \mathbb{R}^{N \times 1}$$

GP prior for model $f = f(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}}) \quad \text{weight prior}$$

$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$$

$$\mathbb{E}[\mathbf{f}] = \mathbb{E}[\Phi \mathbf{w}] = \Phi \mathbb{E}[\mathbf{w}] = \mathbf{0}$$

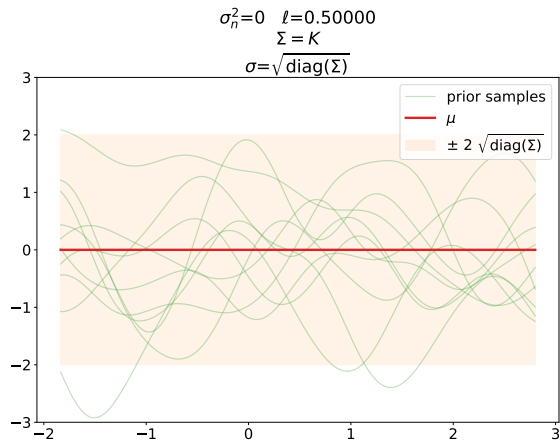
$$\begin{aligned} \text{cov}[\mathbf{f}] &= \mathbb{E}[(\mathbf{f} - \mathbb{E}[\mathbf{f}]) (\mathbf{f} - \mathbb{E}[\mathbf{f}])^\top] \\ &= \Phi \Sigma_{\mathbf{w}} \Phi^\top =: \mathbf{K} \end{aligned}$$

Covariance (kernel) function $\kappa(\cdot, \cdot)$

$$K_{ij} = \phi(\mathbf{x}_i)^\top \Sigma_{\mathbf{w}} \phi(\mathbf{x}_j) =: \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

e.g. $\Sigma_{\mathbf{w}} = \tau^2 \mathbf{I}_D \rightarrow$ scaling factor in κ

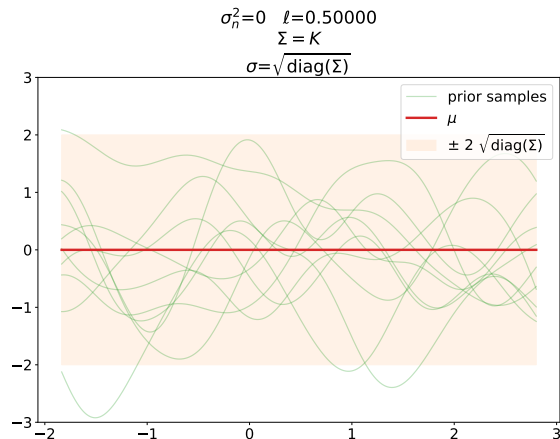
Function space view: the GP prior



The GP as a distribution over *functions* f

$$f \sim \mathcal{GP}(m(\cdot), \kappa(\cdot, \cdot))$$

Function space view: the GP prior



The GP as a distribution over *functions* f

$$f \sim \mathcal{GP}(m(\cdot), \kappa(\cdot, \cdot))$$

$$p(\mathbf{f}_i | \mathbf{x}_i) = \mathcal{N}(\mathbf{m}(\mathbf{x}_i), \kappa(\mathbf{x}_i, \mathbf{x}_i))$$

$$p(\mathbf{f} | \mathbf{X}) = \mathcal{N}(\mathbf{m}(\mathbf{X}), \mathbf{K})$$

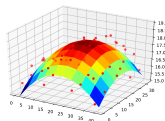
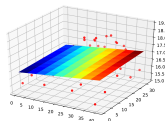
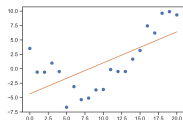
$$\mathbb{E}[\mathbf{f}_i] = \mathbf{m}(\mathbf{x}_i)$$

$$\mathbb{E}[\mathbf{f}] = \mathbf{m}(\mathbf{X})$$

$$\begin{aligned} \text{cov}[\mathbf{f}_i, \mathbf{f}_j] &= \mathbb{E}[(\mathbf{f}_i - \mathbf{m}(\mathbf{x}_i)) (\mathbf{f}_j - \mathbf{m}(\mathbf{x}_j))] \\ &=: \kappa(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

$$\text{cov}[\mathbf{f}] = \mathbf{K}$$

Likelihood



Model noise σ_n^2 in data y .

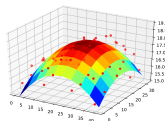
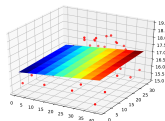
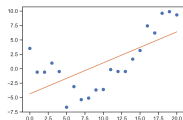
$p(y|\mathbf{x}, \mathbf{w})$ interpretation:

- ▶ distribution $p(y|\dots)$ over y
- ▶ function of \mathbf{w}

"The likelihood function reflects the data we expect to see for each setting of the parameters \mathbf{w} ."

$$p(y|\mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \phi(\mathbf{x}), \sigma_n^2)$$

Likelihood



Model noise σ_n^2 in data y .

$p(y|\mathbf{x}, \mathbf{w})$ interpretation:

- ▶ distribution $p(y|\dots)$ over y
- ▶ function of \mathbf{w}

"The likelihood function reflects the data we expect to see for each setting of the parameters \mathbf{w} ."

$$p(y|\mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \phi(\mathbf{x}), \sigma_n^2)$$

$$\mathbf{f} = \mathbf{w}^\top \phi(\mathbf{x})$$

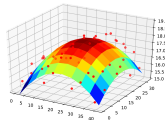
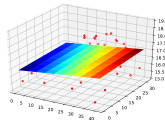
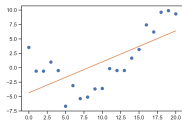
$$y = \mathbf{w}^\top \phi(\mathbf{x}) + \epsilon$$

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

$$\boldsymbol{\theta} = (\mathbf{w}, \boldsymbol{\xi})$$

$$\boldsymbol{\xi} = (\sigma_n^2) \quad \text{hyper parameters}$$

Likelihood



Model noise σ_n^2 in data y .

$p(y|\mathbf{x}, \mathbf{w})$ interpretation:

► distribution $p(y|\dots)$ over y

► function of \mathbf{w}

"The likelihood function reflects the data we expect to see for each setting of the parameters \mathbf{w} ."

$$p(y|\mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \phi(\mathbf{x}), \sigma_n^2)$$

$$\mathbf{f} = \mathbf{w}^\top \phi(\mathbf{x})$$

$$y = \mathbf{w}^\top \phi(\mathbf{x}) + \epsilon$$

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

$$\boldsymbol{\theta} = (\mathbf{w}, \boldsymbol{\xi})$$

$$\boldsymbol{\xi} = (\sigma_n^2) \quad \text{hyper parameters}$$

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\boldsymbol{\Phi} \mathbf{w}, \sigma_n^2 \mathbf{I}_N)$$

$$= \prod_{i=1}^N p(y_i|\mathbf{x}_i, \mathbf{w})$$

$$= \prod_{i=1}^N \frac{1}{\sqrt{2\pi} \sigma_n} \exp\left(-\frac{(y_i - f_i)^2}{2\sigma_n^2}\right)$$

$$= \frac{1}{\sqrt{(2\pi)^N \sigma_n^2}} \exp\left(-\frac{\boldsymbol{\epsilon}^\top \boldsymbol{\epsilon}}{2\sigma_n^2}\right)$$

Bayes' rule

Bayesian inference: infer *posterior* distribution over weights (i.e. models) $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$ by using training data (\mathbf{X}, \mathbf{y})

$$\overbrace{p(\mathbf{w}|\mathbf{X}, \mathbf{y})}^{\text{weight posterior}} = \frac{\overbrace{p(\mathbf{y}|\mathbf{X}, \mathbf{w})}^{\text{likelihood}} \overbrace{p(\mathbf{w})}^{\text{weight prior}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marginal likelihood or evidence}}} = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}}$$

Bayes' rule

Bayesian inference: infer *posterior* distribution over weights (i.e. models) $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$ by using training data (\mathbf{X}, \mathbf{y})

$$\overbrace{p(\mathbf{w}|\mathbf{X}, \mathbf{y})}^{\text{weight posterior}} = \frac{\overbrace{p(\mathbf{y}|\mathbf{X}, \mathbf{w})}^{\text{likelihood}} \overbrace{p(\mathbf{w})}^{\text{weight prior}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marginal likelihood or evidence}}} = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}}$$

More compact notation

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})} = \frac{p(\mathcal{D}|\mathbf{w}) p(\mathbf{w})}{\int p(\mathcal{D}|\mathbf{w}) p(\mathbf{w}) d\mathbf{w}}$$

In simple cases, inference can be performed analytically, e.g. for a Gaussian likelihood.

Posterior predictive distribution

Bayesian model averaging (BMA)

$$\langle w \rangle = \int w p(w) \, dw$$

$$\langle f(w) \rangle = \int f(w) p(w) \, dw$$

Posterior predictive distribution

Bayesian model averaging (BMA)

$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \int \underbrace{p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{w})}_{\text{likelihood}} \underbrace{p(\mathbf{w} | \mathbf{X}, \mathbf{y})}_{\text{weight posterior}} d\mathbf{w} = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$
$$\langle w \rangle = \int w p(w) dw$$
$$\langle f(w) \rangle = \int f(w) p(w) dw$$

Posterior predictive distribution

Bayesian model averaging (BMA)

$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \int \underbrace{p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{w})}_{\text{likelihood}} \underbrace{p(\mathbf{w} | \mathbf{X}, \mathbf{y})}_{\text{weight posterior}} d\mathbf{w} = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$
$$\langle w \rangle = \int w p(w) dw$$
$$\langle f(w) \rangle = \int f(w) p(w) dw$$

Predictive mean $\boldsymbol{\mu}_*$ and cov. $\boldsymbol{\Sigma}_*$

$$\begin{aligned}\boldsymbol{\mu}_* &= \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{y} \\ &= \mathbf{K}_* \boldsymbol{\alpha}\end{aligned}$$

$$\begin{aligned}\mathbf{K}_* &= \kappa(\mathbf{X}_*, \mathbf{X}) \\ \mathbf{K}_{**} &= \kappa(\mathbf{X}_*, \mathbf{X}_*)\end{aligned}$$

Posterior predictive distribution

Bayesian model averaging (BMA)

$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \int \underbrace{p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{w})}_{\text{likelihood}} \underbrace{p(\mathbf{w} | \mathbf{X}, \mathbf{y})}_{\text{weight posterior}} d\mathbf{w} = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$
$$\langle w \rangle = \int w p(w) dw$$
$$\langle f(w) \rangle = \int f(w) p(w) dw$$

Predictive mean $\boldsymbol{\mu}_*$ and cov. $\boldsymbol{\Sigma}_*$

$$\begin{aligned}\boldsymbol{\mu}_* &= \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{y} \\ &= \mathbf{K}_* \boldsymbol{\alpha}\end{aligned}$$

$$\mu_* = \sum_{j=1}^N \alpha_j \kappa(\mathbf{x}_*, \mathbf{x}_j)$$

$$\mathbf{K}_* = \kappa(\mathbf{X}_*, \mathbf{X})$$

$$\mathbf{K}_{**} = \kappa(\mathbf{X}_*, \mathbf{X}_*)$$

Posterior predictive distribution

Bayesian model averaging (BMA)

$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \int \underbrace{p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{w})}_{\text{likelihood}} \underbrace{p(\mathbf{w} | \mathbf{X}, \mathbf{y})}_{\text{weight posterior}} d\mathbf{w} = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$
$$\langle w \rangle = \int w p(w) dw$$
$$\langle f(w) \rangle = \int f(w) p(w) dw$$

Predictive mean $\boldsymbol{\mu}_*$ and cov. $\boldsymbol{\Sigma}_*$

$$\begin{aligned}\boldsymbol{\mu}_* &= \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{y} \\ &= \mathbf{K}_* \boldsymbol{\alpha}\end{aligned}$$

$$\mu_* = \sum_{j=1}^N \alpha_j \kappa(\mathbf{x}_*, \mathbf{x}_j)$$

$$\boldsymbol{\Sigma}_* = \mathbf{K}_{**} - \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{K}_*^\top$$

$$\mathbf{K}_* = \kappa(\mathbf{X}_*, \mathbf{X})$$

$$\mathbf{K}_{**} = \kappa(\mathbf{X}_*, \mathbf{X}_*)$$

Posterior predictive distribution

Bayesian model averaging (BMA)

$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \int \underbrace{p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{w})}_{\text{likelihood}} \underbrace{p(\mathbf{w} | \mathbf{X}, \mathbf{y})}_{\text{weight posterior}} d\mathbf{w} = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$
$$\langle w \rangle = \int w p(w) dw$$
$$\langle f(w) \rangle = \int f(w) p(w) dw$$

Predictive mean $\boldsymbol{\mu}_*$ and cov. $\boldsymbol{\Sigma}_*$

$$\begin{aligned}\boldsymbol{\mu}_* &= \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{y} \\ &= \mathbf{K}_* \boldsymbol{\alpha}\end{aligned}$$

$$\mu_* = \sum_{j=1}^N \alpha_j \kappa(\mathbf{x}_*, \mathbf{x}_j)$$

$$\boldsymbol{\Sigma}_* = \mathbf{K}_{**} - \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{K}_*^\top$$

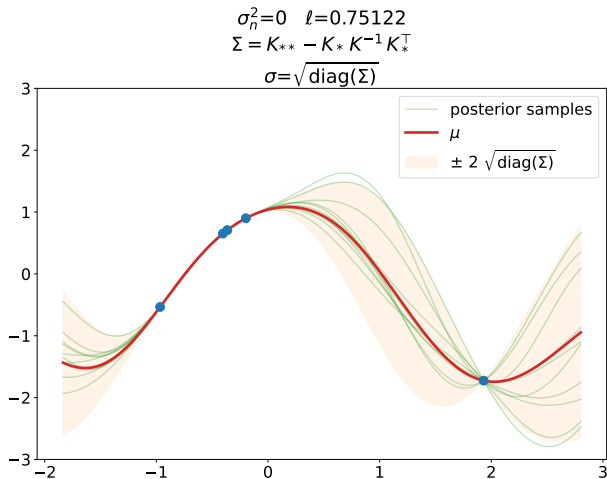
$$\mathbf{K}_* = \kappa(\mathbf{X}_*, \mathbf{X})$$

$$\mathbf{K}_{**} = \kappa(\mathbf{X}_*, \mathbf{X}_*)$$

Non-parametric model: $\boldsymbol{\mu} = \mathbf{K} \boldsymbol{\alpha}$

- ▶ $K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$
- ▶ $\mathbf{K} \in \mathbb{R}^{N \times N}$ contains info about whole training inputs $\mathbf{X} \in \mathbb{R}^{N \times D}$
- ▶ weights $\boldsymbol{\alpha} \in \mathbb{R}^N$ contain info about (\mathbf{X}, \mathbf{y})
- ▶ large data sets (large N) make vanilla GPs costly

Posterior predictive with $\sigma_n^2 = 0$



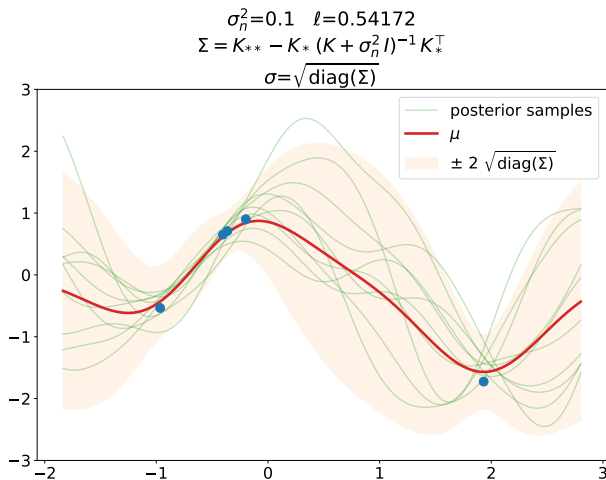
$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$

$$\begin{aligned} \boldsymbol{\mu}_* &= \mathbf{K}_* \mathbf{K}^{-1} \mathbf{y} \\ &= \mathbf{K}_* \boldsymbol{\alpha} \end{aligned}$$

$$\boldsymbol{\mu}_* = \sum_j \alpha_j \kappa(\mathbf{x}_*, \mathbf{x}_j)$$

$$\boldsymbol{\Sigma}_* = \mathbf{K}_{**} - \mathbf{K}_* \mathbf{K}^{-1} \mathbf{K}_*^\top$$

Posterior predictive with $\sigma_n^2 > 0$



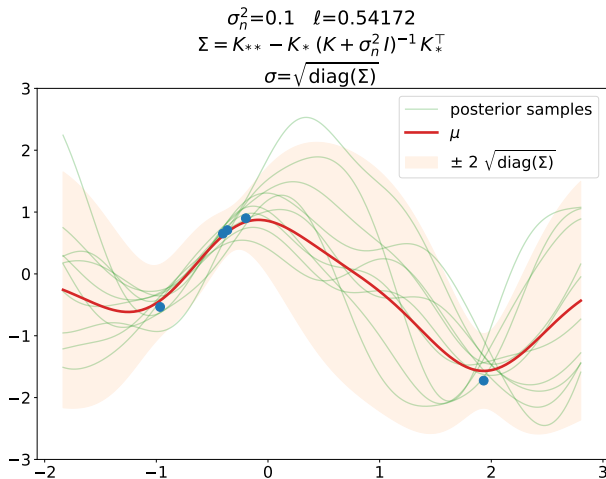
$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$

$$\begin{aligned} \boldsymbol{\mu}_* &= \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{y} \\ &= \mathbf{K}_* \boldsymbol{\alpha} \end{aligned}$$

$$\mu_* = \sum_j \alpha_j \kappa(\mathbf{x}_*, \mathbf{x}_j)$$

$$\boldsymbol{\Sigma}_* = \mathbf{K}_{**} - \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{K}_*^\top$$

Posterior predictive with $\sigma_n^2 > 0$



$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$

$$\begin{aligned} \boldsymbol{\mu}_* &= \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{y} \\ &= \mathbf{K}_* \boldsymbol{\alpha} \end{aligned}$$

$$\mu_* = \sum_j \alpha_j \kappa(\mathbf{x}_*, \mathbf{x}_j)$$

$$\boldsymbol{\Sigma}_* = \mathbf{K}_{**} - \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{K}_*^\top$$

Data noise σ_n^2 : transform
interpolation \rightarrow regression, same
effect as a regularization term in NN
training

Function space view of GPs: the joint

We rewrite the prior: Divide data into "train" \mathbf{f} and "test/prediction" \mathbf{f}_*

$$M = N + N_*$$

$$\mathbf{X} \in \mathbb{R}^{M \times D} \rightarrow (\mathbf{X} \in \mathbb{R}^{N \times D}, \mathbf{X}_* \in \mathbb{R}^{N_* \times D})$$

$$\mathbf{f} \in \mathbb{R}^M \rightarrow (\mathbf{f} \in \mathbb{R}^N, \mathbf{f}_* \in \mathbb{R}^{N_*})$$

Function space view of GPs: the joint

We rewrite the prior: Divide data into "train" \mathbf{f} and "test/prediction" \mathbf{f}_*

$$M = N + N_*$$

$$\mathbf{X} \in \mathbb{R}^{M \times D} \rightarrow (\mathbf{X} \in \mathbb{R}^{N \times D}, \mathbf{X}_* \in \mathbb{R}^{N_* \times D})$$

$$\mathbf{f} \in \mathbb{R}^M \rightarrow (\mathbf{f} \in \mathbb{R}^N, \mathbf{f}_* \in \mathbb{R}^{N_*})$$

write the prior as joint over concat. $(\mathbf{f}, \mathbf{f}_*)$

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}(\mathbf{X}) \\ \mathbf{m}(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} \text{cov}[\mathbf{f}] & \mathbf{K}_*^\top \\ \mathbf{K}_* & \text{cov}[\mathbf{f}_*] \end{bmatrix} \right)$$

$$\sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}(\mathbf{X}) \\ \mathbf{m}(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} \mathbf{K} & \mathbf{K}_*^\top \\ \mathbf{K}_* & \mathbf{K}_{**} \end{bmatrix} \right)$$

$$\sim p(\mathbf{f}, \mathbf{f}_* | \mathbf{X}, \mathbf{X}_*)$$

Function space view of GPs: the joint

We rewrite the prior: Divide data into "train" \mathbf{f} and "test/prediction" \mathbf{f}_*

$$M = N + N_*$$

$$\mathbf{X} \in \mathbb{R}^{M \times D} \rightarrow (\mathbf{X} \in \mathbb{R}^{N \times D}, \mathbf{X}_* \in \mathbb{R}^{N_* \times D})$$

$$\mathbf{f} \in \mathbb{R}^M \rightarrow (\mathbf{f} \in \mathbb{R}^N, \mathbf{f}_* \in \mathbb{R}^{N_*})$$

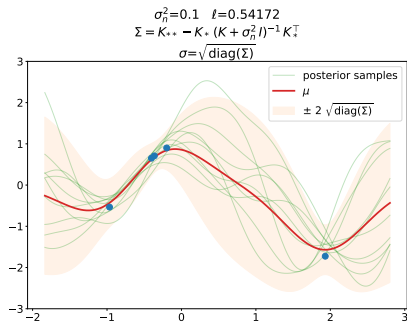
write the prior as joint over concat. $(\mathbf{f}, \mathbf{f}_*)$

$$\begin{aligned} \begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} &\sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}(\mathbf{X}) \\ \mathbf{m}(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} \text{cov}[\mathbf{f}] & \mathbf{K}_*^\top \\ \mathbf{K}_* & \text{cov}[\mathbf{f}_*] \end{bmatrix} \right) \\ &\sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}(\mathbf{X}) \\ \mathbf{m}(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} \mathbf{K} & \mathbf{K}_*^\top \\ \mathbf{K}_* & \mathbf{K}_{**} \end{bmatrix} \right) \\ &\sim p(\mathbf{f}, \mathbf{f}_* | \mathbf{X}, \mathbf{X}_*) \end{aligned}$$

For noisy $\mathbf{y} = \mathbf{f} + \boldsymbol{\epsilon}$, we have

$$\begin{aligned} \begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} &\sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}(\mathbf{X}) \\ \mathbf{m}(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} \text{cov}[\mathbf{y}] & \mathbf{K}_*^\top \\ \mathbf{K}_* & \text{cov}[\mathbf{f}_*] \end{bmatrix} \right) \\ &\sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}(\mathbf{X}) \\ \mathbf{m}(\mathbf{X}_*) \end{bmatrix}, \begin{bmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I}_N & \mathbf{K}_*^\top \\ \mathbf{K}_* & \mathbf{K}_{**} \end{bmatrix} \right) \\ &\sim p(\mathbf{y}, \mathbf{f}_* | \mathbf{X}, \mathbf{X}_*) \end{aligned}$$

Function space view of GPs: posterior predictive



Transform the joint $p(\mathbf{y}, \mathbf{f}_* | \mathbf{X}, \mathbf{X}_*)$ into the posterior predictive $p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y})$ by conditioning on (\mathbf{X}, \mathbf{y}) ("training data").

$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$$

$$\boldsymbol{\mu}_* = \mathbf{m}(\mathbf{X}_*) + \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} (\mathbf{y} - \mathbf{m}(\mathbf{X}))$$

$$\boldsymbol{\Sigma}_* = \text{cov}[\mathbf{f}_*] = \mathbf{K}_{**} - \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{K}_*^\top$$

Same result as the posterior predictive obtained from Bayes' rule + model averaging. Here we also have a mean function $\mathbf{m}(\cdot) \neq 0$.

GP hyperparameter optimization

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

Bayes' rule

$$\begin{aligned} \overbrace{p(\mathbf{w}|\mathbf{X}, \mathbf{y})}^{\text{weight posterior}} &= \frac{\overbrace{p(\mathbf{y}|\mathbf{X}, \mathbf{w})}^{\text{likelihood}} \overbrace{p(\mathbf{w})}^{\text{weight prior}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marginal likelihood or evidence}}} \\ &= \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}} \end{aligned}$$

GP hyperparameter optimization

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

Marginal likelihood as function of hyperparameters $\xi = (\ell, \sigma_n^2)$.

Bayes' rule

$$\begin{aligned} \overbrace{p(\mathbf{w}|\mathbf{X}, \mathbf{y})}^{\text{weight posterior}} &= \frac{\overbrace{p(\mathbf{y}|\mathbf{X}, \mathbf{w})}^{\text{likelihood}} \overbrace{p(\mathbf{w})}^{\text{weight prior}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marginal likelihood or evidence}}} \\ &= \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}} \end{aligned}$$

GP hyperparameter optimization

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

Marginal likelihood as function of hyperparameters $\boldsymbol{\xi} = (\ell, \sigma_n^2)$. Because of $\mathbf{f} = \Phi \mathbf{w}$, $\int \dots d\mathbf{w} \rightarrow \int \dots d\mathbf{f}$

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

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\xi}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{f}, \boldsymbol{\xi}) p(\mathbf{f}|\mathbf{X}, \boldsymbol{\xi}) d\mathbf{f}$$

Bayes' rule

$$\begin{aligned} \overbrace{p(\mathbf{w}|\mathbf{X}, \mathbf{y})}^{\text{weight posterior}} &= \frac{\overbrace{p(\mathbf{y}|\mathbf{X}, \mathbf{w})}^{\text{likelihood}} \overbrace{p(\mathbf{w})}^{\text{weight prior}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marginal likelihood or evidence}}} \\ &= \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}} \end{aligned}$$

GP hyperparameter optimization

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

Marginal likelihood as function of hyperparameters $\boldsymbol{\xi} = (\ell, \sigma_n^2)$. Because of $\mathbf{f} = \Phi \mathbf{w}$, $\int \dots d\mathbf{w} \rightarrow \int \dots d\mathbf{f}$

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

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\xi}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{f}, \boldsymbol{\xi}) p(\mathbf{f}|\mathbf{X}, \boldsymbol{\xi}) d\mathbf{f}$$

(negative) log marginal likelihood

$$\ln p(\mathbf{y}|\mathbf{X}, \boldsymbol{\xi}) = -\frac{1}{2} \left[\underbrace{\mathbf{y}^\top (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{y}}_{\text{model fit}} + \underbrace{\ln |\mathbf{K} + \sigma_n^2 \mathbf{I}_N|}_{\text{model complexity}} + N \ln(2\pi) \right]$$

Bayes' rule

$$\begin{aligned} \overbrace{p(\mathbf{w}|\mathbf{X}, \mathbf{y})}^{\text{weight posterior}} &= \frac{\overbrace{p(\mathbf{y}|\mathbf{X}, \mathbf{w})}^{\text{likelihood}} \overbrace{p(\mathbf{w})}^{\text{weight prior}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marginal likelihood or evidence}}} \\ &= \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}} \end{aligned}$$

GP hyperparameter optimization

Bayes' rule

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

Marginal likelihood as function of hyperparameters $\boldsymbol{\xi} = (\ell, \sigma_n^2)$. Because of $\mathbf{f} = \Phi \mathbf{w}$, $\int \dots d\mathbf{w} \rightarrow \int \dots d\mathbf{f}$

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

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\xi}) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{f}, \boldsymbol{\xi}) p(\mathbf{f}|\mathbf{X}, \boldsymbol{\xi}) d\mathbf{f}$$

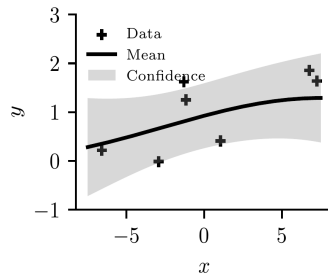
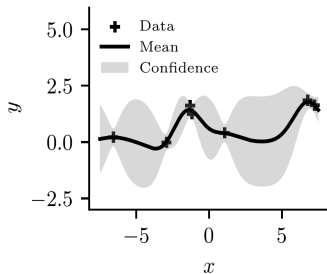
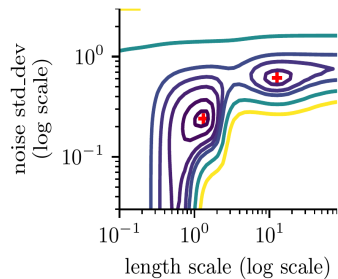
$$\begin{aligned} \overbrace{p(\mathbf{w}|\mathbf{X}, \mathbf{y})}^{\text{weight posterior}} &= \frac{\overbrace{p(\mathbf{y}|\mathbf{X}, \mathbf{w})}^{\text{likelihood}} \overbrace{p(\mathbf{w})}^{\text{weight prior}}}{\underbrace{p(\mathbf{y}|\mathbf{X})}_{\text{marginal likelihood or evidence}}} \\ &= \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}} \end{aligned}$$

(negative) log marginal likelihood

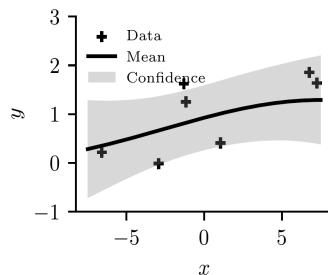
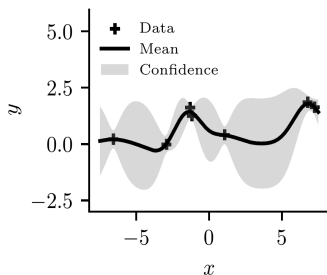
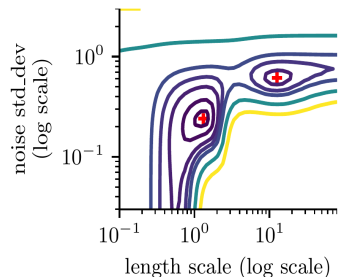
$$\ln p(\mathbf{y}|\mathbf{X}, \boldsymbol{\xi}) = -\frac{1}{2} \left[\underbrace{\mathbf{y}^\top (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{y}}_{\text{model fit}} + \underbrace{\ln |\mathbf{K} + \sigma_n^2 \mathbf{I}_N|}_{\text{model complexity}} + N \ln(2\pi) \right]$$

$$\boldsymbol{\xi}^* = \arg \max_{\boldsymbol{\xi}} \ln p(\mathbf{y}|\mathbf{X}, \boldsymbol{\xi}) = \arg \min_{\boldsymbol{\xi}} (-\ln p(\mathbf{y}|\mathbf{X}, \boldsymbol{\xi}))$$

Multiple minima of negative log marginal likelihood



Multiple minima of negative log marginal likelihood



Multiple minima: explain data in different ways

- ▶ small length scale ℓ , flexible model, low variance $\sigma_n^2 \rightarrow$ good model fit but complex model
- ▶ large length scale ℓ , "stiff"/low curvature model, high variance $\sigma_n^2 \rightarrow$ worse model fit but low model complexity

Relation to uncertainty quantification

Different kinds of uncertainty:

- ▶ epistemic / model uncertainty: weight posterior $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$ and
$$\text{cov}[\mathbf{f}_*] = \mathbf{\Sigma}_* = \mathbf{K}_{**} - \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{K}_*^\top$$
- ▶ aleatoric / data uncertainty: σ_n^2

Relation to uncertainty quantification

Different kinds of uncertainty:

- ▶ epistemic / model uncertainty: weight posterior $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$ and $\text{cov}[\mathbf{f}_*] = \Sigma_* = \mathbf{K}_{**} - \mathbf{K}_* (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{K}_*^\top$
- ▶ aleatoric / data uncertainty: σ_n^2

Distinction between "noise-free/noiseless prediction" of \mathbf{f}_* when using Σ_*

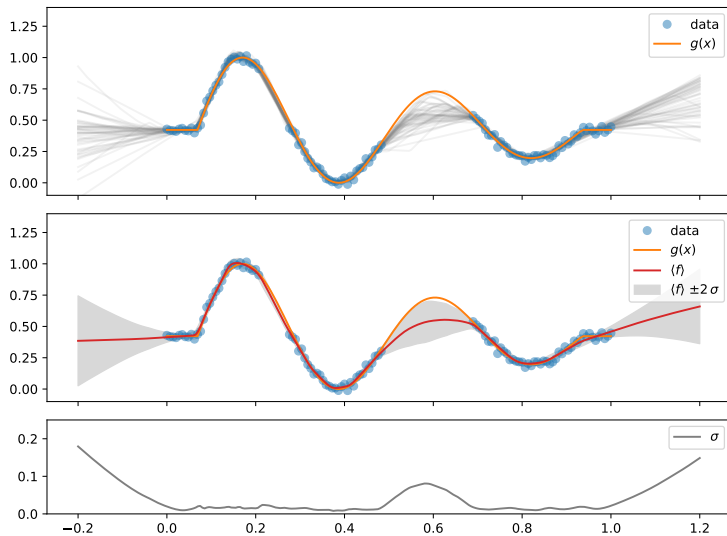
$$p(\mathbf{f}_*|\mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_*, \Sigma_*)$$

and "noisy predictions" of \mathbf{y}_* where we use $\Sigma_* + \sigma_n^2 \mathbf{I}_N$

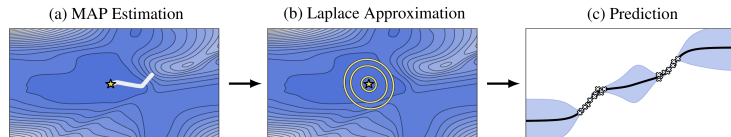
$$p(\mathbf{y}_*|\mathbf{X}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_*, \Sigma_* + \sigma_n^2 \mathbf{I}_N)$$

In both cases we have the same $\boldsymbol{\mu}_*$

Approximate $p(\mathbf{w}|\mathcal{D})$: NN ensembles for UQ



Approximate $p(\mathbf{w}|\mathcal{D})$: Laplace approximation for UQ



Post-processing step after NN training (= MAP estimate): $\mathbf{w}^* = \arg \min_{\mathbf{w}} -\ln p(\mathbf{w}|\mathcal{D})$

$$-\ln p(\mathbf{w}|\mathcal{D}) = -\ln \left(\frac{p(\mathcal{D}|\mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})} \right) = \overbrace{-\ln p(\mathcal{D}|\mathbf{w}) - \ln p(\mathbf{w})}^{\text{NN loss } L(\mathbf{w}) = \text{NLL} + \text{regularizer}} + \ln p(\mathcal{D})$$

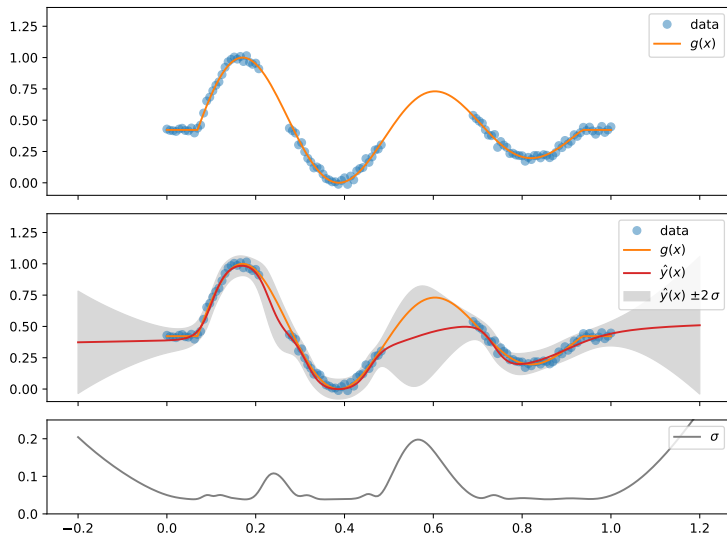
With gradient $\mathbf{g} = \nabla L|_{\mathbf{w}^*}$, Hessian $\mathbf{H} = \partial^2 L|_{\mathbf{w}^*}$ and $\mathbf{h} = \mathbf{w} - \mathbf{w}^*$, approximate loss to 2nd order

$$L(\mathbf{w}) \approx L(\mathbf{w}^*) + \underbrace{\mathbf{g}^\top}_{=0} \mathbf{h} + \frac{1}{2} \mathbf{h}^\top \mathbf{H} \mathbf{h}$$

Approximate posterior probability distribution over \mathbf{w} (i.e. over models)

$$p(\mathbf{w}|\mathcal{D}) \approx \mathcal{N}(\mathbf{w}^*, \mathbf{\Sigma}) \quad \text{where } \mathbf{\Sigma} = \mathbf{H}^{-1}$$

Approximate $p(\mathbf{w}|\mathcal{D})$: Laplace approximation for UQ

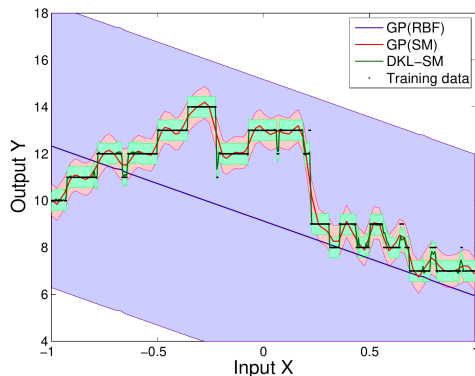


Kernel learning with NNs

(deep) kernel learning: more flexible
kernels via NNs: use base kernel + NN
features:

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp \left(-\frac{\|\mathbf{h}_\gamma(\mathbf{x}_i) - \mathbf{h}_\gamma(\mathbf{x}_j)\|_2^2}{2\ell^2} \right)$$

with $\mathbf{h}_\gamma(\mathbf{x}_i)$ an NN embedding ("feature extractor") and γ the NN parameters (weights, biases), optimize $\xi = (\gamma, \ell, \sigma_n^2)$ jointly using log marginal likelihood



Software

- ▶ <https://scikit-learn.org>, uses *numpy*
 - ▶ `sklearn.gaussian_process.GaussianProcessRegressor`
 - ▶ `sklearn.kernel_ridge.KernelRidge`
- ▶ <https://gpytorch.ai>: *PyTorch*-based, lots of advanced features, approximate methods for scaling GPs, API flexible but complex, GPU support via *PyTorch*
 - ▶ Define a mean function, since $f \sim \mathcal{GP}(m(\cdot), \kappa(\cdot, \cdot))$, so far we had $m(\cdot) = 0$
 - ▶ GPs are *non-parametric* models, $\mathbf{K} \in \mathbb{R}^{N \times N}$, $\dim \boldsymbol{\alpha} = N$, Cholesky decomposition for $\boldsymbol{\alpha} = (\mathbf{K} + \sigma_n^2 \mathbf{I}_N)^{-1} \mathbf{y}$ is $\mathcal{O}(N^3)$, **lots** of approximate methods, such as KISS-GP (a.k.a. SKI = structured kernel interpolation) for improved scaling
 - ▶ variational GPs for approximate inference of $p(\mathbf{w}|\mathcal{D})$, e.g. for non-Gaussian likelihoods
 - ▶ GP theory is for $f: \mathbb{R}^D \rightarrow \mathbb{R}$, *GPYtorch* supports multi-output GPs for $f: \mathbb{R}^D \rightarrow \mathbb{R}^M$
- ▶ <https://github.com/dfm/tinygp>: basic (educational) code, GPU support via *JAX*
- ▶ <https://github.com/JaxGaussianProcesses>, similar to *tinygp* but more features, GPU support via *JAX*
- ▶ <https://github.com/SheffieldML/GPy>, uses *numpy* + *Cython*
- ▶ Laplace approximation: <https://github.com/AlexImmer/Laplace>, *PyTorch*

Resources

- ▶ The Book: C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006 (<http://gaussianprocess.org/gpml>)
- ▶ K. P. Murphy. *Probabilistic Machine Learning: An introduction*. MIT Press, 2022, K. P. Murphy. *Probabilistic Machine Learning: Advanced Topics (draft version)*. MIT Press, 2022 (<https://probml.github.io/pml-book>)
- ▶ M. Kanagawa et al. *Gaussian Processes and Kernel Methods: A Review on Connections and Equivalences*. 2018. URL: <http://arxiv.org/abs/1807.02582>
- ▶ shameless plug: https://elcorto.github.io/gp_playground
- ▶ UQ in classification problems: P. Steinbach et al. “Machine Learning State-of-the-Art with Uncertainties”. In: *ICLR* (2022)
- ▶ J. Gawlikowski et al. *A Survey of Uncertainty in Deep Neural Networks*. 2022. URL: <http://arxiv.org/abs/2107.03342> (visited on 07/12/2022)