

ENM53 I: Data-driven modeling and probabilistic scientific computing

Lecture #19: Gaussian processes

Paris Perdikaris
April 14, 2020



Using Gaussian processes for nonlinear regression

Imagine observing a data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)_{i=1}^n\} = (\mathbf{X}, \mathbf{y})$.

Model:

$$y_i = f(\mathbf{x}_i) + \epsilon_i$$

$$f \sim \text{GP}(\cdot|0, K)$$

$$\epsilon_i \sim \text{N}(\cdot|0, \sigma^2)$$

Prior on f is a GP, likelihood is Gaussian, therefore posterior on f is also a GP.

We can use this to make predictions

$$p(y_*|\mathbf{x}_*, \mathcal{D}) = \int p(y_*|\mathbf{x}_*, f, \mathcal{D}) p(f|\mathcal{D}) df$$

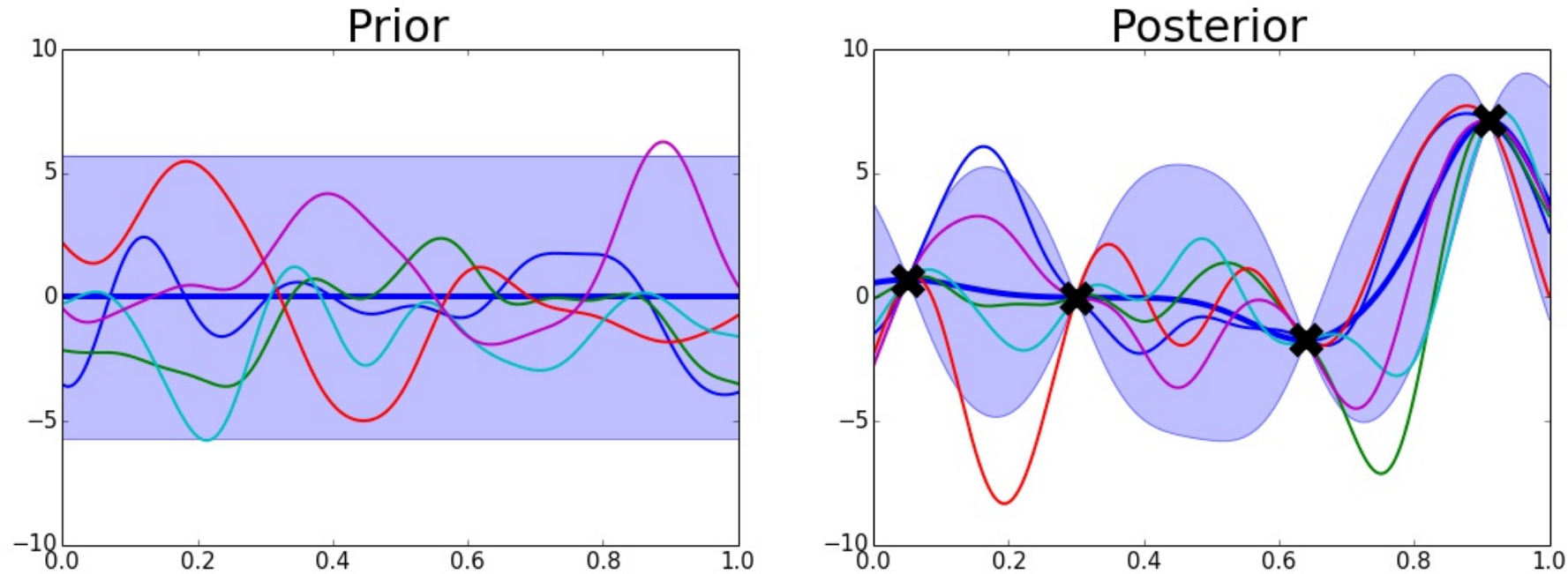
We can also compute the marginal likelihood (evidence) and use this to compare or tune covariance functions

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

Data-driven modeling with Gaussian processes

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

$$f \sim \mathcal{GP}(0, k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}))$$



Training via maximizing the marginal likelihood

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2} \log |\mathbf{K} + \sigma_{\epsilon}^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^T (\mathbf{K} + \sigma_{\epsilon}^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N}{2} \log 2\pi$$

Prediction via conditioning on available data

$$p(f_* | \mathbf{y}, \mathbf{X}, \mathbf{x}_*) = \mathcal{N}(f_* | \mu_*, \sigma_*^2),$$

$$\mu_*(\mathbf{x}_*) = \mathbf{k}_{*N} (\mathbf{K} + \sigma_{\epsilon}^2 \mathbf{I})^{-1} \mathbf{y},$$

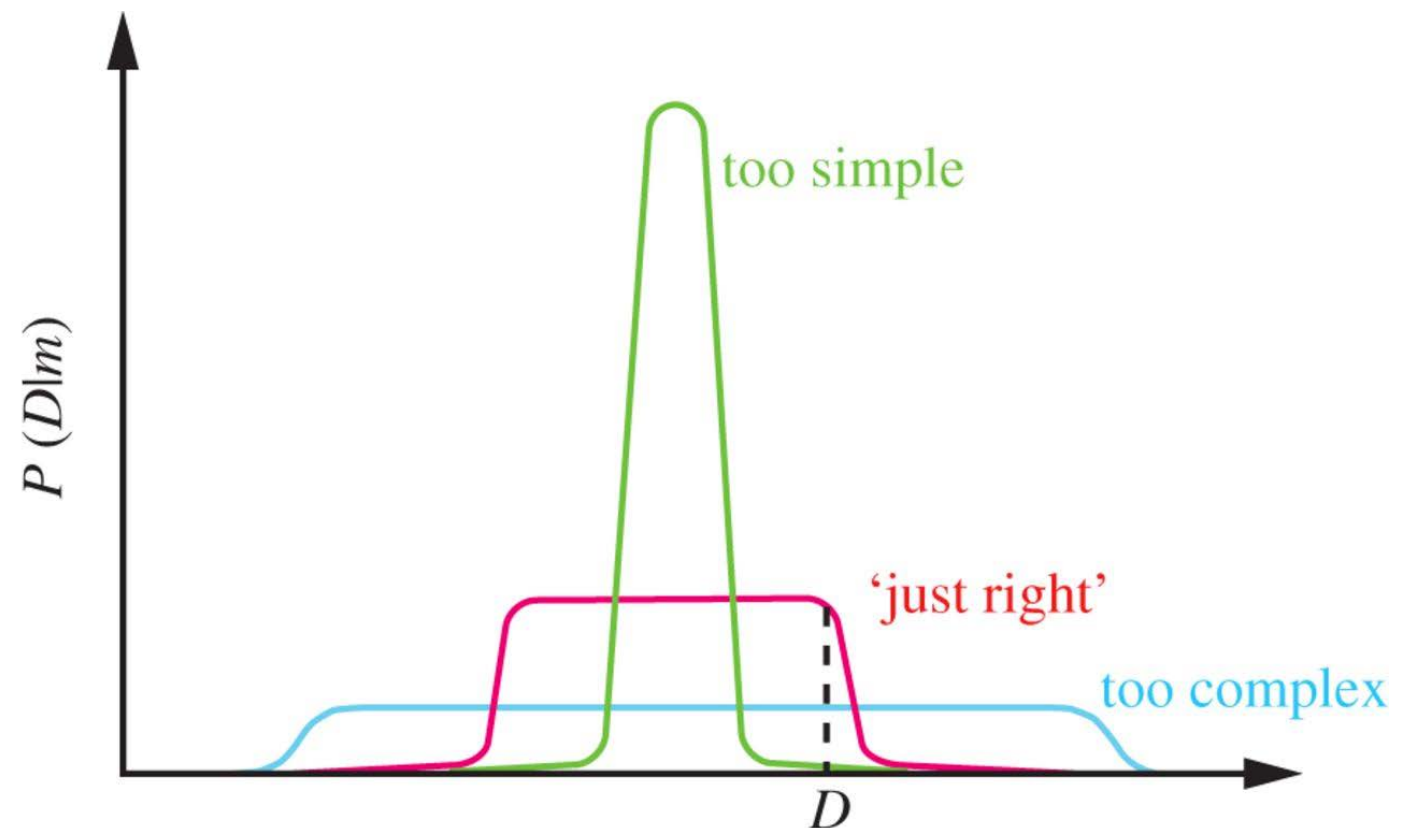
$$\sigma_*^2(\mathbf{x}_*) = \mathbf{k}_{**} - \mathbf{k}_{*N} (\mathbf{K} + \sigma_{\epsilon}^2 \mathbf{I})^{-1} \mathbf{k}_{N*},$$

Occam's razor

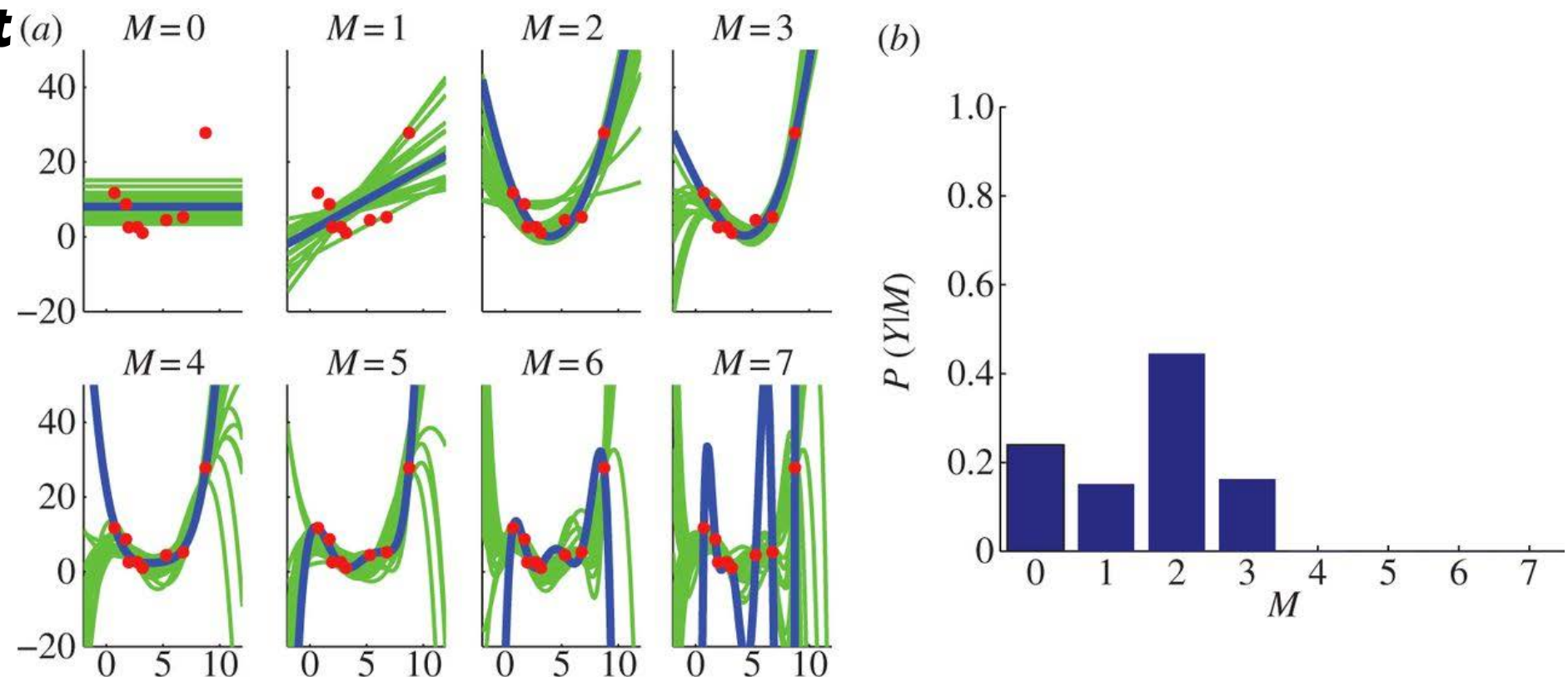
William of Ockham (~1285-1347 A.D)



“plurality should not be posited without necessity.”



all possible datasets of size n

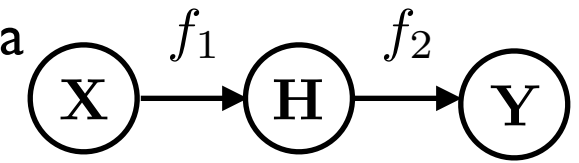


Ghahramani, Z. (2013). Bayesian non-parametrics and the probabilistic approach to modelling. *Phil. Trans. R. Soc. A*, 371(1984), 20110553.

Challenges, limitations, and recent progress

Discontinuities and non-stationarity: GPs struggle with discontinuous data

Use warping functions to transform into a jointly stationary input space



- Log, sigmoid, betaCDF —> “Warped GPs” *Snelson, E., C.E. Rasmussen, and Z.Ghahramani. "Warped gaussian processes."*
- Neural networks —> “Manifold GPs” *Calandra, R., et al. "Manifold Gaussian processes for regression."*
- Gaussian processes —> “Deep GPs” *Damianou, A. C., and N.D. Lawrence. "Deep Gaussian processes."*

Theoretical guarantees: Accuracy, convergence rates, posterior consistency, contraction rates, etc.

Approximation theory in Reproducing Kernel Hilbert Spaces

Stuart, A.M., and A.L.Teckentrup. "Posterior consistency for Gaussian process approximations of Bayesian posterior distributions."

Scalability: GPs suffer from a cubic scaling with the data

Low-rank approximations to the covariance

Snelson, E., and Z. Ghahramani. "Sparse Gaussian processes using pseudo-inputs."

Frequency-domain learning algorithms

Perdikaris P., D.Venturi, G.E. Karniadakis “Multi-fidelity information fusion algorithms for high dimensional systems and massive data-sets”.

Stochastic variational inference

Cheng, C., and B. Boots. "Variational Inference for Gaussian Process Models with Linear Complexity." NIPS, 2017.

Smart linear algebra and GPU acceleration

Wang, K. e. al.. “Exact Gaussian Processes on a Million Data Points”, 2019.

High-dimensions: Tensor product kernels suffer from the curse of dimensionality, i.e. they require an exponentially increasing amount of training data

Data-driven additive kernels

Perdikaris P., D.Venturi, G.E. Karniadakis “Multi-fidelity information fusion algorithms for high dimensional systems and massive data-sets”.

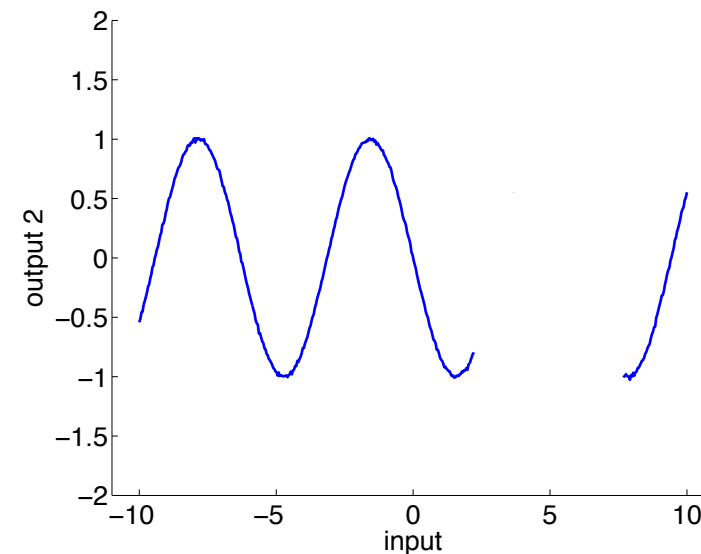
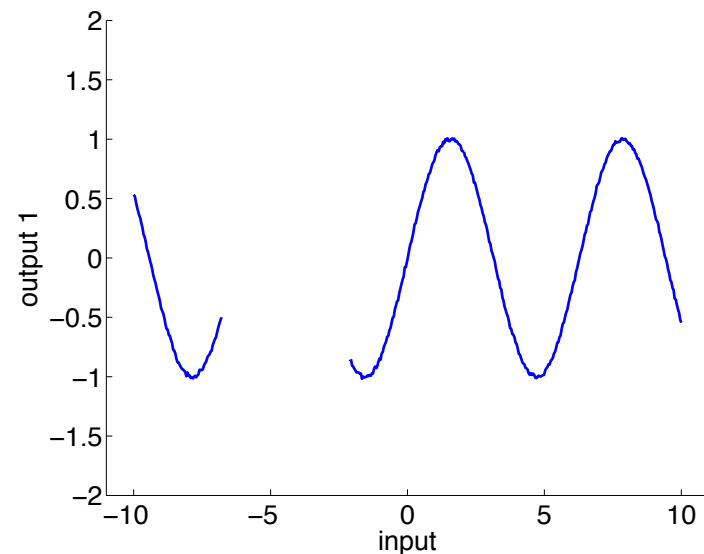
Unsupervised dimensionality-reduction (GPLVM, variational auto-encoders)

Lawrence, N.D. "Gaussian process latent variable models for visualisation of high dimensional data."

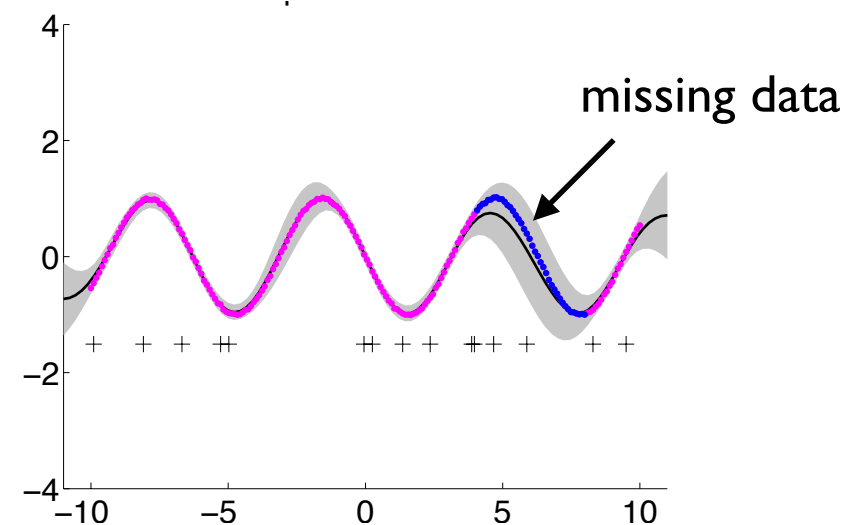
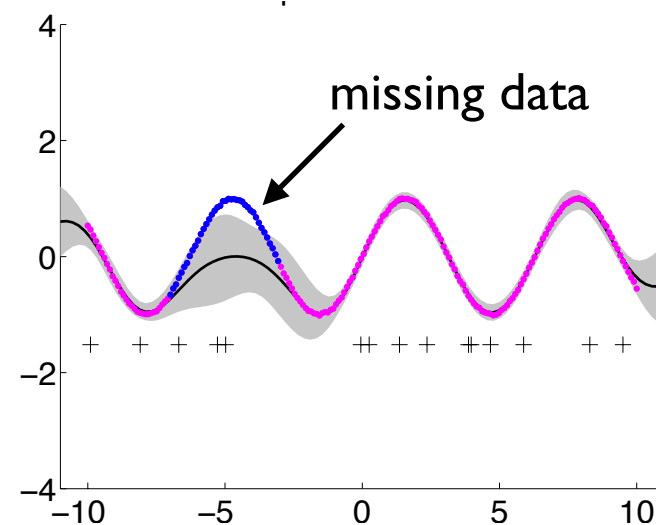
Multi-output Gaussian process regression

Learn two correlated tasks (outputs) with lots of data + missing data

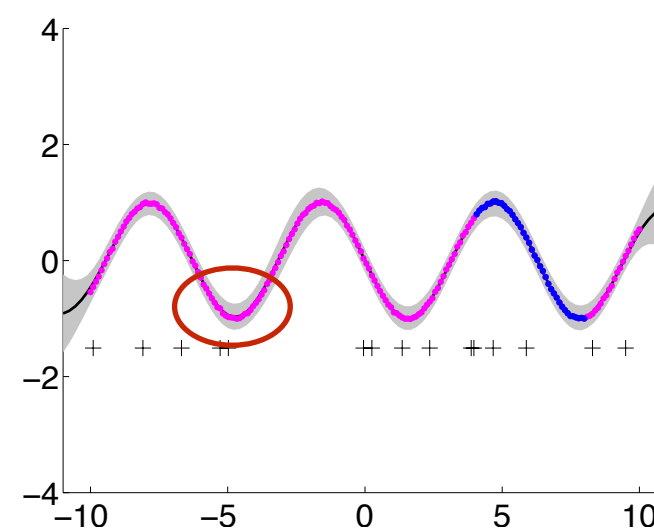
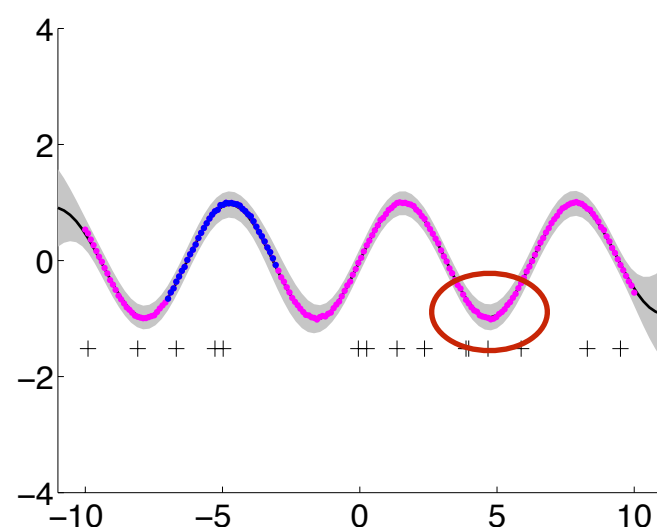
$$\text{output 2} = -\text{output 1}$$



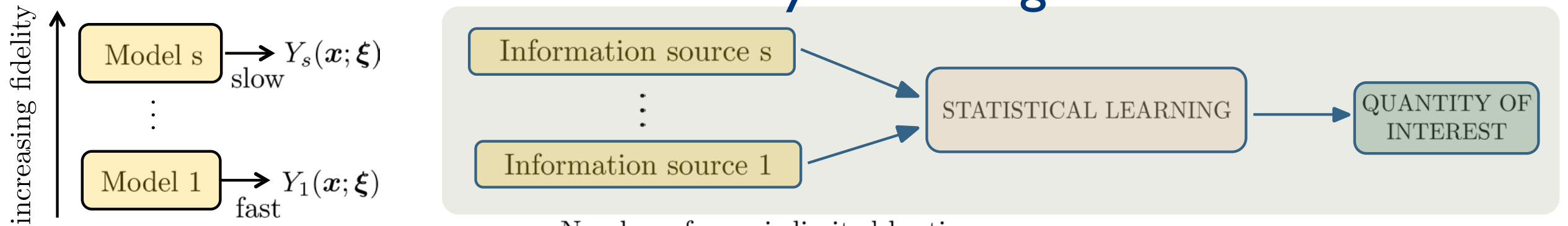
independent
learning



COGP



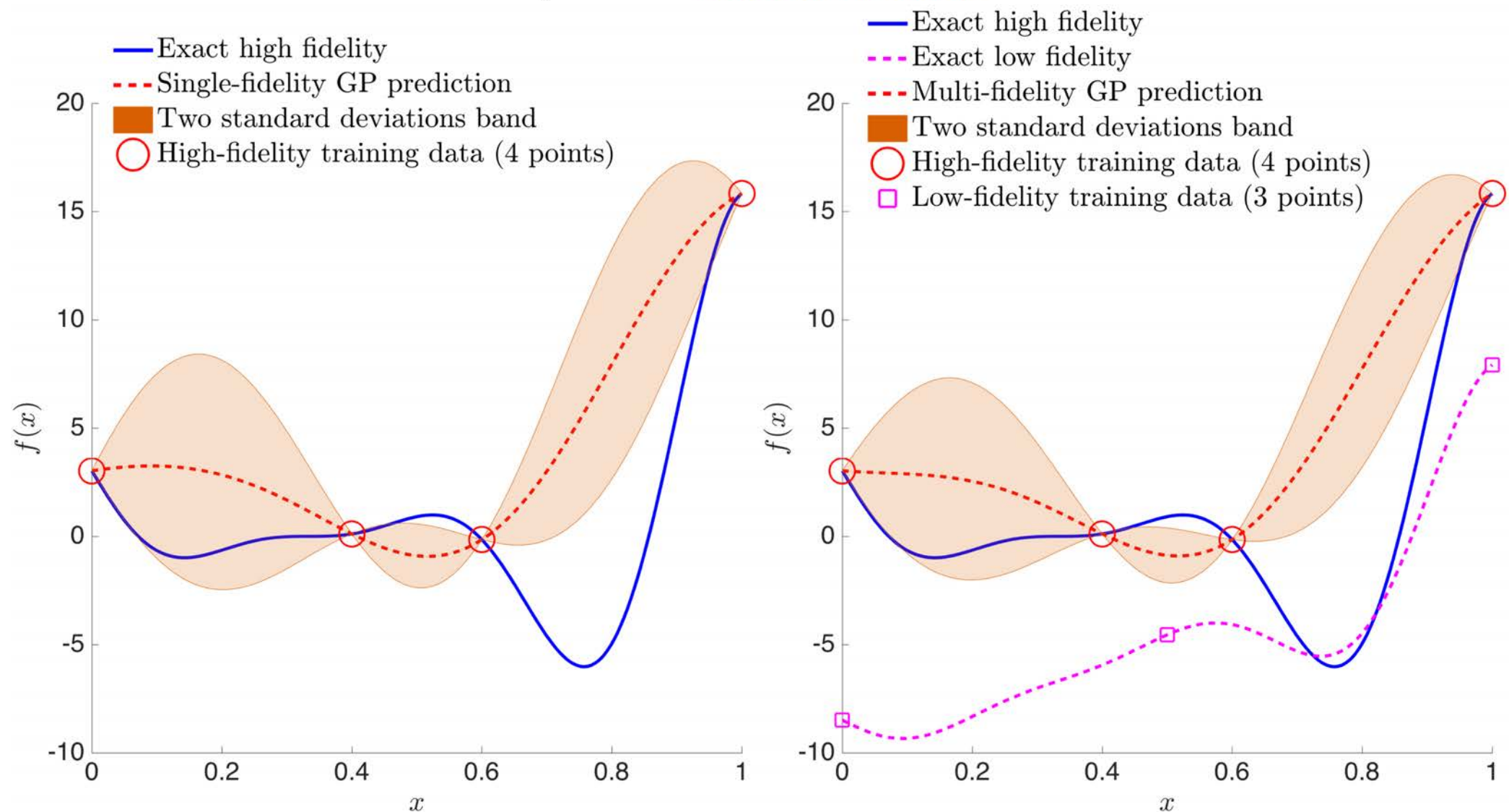
Multi-fidelity modeling



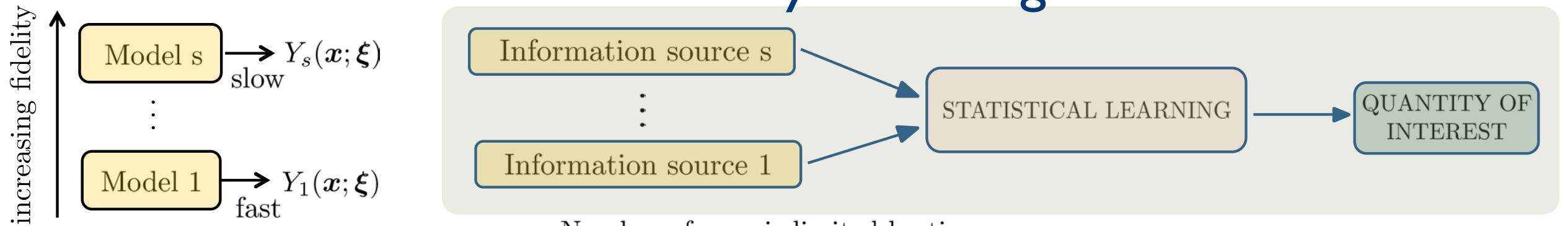
Number of runs is limited by time and computational resources

We cannot compute at all $(x; \xi)$

Prediction of $Z_i(x) = \mathbb{E}[f(Y_i(x; \xi))]$ is a problem of **statistical inference**



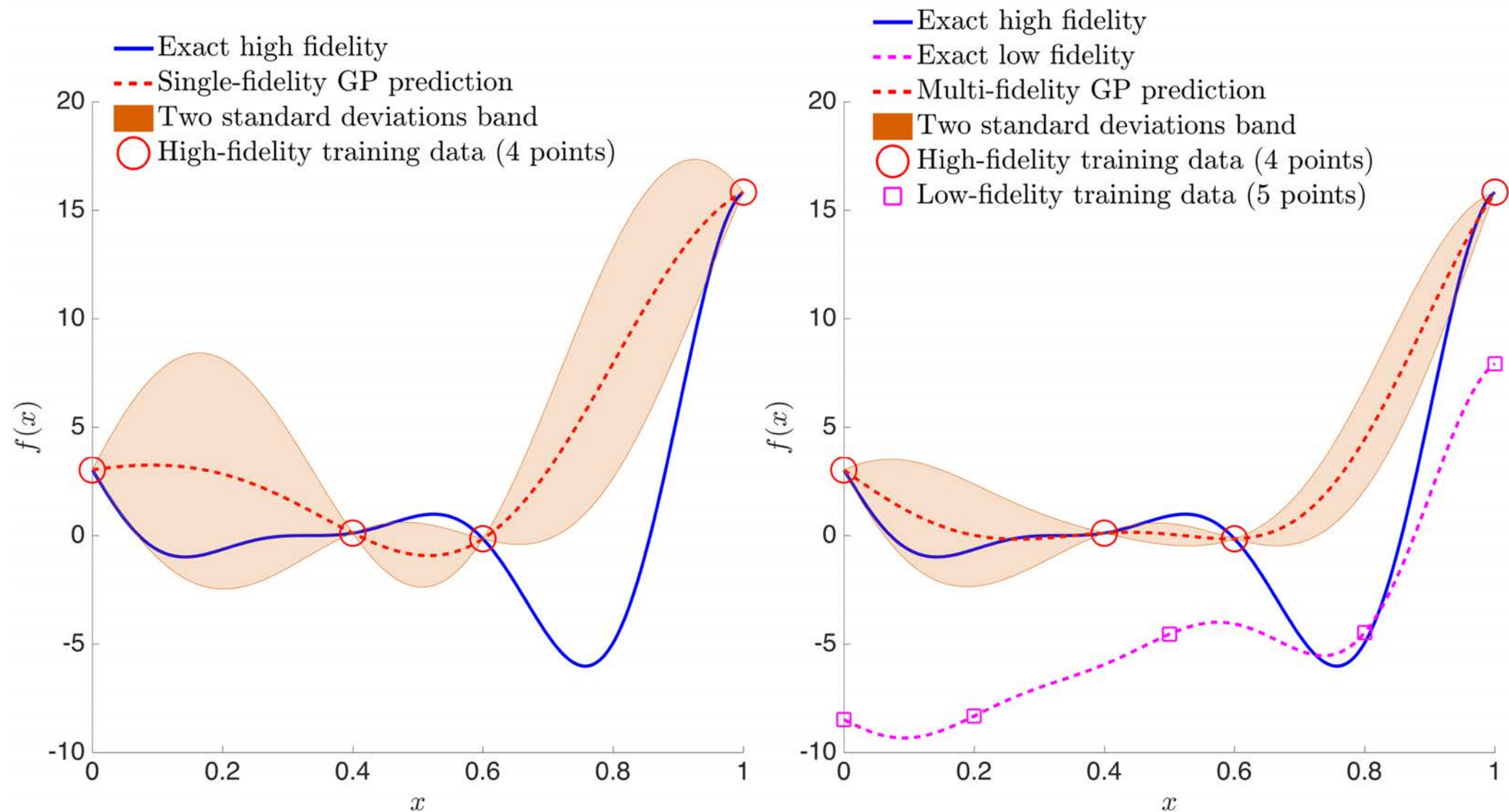
Multi-fidelity modeling



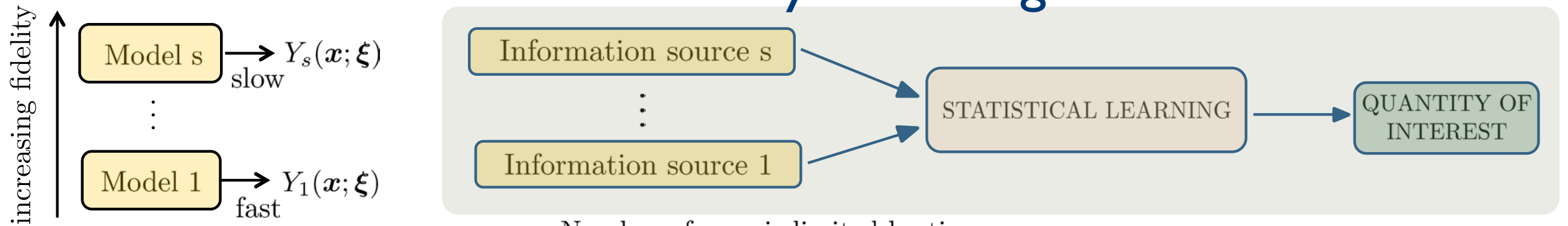
Number of runs is limited by time and computational resources

We cannot compute at all $(x; \xi)$

Prediction of $Z_i(x) = \mathbb{E}[f(Y_i(x; \xi))]$ is a problem of **statistical inference**



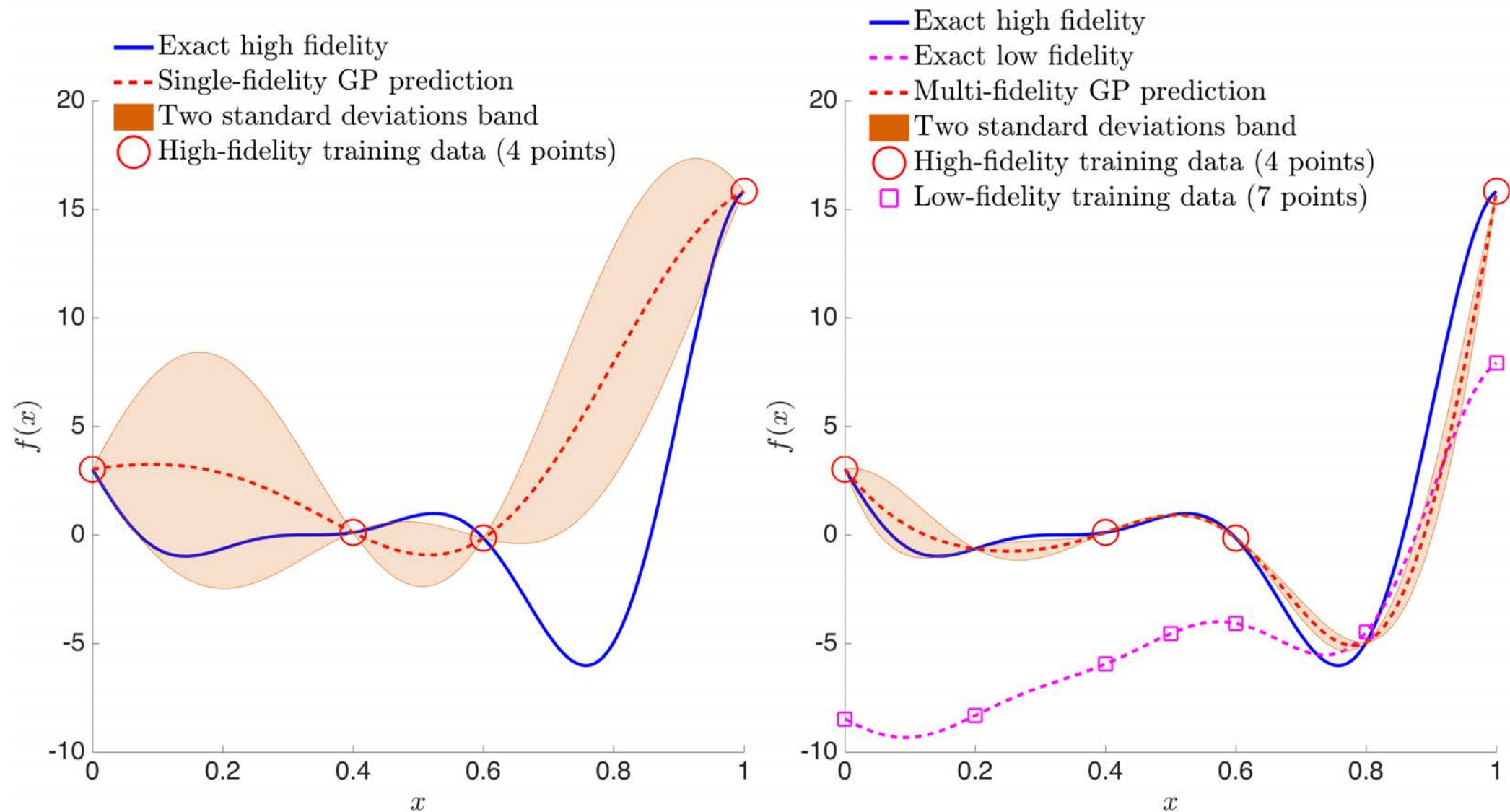
Multi-fidelity modeling



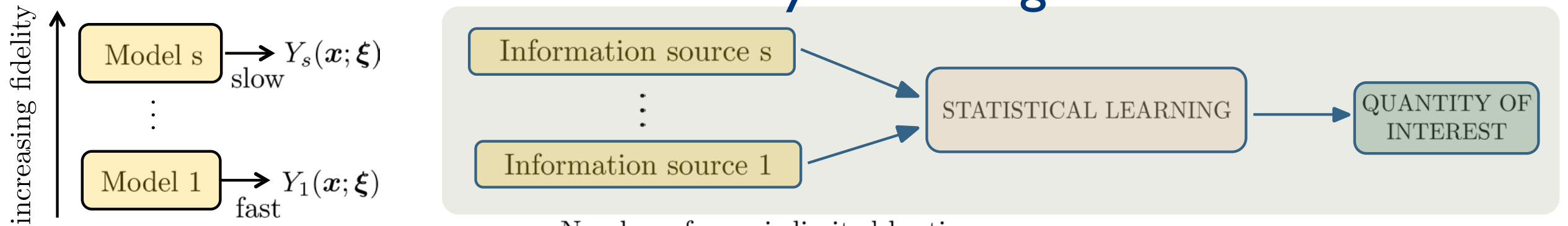
Number of runs is limited by time and computational resources

We cannot compute at all $(\mathbf{x}; \xi)$

Prediction of $Z_i(\mathbf{x}) = \mathbb{E}[f(Y_i(\mathbf{x}; \xi))]$ is a problem of **statistical inference**



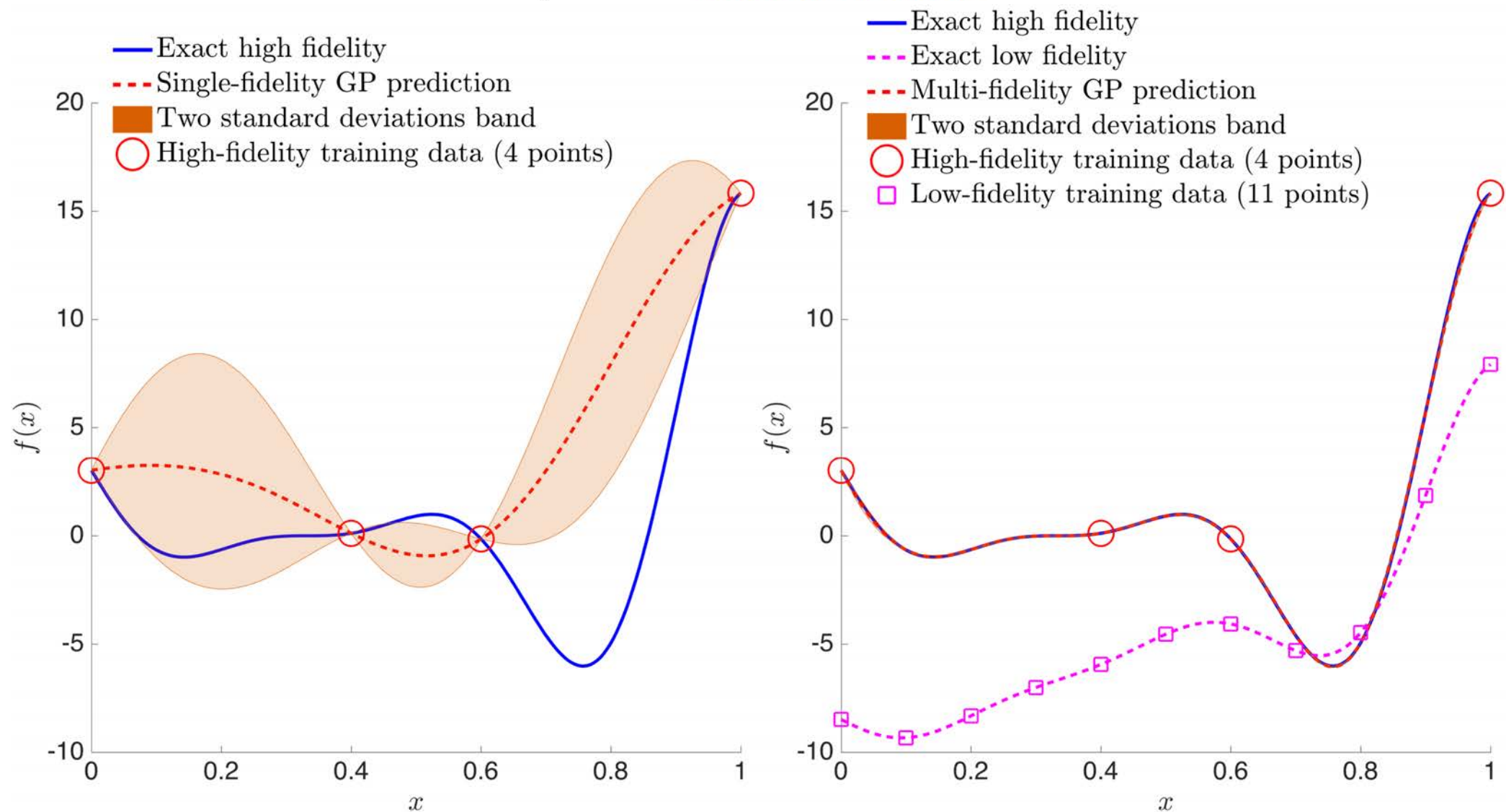
Multi-fidelity modeling



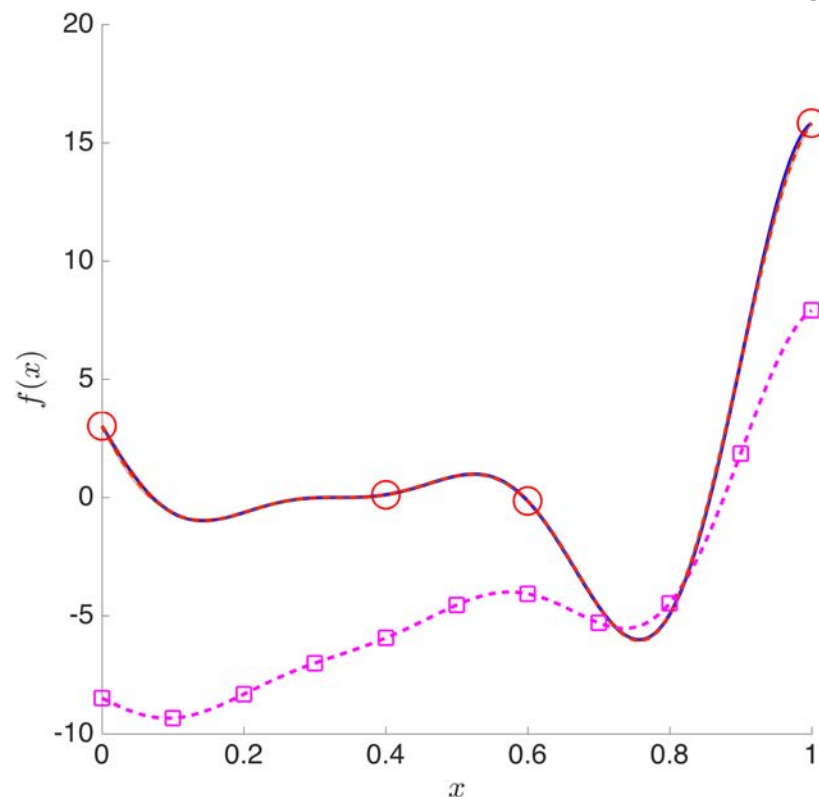
Number of runs is limited by time and computational resources

We cannot compute at all $(\mathbf{x}; \xi)$

Prediction of $Z_i(\mathbf{x}) = \mathbb{E}[f(Y_i(\mathbf{x}; \xi))]$ is a problem of **statistical inference**



Multi-fidelity modeling



Multi-fidelity observations:

$$\mathbf{y}_L = f_L(\mathbf{x}_L) + \epsilon_L$$

$$\mathbf{y}_H = f_H(\mathbf{x}_H) + \epsilon_H$$

Probabilistic model:

$$f_H(\mathbf{x}) = \rho f_L(\mathbf{x}) + \delta(\mathbf{x})$$

$$f_L(\mathbf{x}) \sim \mathcal{GP}(0, k_L(\mathbf{x}, \mathbf{x}'; \theta_L))$$

$$\delta(\mathbf{x}) \sim \mathcal{GP}(0, k_H(\mathbf{x}, \mathbf{x}'; \theta_H))$$

$$\epsilon_L \sim \mathcal{N}(0, \sigma_{\epsilon_L}^2 \mathbf{I})$$

$$\epsilon_H \sim \mathcal{N}(0, \sigma_{\epsilon_H}^2 \mathbf{I})$$

Training:

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_L \\ \mathbf{y}_H \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} k_L(\mathbf{x}_L, \mathbf{x}_L'; \theta_L) + \sigma_{\epsilon_L}^2 \mathbf{I} & \rho k_L(\mathbf{x}_L, \mathbf{x}_H'; \theta_L) \\ \rho k_L(\mathbf{x}_H, \mathbf{x}_L'; \theta_L) & \rho^2 k_L(\mathbf{x}_H, \mathbf{x}_H'; \theta_L) + k_H(\mathbf{x}_H, \mathbf{x}_H'; \theta_H) + \sigma_{\epsilon_H}^2 \mathbf{I} \end{bmatrix} \right)$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_L \\ \mathbf{x}_H \end{bmatrix} \quad -\log p(\mathbf{y}|\mathbf{X}, \theta_L, \theta_H, \rho, \sigma_{\epsilon_L}^2, \sigma_{\epsilon_H}^2) = \frac{1}{2} \log |\mathbf{K}| + \frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} - \frac{N_L + N_H}{2} \log 2\pi$$

Prediction:

$$p(f(\mathbf{x}^*)|\mathbf{y}, \mathbf{X}, \mathbf{x}^*) \sim \mathcal{N}(f(\mathbf{x}^*)|\mu(\mathbf{x}^*), \sigma^2(\mathbf{x}^*))$$

$$\mu(\mathbf{x}^*) = \mathbf{k}(\mathbf{x}^*, \mathbf{X}) \mathbf{K}^{-1} \mathbf{y}$$

$$\sigma(\mathbf{x}^*) = \mathbf{k}(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{k}(\mathbf{x}^*, \mathbf{X}) \mathbf{K}^{-1} \mathbf{k}(\mathbf{X}, \mathbf{x}^*)$$

M.C Kennedy, and A. O'Hagan. *Predicting the output from a complex computer code when fast approximations are available*, 2000.

Demo code: <https://github.com/PredictiveIntelligenceLab/GPTutorial>