

DATA, ESTIMATION & INFERENCE

GAUSSIAN PROCESSES & BAYESIAN DEEP LEARNING



MICHAELMAS TERM, 2020
University of Oxford

Tim G. J. Rudner

`tim.rudner@cs.ox.ac.uk`

Syllabus: <https://www.notion.so/Data-Estimation-and-Inference-Part-2-d177bdb310e74f9aafd034c35652e4f3>

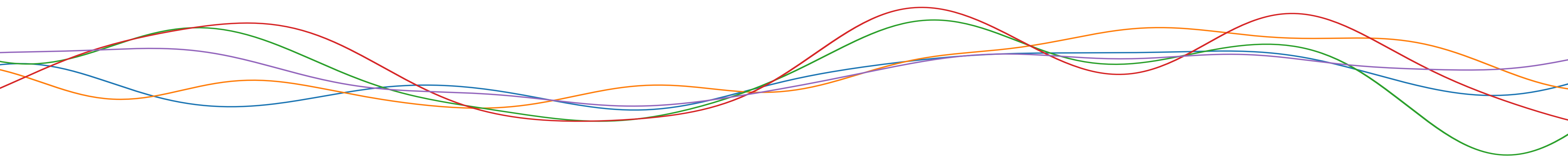
OVERVIEW

Syllabus:

<https://www.notion.so/Data-Estimation-and-Inference-Part-2-d177bdb310e74f9aafd034c35652e4f3>

Today:

- Weight-space function-space duality
- Basics of Gaussian processes (GPs)
- Implementing GPs (if time permits)



BAYESIAN INFERENCE IN MACHINE LEARNING

BAYESIAN INFERENCE IN MACHINE LEARNING

How to represent (stochastic/random) processes in the real world probabilistically?



BAYESIAN INFERENCE IN MACHINE LEARNING

How to represent (stochastic/random) processes in the real world probabilistically?



DISTRIBUTIONS OVER FUNCTIONS

DISTRIBUTIONS OVER FUNCTIONS

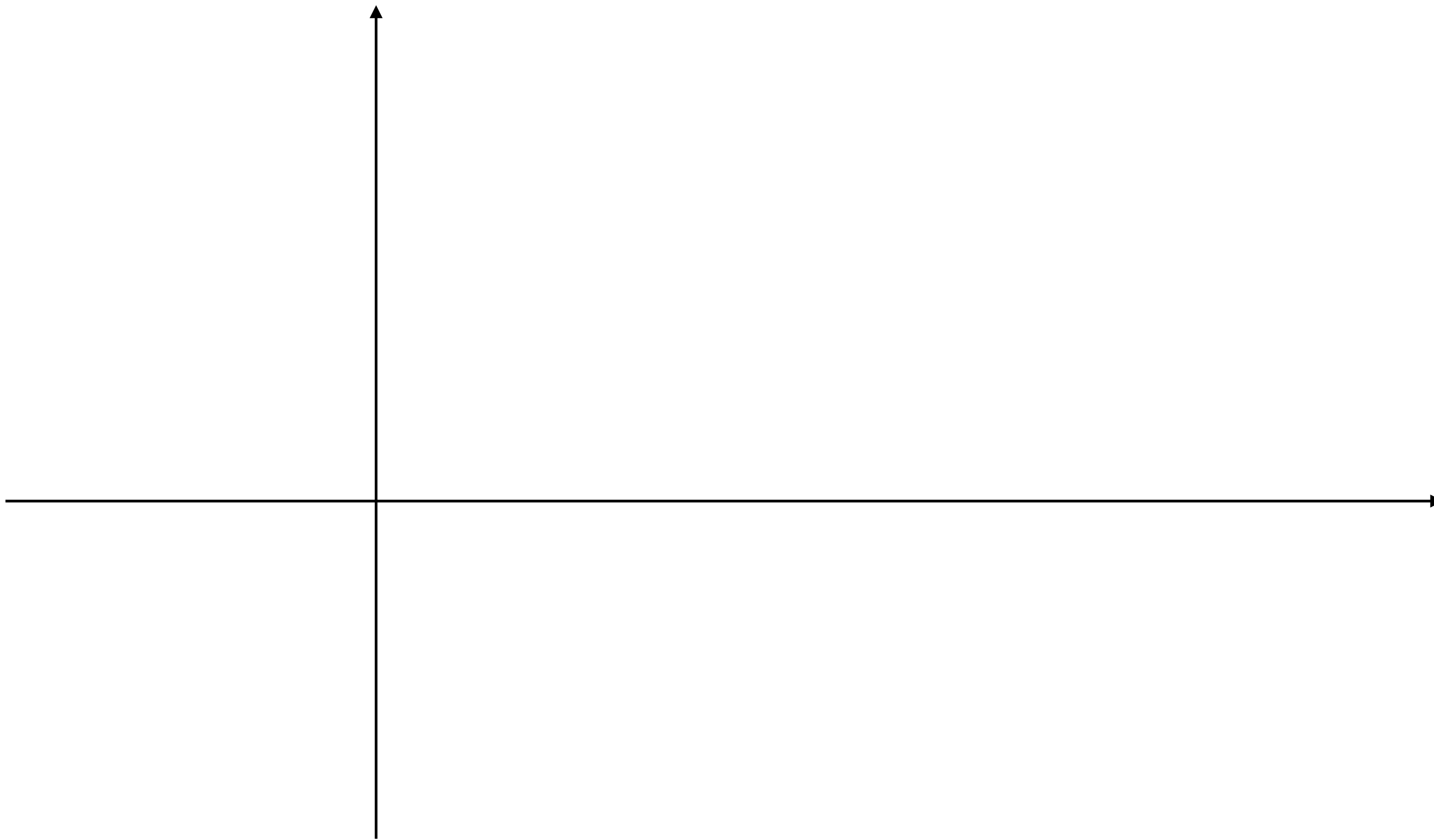
Q: What is a stochastic function?

Q: What is a distribution over functions?

Consider

$$f(x; w) = wx, \quad w \sim \mathcal{N}(w|0, 1)$$

[Link to visualization](#)



WEIGHT-SPACE FUNCTION-SPACE DUALITY

WEIGHT-SPACE FUNCTION-SPACE DUALITY

Probabilistic model:

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}, \quad y = f(\mathbf{x}) + \varepsilon,$$

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p) \quad \varepsilon \sim \mathcal{N}(0, \sigma_n^2)$$

Bayesian inference:

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}},$$

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

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

WEIGHT-SPACE FUNCTION-SPACE DUALITY

Probabilistic model:

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}, \quad y = f(\mathbf{x}) + \varepsilon,$$
$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p) \quad \varepsilon \sim \mathcal{N}(0, \sigma_n^2)$$

Bayesian inference:

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}, \quad p(\mathbf{w}|\mathbf{y}, X) = \frac{p(\mathbf{y}|X, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|X)}$$

Posterior over weights:

$$p(\mathbf{w}|X, \mathbf{y}) \propto \exp\left(-\frac{1}{2\sigma_n^2}(\mathbf{y} - X^\top \mathbf{w})^\top (\mathbf{y} - X^\top \mathbf{w})\right) \exp\left(-\frac{1}{2}\mathbf{w}^\top \Sigma_p^{-1} \mathbf{w}\right)$$
$$\propto \exp\left(-\frac{1}{2}(\mathbf{w} - \bar{\mathbf{w}})^\top \left(\frac{1}{\sigma_n^2} X X^\top + \Sigma_p^{-1}\right) (\mathbf{w} - \bar{\mathbf{w}})\right)$$
$$\bar{\mathbf{w}} = \sigma_n^{-2} (\sigma_n^{-2} X X^\top + \Sigma_p^{-1})^{-1} X \mathbf{y}$$

WEIGHT-SPACE FUNCTION-SPACE DUALITY

Probabilistic model: $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}, \quad y = f(\mathbf{x}) + \varepsilon,$
 $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p) \quad \varepsilon \sim \mathcal{N}(0, \sigma_n^2)$

Posterior over weights:

$$p(\mathbf{w}|X, \mathbf{y}) \propto \exp\left(-\frac{1}{2\sigma_n^2}(\mathbf{y} - X^\top \mathbf{w})^\top (\mathbf{y} - X^\top \mathbf{w})\right) \exp\left(-\frac{1}{2}\mathbf{w}^\top \Sigma_p^{-1} \mathbf{w}\right)$$

$$\propto \exp\left(-\frac{1}{2}(\mathbf{w} - \bar{\mathbf{w}})^\top \left(\frac{1}{\sigma_n^2} X X^\top + \Sigma_p^{-1}\right) (\mathbf{w} - \bar{\mathbf{w}})\right)$$

$$\bar{\mathbf{w}} = \sigma_n^{-2} (\sigma_n^{-2} X X^\top + \Sigma_p^{-1})^{-1} X \mathbf{y} \quad A = \sigma_n^{-2} X X^\top + \Sigma_p^{-1}$$

$$p(\mathbf{w}|X, \mathbf{y}) \sim \mathcal{N}\left(\bar{\mathbf{w}} = \frac{1}{\sigma_n^2} A^{-1} X \mathbf{y}, A^{-1}\right)$$

WEIGHT-SPACE FUNCTION-SPACE DUALITY

Probabilistic model: $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}, \quad y = f(\mathbf{x}) + \varepsilon,$
 $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p) \quad \varepsilon \sim \mathcal{N}(0, \sigma_n^2)$

Posterior over weights:

$$\bar{\mathbf{w}} = \sigma_n^{-2} (\sigma_n^{-2} X X^\top + \Sigma_p^{-1})^{-1} X \mathbf{y} \quad A = \sigma_n^{-2} X X^\top + \Sigma_p^{-1}$$

$$p(\mathbf{w} | X, \mathbf{y}) \sim \mathcal{N}(\bar{\mathbf{w}} = \frac{1}{\sigma_n^2} A^{-1} X \mathbf{y}, A^{-1})$$

Posterior predictive distribution: Let $f_* \triangleq f(\mathbf{x}_*)$ at \mathbf{x}_*

$$\begin{aligned} p(f_* | \mathbf{x}_*, X, \mathbf{y}) &= \int p(f_* | \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} | X, \mathbf{y}) d\mathbf{w} \\ &= \mathcal{N}\left(\frac{1}{\sigma_n^2} \mathbf{x}_*^\top A^{-1} X \mathbf{y}, \mathbf{x}_*^\top A^{-1} \mathbf{x}_*\right) \end{aligned}$$

WEIGHT-SPACE FUNCTION-SPACE DUALITY

Probabilistic model: $f(\mathbf{x}) = \phi(\mathbf{x})^\top \mathbf{w}$
 $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p) \quad \varepsilon \sim \mathcal{N}(0, \sigma_n^2)$

Posterior predictive distribution:

$$f_* | \mathbf{x}_*, X, \mathbf{y} \sim \mathcal{N}\left(\frac{1}{\sigma_n^2} \phi(\mathbf{x}_*)^\top A^{-1} \Phi \mathbf{y}, \phi(\mathbf{x}_*)^\top A^{-1} \phi(\mathbf{x}_*)\right)$$

with $\Phi = \Phi(X)$ and $A = \sigma_n^{-2} \Phi \Phi^\top + \Sigma_p^{-1}$

After application of Woodbury matrix inversion lemma:

$$f_* | \mathbf{x}_*, X, \mathbf{y} \sim \mathcal{N}\left(\phi_*^\top \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \mathbf{y}, \right. \\ \left. \phi_*^\top \Sigma_p \phi_* - \phi_*^\top \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \Phi^\top \Sigma_p \phi_*\right)$$

$$\phi(\mathbf{x}_*) = \phi_* \quad K = \Phi^\top \Sigma_p \Phi$$

WEIGHT-SPACE FUNCTION-SPACE DUALITY

Gaussian Processes:

Definition 2.1 A Gaussian process is a collection of random variables, any finite number of which have a joint Gaussian distribution. \square

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

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

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$

Bayesian linear model:

$$f(\mathbf{x}) = \phi(\mathbf{x})^\top \mathbf{w} \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p)$$

$$\mathbb{E}[f(\mathbf{x})] = \phi(\mathbf{x})^\top \mathbb{E}[\mathbf{w}] = 0,$$

$$\mathbb{E}[f(\mathbf{x})f(\mathbf{x}')] = \phi(\mathbf{x})^\top \mathbb{E}[\mathbf{w}\mathbf{w}^\top] \phi(\mathbf{x}') = \phi(\mathbf{x})^\top \Sigma_p \phi(\mathbf{x}')$$

WEIGHT-SPACE FUNCTION-SPACE DUALITY

Gaussian Process prior:

$$\mathbf{f}_* \sim \mathcal{N}(\mathbf{0}, K(X_*, X_*))$$

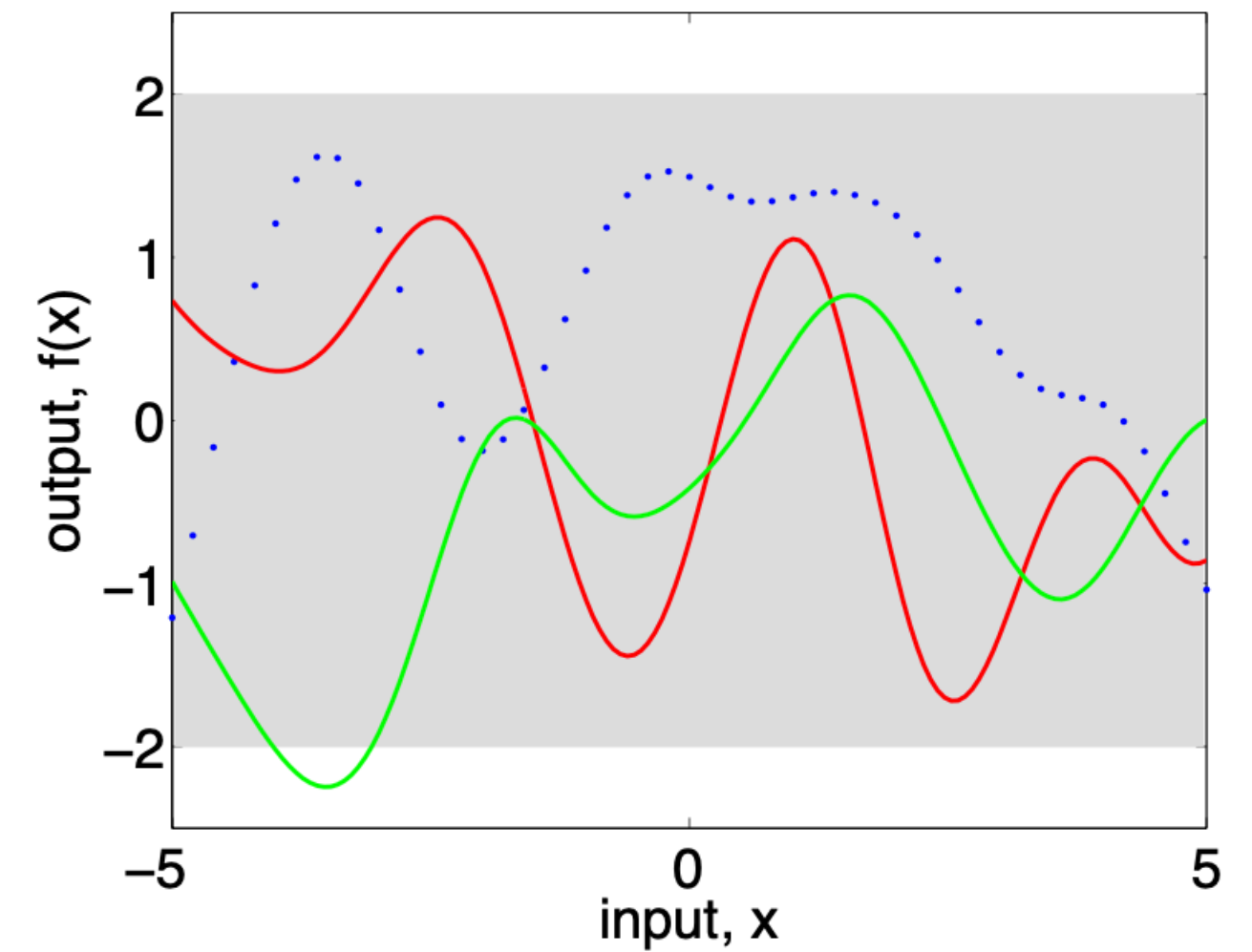
Gaussian Process joint:

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(X, X) & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right)$$

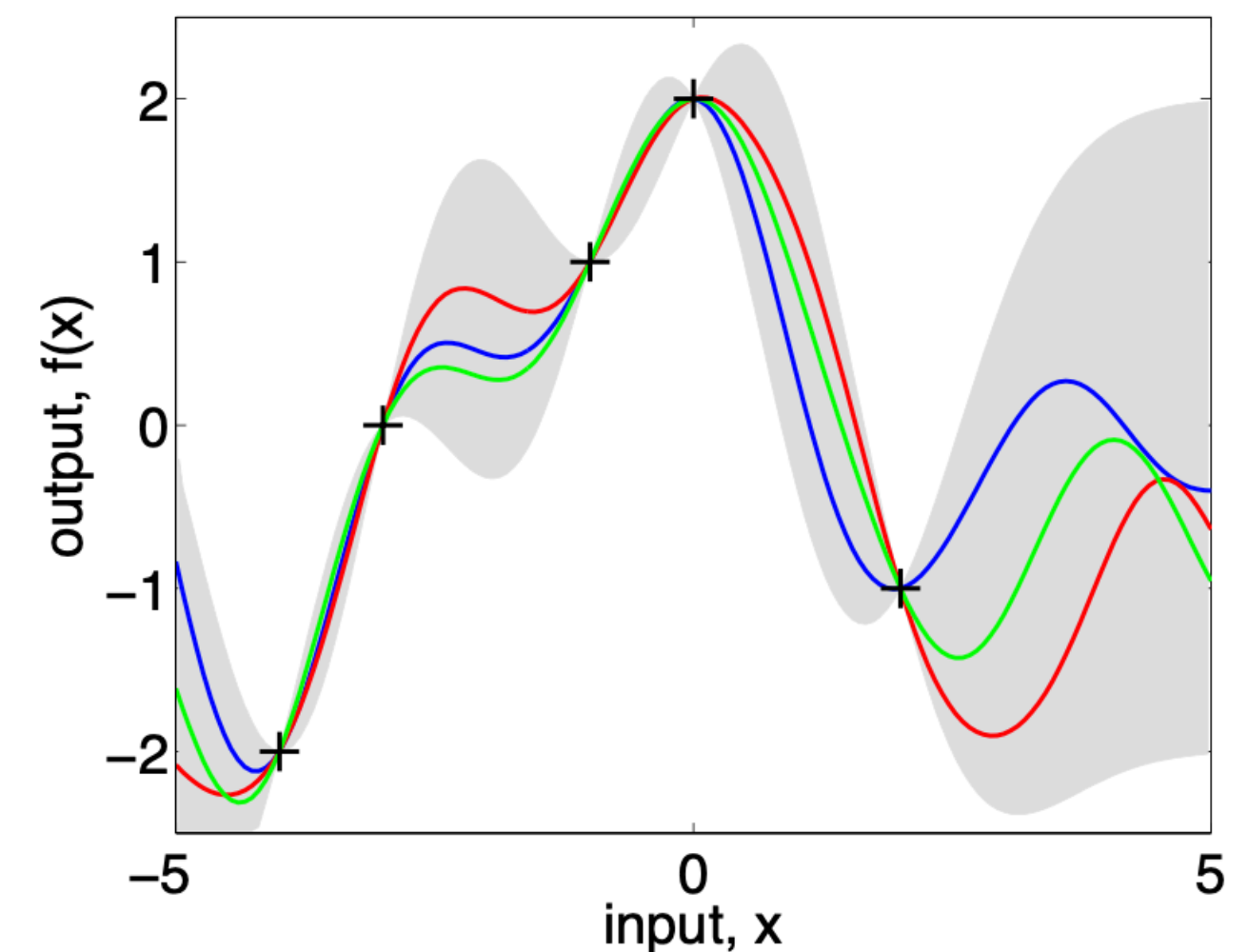
Gaussian Process posterior:

$$\mathbf{f}_* | X_*, X, \mathbf{f} \sim \mathcal{N}(K(X_*, X)K(X, X)^{-1}\mathbf{f},$$

$$K(X_*, X_*) - K(X_*, X)K(X, X)^{-1}K(X, X_*))$$



(a), prior



(b), posterior

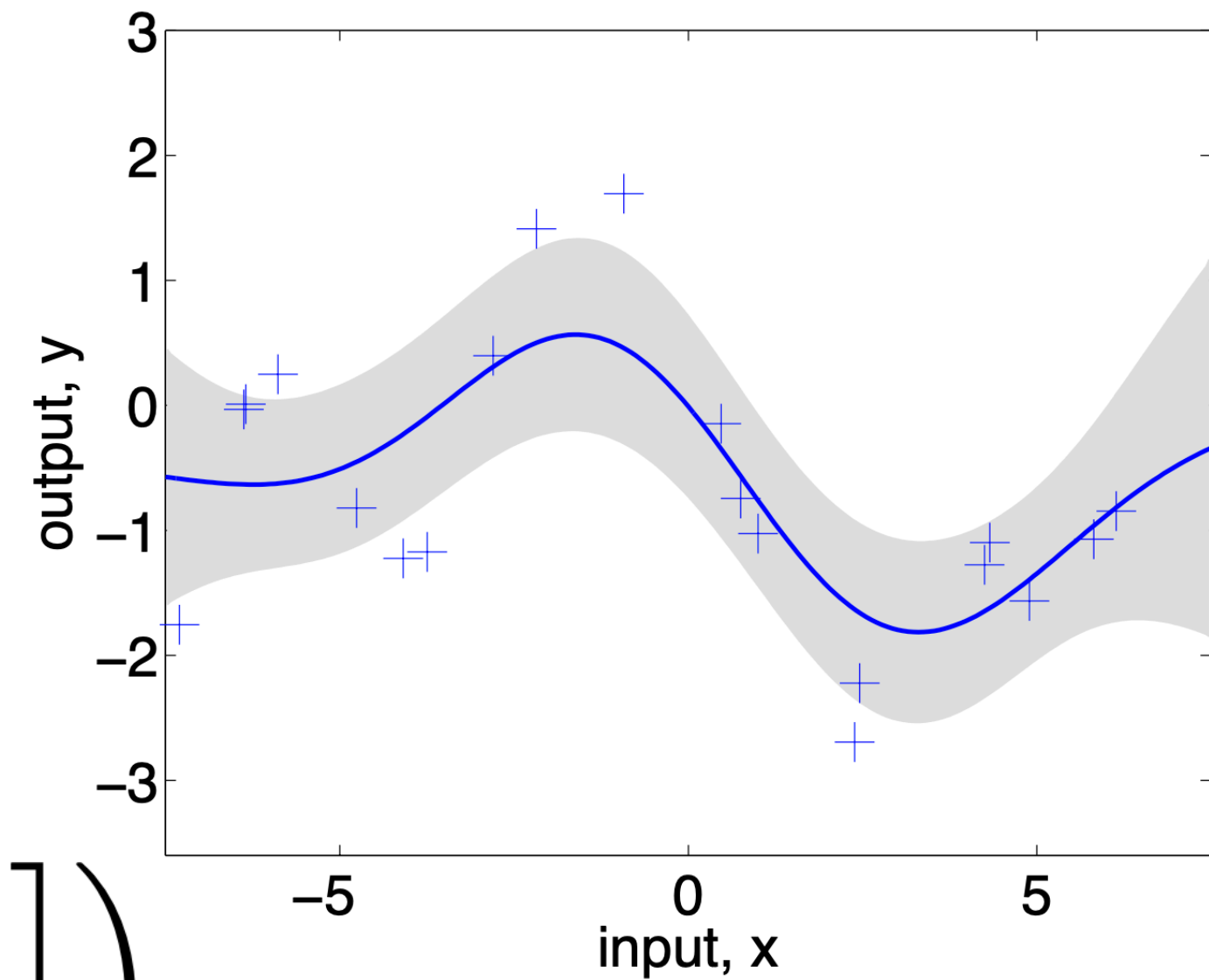
WEIGHT-SPACE FUNCTION-SPACE DUALITY

Gaussian Process prior:

$$\mathbf{f}_* \sim \mathcal{N}(\mathbf{0}, K(X_*, X_*))$$

Gaussian Process joint:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(X, X) + \sigma_n^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right)$$



Gaussian Process posterior:

$$\mathbf{f}_* | X, \mathbf{y}, X_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)), \text{ where}$$

$$\bar{\mathbf{f}}_* \triangleq \mathbb{E}[\mathbf{f}_* | X, \mathbf{y}, X_*] = K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} \mathbf{y},$$

$$\text{cov}(\mathbf{f}_*) = K(X_*, X_*) - K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} K(X, X_*)$$

WEIGHT-SPACE FUNCTION-SPACE DUALITY

Gaussian Process posterior predictive distribution:

$\mathbf{f}_* | X, \mathbf{y}, X_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)),$ where

$$\bar{\mathbf{f}}_* \triangleq \mathbb{E}[\mathbf{f}_* | X, \mathbf{y}, X_*] = K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} \mathbf{y},$$

$$\text{cov}(\mathbf{f}_*) = K(X_*, X_*) - K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} K(X, X_*)$$

Bayesian linear regression posterior predictive distribution:

$f_* | \mathbf{x}_*, X, \mathbf{y} \sim \mathcal{N}(\phi_*^\top \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \mathbf{y},$

$$\phi_*^\top \Sigma_p \phi_* - \phi_*^\top \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \Phi^\top \Sigma_p \phi_*)$$

$$\phi(\mathbf{x}_*) = \phi_* \quad K = \Phi^\top \Sigma_p \Phi$$

IMPLEMENTING GAUSSIAN PROCESS REGRESSION

IMPLEMENTING GAUSSIAN PROCESS REGRESSION

Gaussian Process prior:

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

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

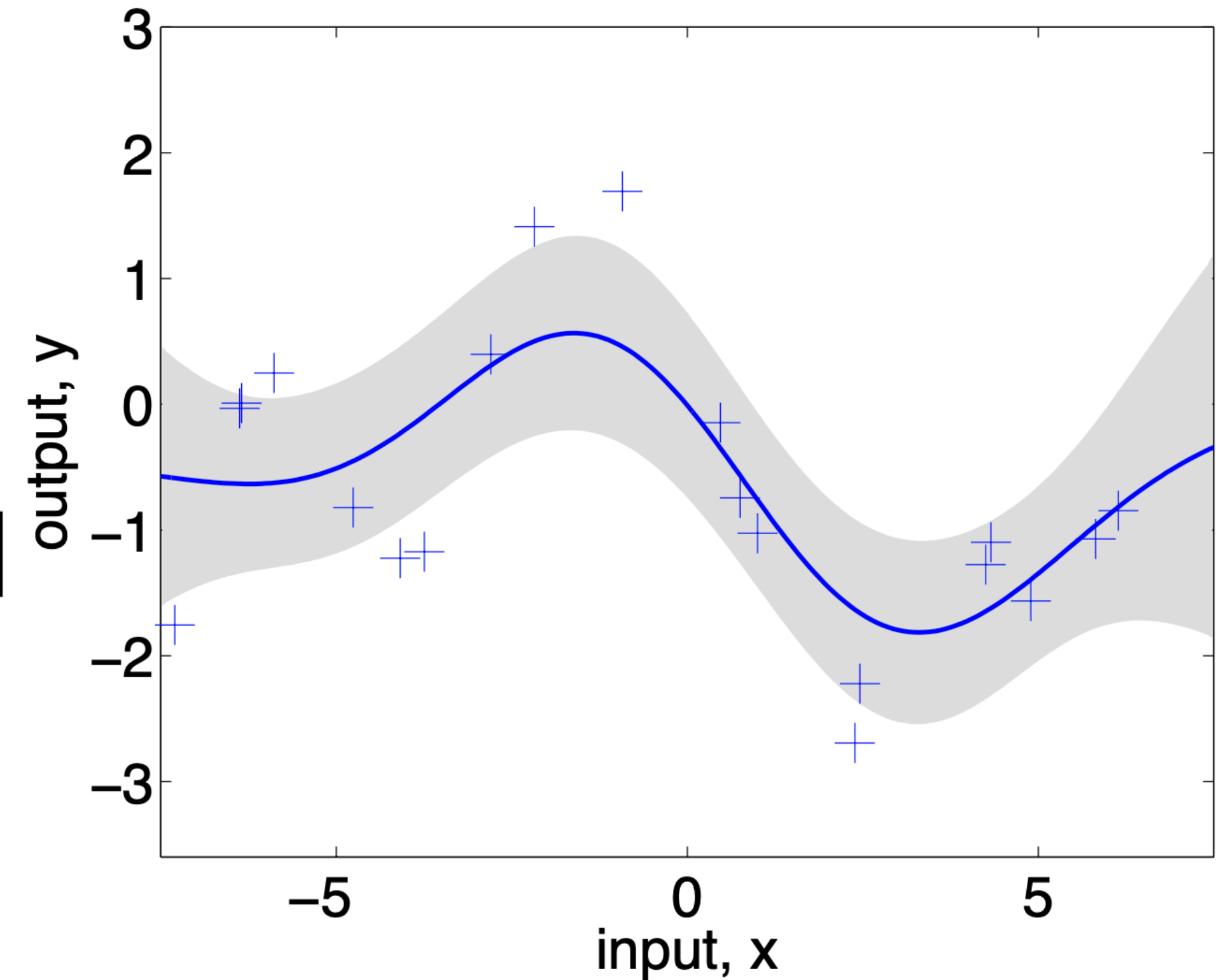
$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$

Gaussian Process posterior:

$$\mathbf{f}_* | X, \mathbf{y}, X_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)), \text{ where}$$

$$\bar{\mathbf{f}}_* \triangleq \mathbb{E}[\mathbf{f}_* | X, \mathbf{y}, X_*] = K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} \mathbf{y},$$

$$\text{cov}(\mathbf{f}_*) = K(X_*, X_*) - K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} K(X, X_*)$$



IMPLEMENTING GAUSSIAN PROCESS REGRESSION

Marginal likelihood:

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

$$\log p(\mathbf{f}|X) = -\frac{1}{2}\mathbf{f}^\top K^{-1}\mathbf{f} - \frac{1}{2}\log |K| - \frac{n}{2}\log 2\pi.$$

$$\log p(\mathbf{y}|X) = -\frac{1}{2}\mathbf{y}^\top (K + \sigma_n^2 I)^{-1}\mathbf{y} - \frac{1}{2}\log |K + \sigma_n^2 I| - \frac{n}{2}\log 2\pi$$

input: X (inputs), \mathbf{y} (targets), k (covariance function), σ_n^2 (noise level),
 \mathbf{x}_* (test input)

2: $L := \text{cholesky}(K + \sigma_n^2 I)$
 $\boldsymbol{\alpha} := L^\top \setminus (L \setminus \mathbf{y})$ } predictive mean eq. (2.25)

4: $\bar{f}_* := \mathbf{k}_*^\top \boldsymbol{\alpha}$
 $\mathbf{v} := L \setminus \mathbf{k}_*$ } predictive variance eq. (2.26)

6: $\mathbb{V}[f_*] := k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{v}^\top \mathbf{v}$
 $\log p(\mathbf{y}|X) := -\frac{1}{2}\mathbf{y}^\top \boldsymbol{\alpha} - \sum_i \log L_{ii} - \frac{n}{2}\log 2\pi$ eq. (2.30)

8: **return:** \bar{f}_* (mean), $\mathbb{V}[f_*]$ (variance), $\log p(\mathbf{y}|X)$ (log marginal likelihood)

IMPLEMENTING GAUSSIAN PROCESS REGRESSION

Implementation checklist:

1) GP prior, including kernel function

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

2) Posterior predictive distribution:

$$\mathbf{f}_* | X, \mathbf{y}, X_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)), \text{ where}$$

$$\bar{\mathbf{f}}_* \triangleq \mathbb{E}[\mathbf{f}_* | X, \mathbf{y}, X_*] = K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} \mathbf{y},$$

$$\text{cov}(\mathbf{f}_*) = K(X_*, X_*) - K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} K(X, X_*)$$

3) Marginal likelihood:

$$\log p(\mathbf{y} | X) = -\frac{1}{2} \mathbf{y}^\top (K + \sigma_n^2 I)^{-1} \mathbf{y} - \frac{1}{2} \log |K + \sigma_n^2 I| - \frac{n}{2} \log 2\pi$$

LAB ASSIGNMENT

Course page: `http://www.robots.ox.ac.uk/~mosb/aims_cdt`

Goal: predict missing sensor measurements using GP regression