

Applications of Gaussian Processes

EN.540.782: Statistical Uncertainty Quantification

N. Wichrowski

Johns Hopkins University

April 8, 2021

Outline

- 1 Review: Gaussian Processes
- 2 Bayesian Optimization
 - The Algorithm
 - Acquisition Functions
 - Practical Considerations
- 3 Mixed-Input Regression
- 4 Conclusion

Review: Gaussian Processes

Definition

A Gaussian Process (GP) is a collection of random variables $\{X_\alpha \mid \alpha \in \mathcal{X}\}$ any finite number of which have a joint Gaussian distribution.

- A GP is completely specified by:
 - ▶ mean function: $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$
 - ▶ covariance kernel: $k(\mathbf{x}, \mathbf{x}') = \text{Cov}[f(\mathbf{x}), f(\mathbf{x}')]]$
- Use conditional distributions to predict values at new inputs:
 - ▶ Observations: inputs $X \in \mathbb{R}^{n \times d}$, outputs $\mathbf{y} \in \mathbb{R}^n$.
 - ▶ For new inputs $X_\star \in \mathbb{R}^{n_\star \times d}$,

$$\mathbb{E}[\mathbf{y}_\star | X_\star, X, \mathbf{y}] = K(X_\star, X)K(X, X)^{-1}\mathbf{y}$$

$$\text{Cov}[\mathbf{y}_\star | X_\star, X, \mathbf{y}] = K(X_\star, X_\star) - K(X_\star, X)K(X, X)^{-1}K(X, X_\star)$$

Problem Setting

- Objective $f : \mathcal{X} \rightarrow \mathbb{R}$ on domain $\mathcal{X} \subseteq \mathbb{R}^d$.
- We seek a global minimizer

$$\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}).$$

- Bayesian optimization (BO) is most appropriate when:
 - ▶ Evaluating f is expensive (and possibly noisy).
 - ▶ We have no derivative information.

Basic BO Algorithm

Algorithm 1: Bayesian Optimization

Data: objective function $f : \mathcal{X} \rightarrow \mathbb{R}$; initial sample $\mathcal{D}_n = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$

for $i \in \{n+1, \dots, N\}$ **do**

 Fit a GPR model to data \mathcal{D}_{i-1} ;

 Choose a new point $\mathbf{x}_i \in \mathcal{X}$;

 Evaluate $y_i \leftarrow f(\mathbf{x}_i)$;

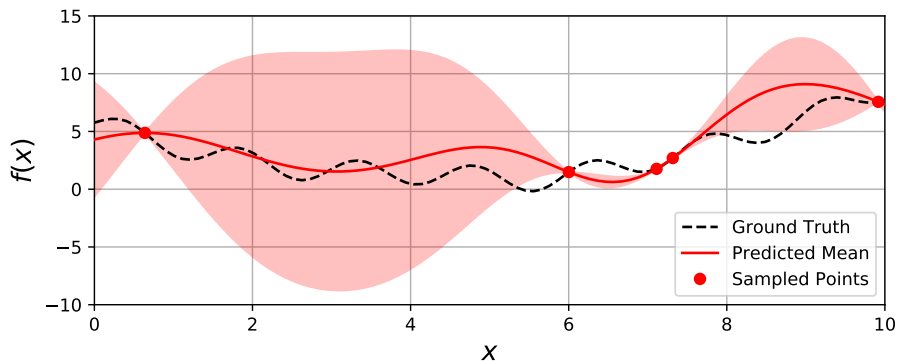
 Update sample $\mathcal{D}_i \leftarrow \mathcal{D}_{i-1} \cup \{(\mathbf{x}_i, y_i)\}$;

end

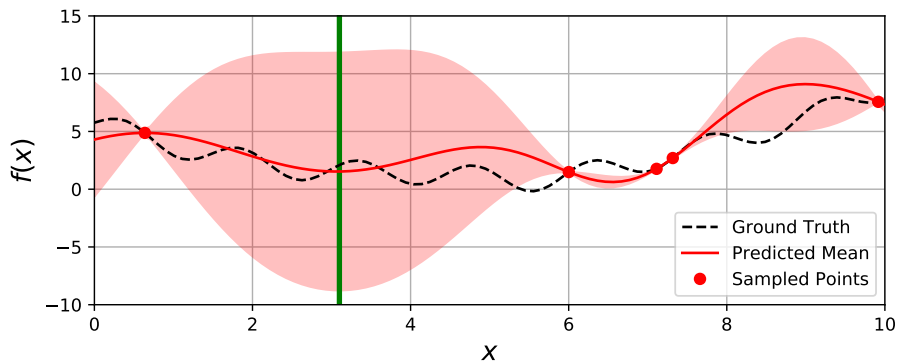
$j \leftarrow \underset{1 \leq i \leq N}{\operatorname{argmin}} y_i$;

Result: approximate global minimizer \mathbf{x}_j

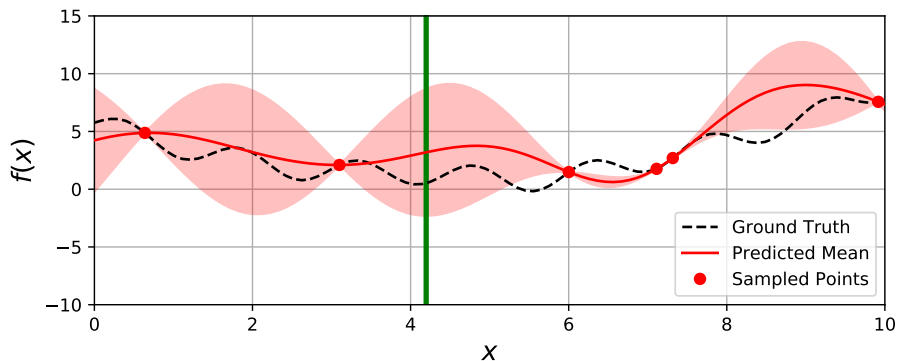
An Illustrated Example



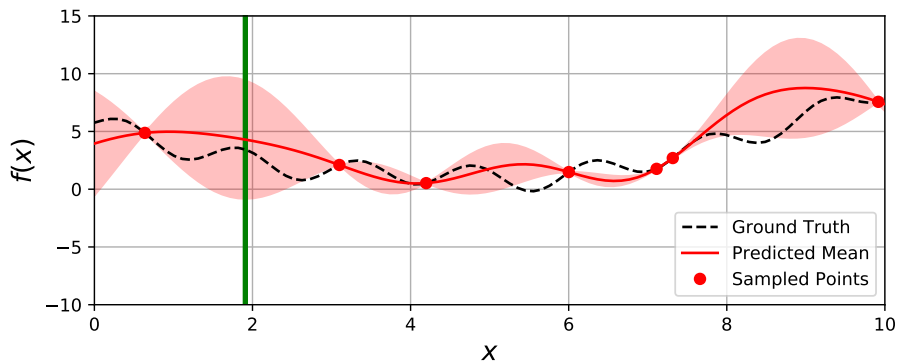
An Illustrated Example



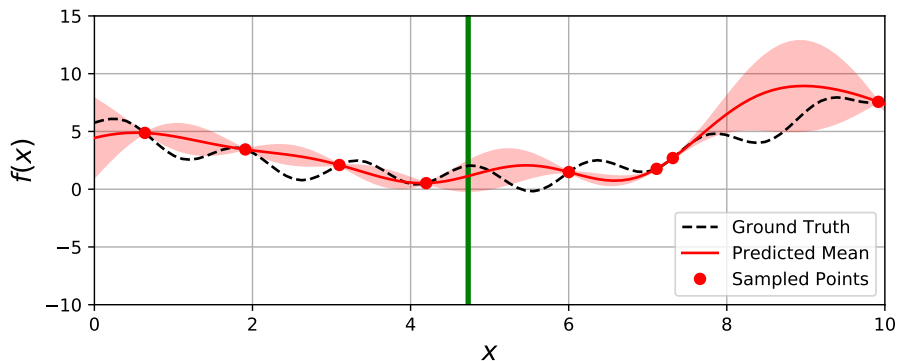
An Illustrated Example



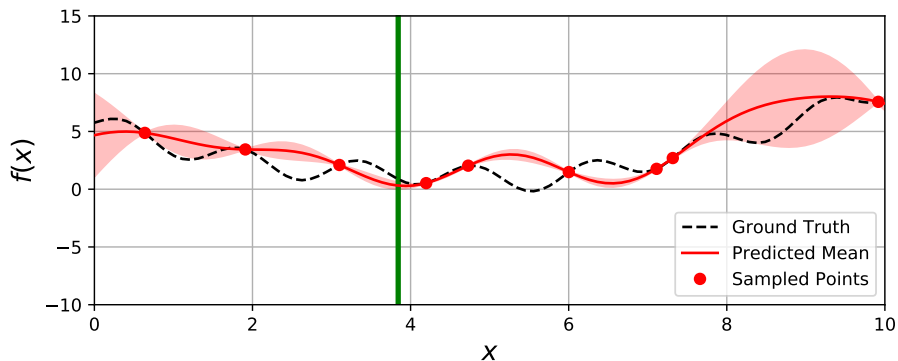
An Illustrated Example



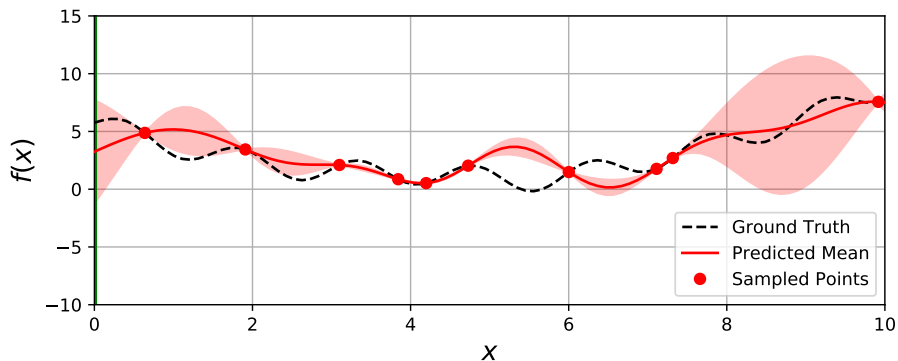
An Illustrated Example



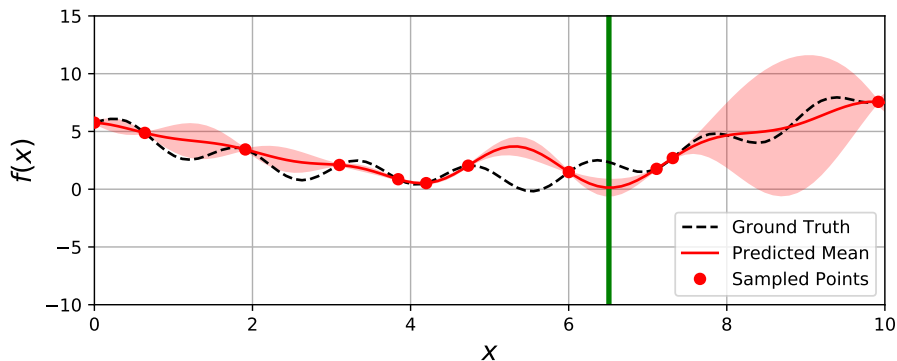
An Illustrated Example



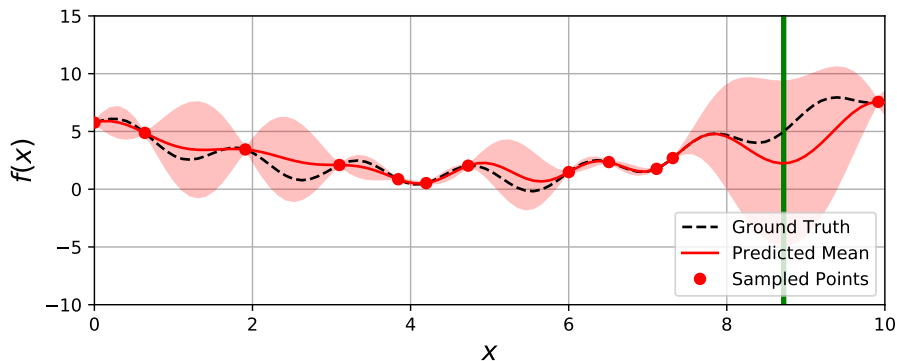
An Illustrated Example



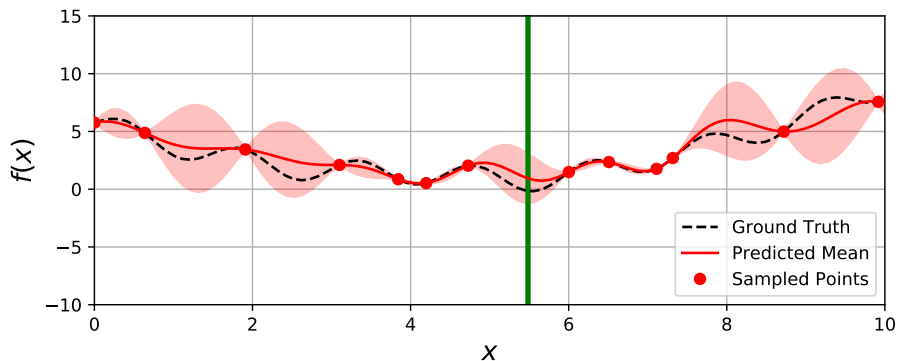
An Illustrated Example



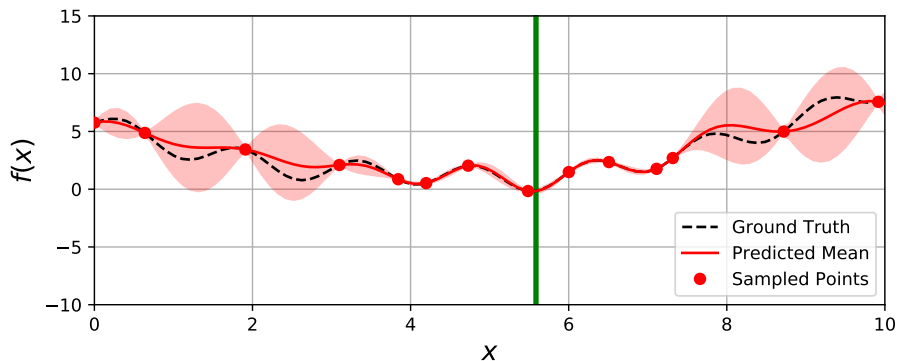
An Illustrated Example



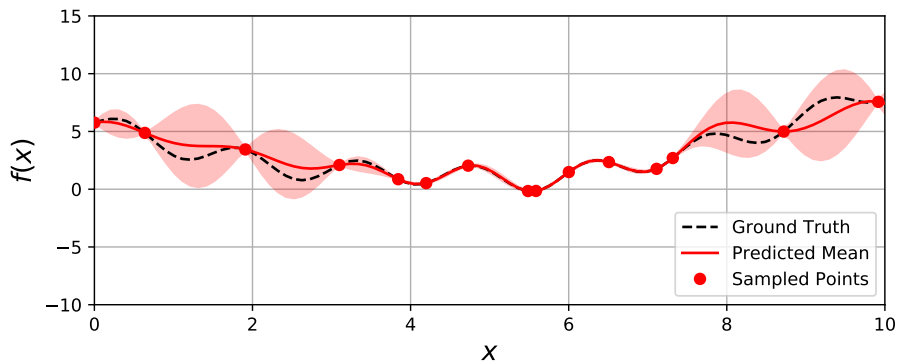
An Illustrated Example



An Illustrated Example



An Illustrated Example



Selecting the Next Sample Point

- Use current knowledge to inform our choice.
- Exploration vs. Exploitation
- An acquisition function measures utility of sampling at $\mathbf{x} \in \mathcal{X}$:

$$\mathbf{x}_{i+1} = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \alpha(\mathbf{x}; \mathcal{D}_i)$$

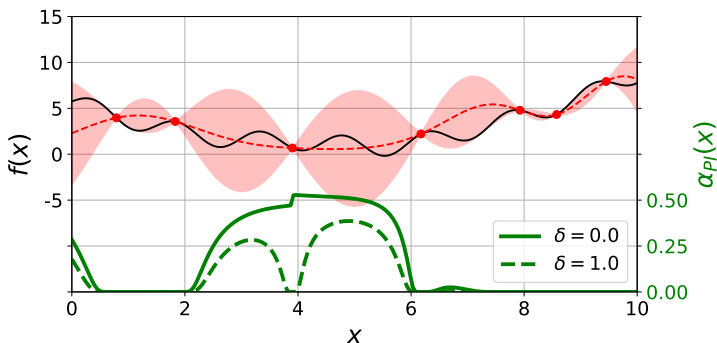
- How is this better? Optimizing α does not require evaluating f .

Acquisition Function: Probability of Improvement

- Which inputs are likely to improve the current best observation?
- Gaussian distribution provides an analytical expression:

$$\alpha_{PI}(\mathbf{x}; \mathcal{D}) = \mathbb{P}[f(\mathbf{x}) < \tau \mid \mathcal{D}] = \Phi\left(\frac{\tau - \mu(\mathbf{x})}{\sigma(\mathbf{x})}\right)$$

- Threshold τ usually written as $\tau = y_{\min} - \delta$ for $\delta \geq 0$.

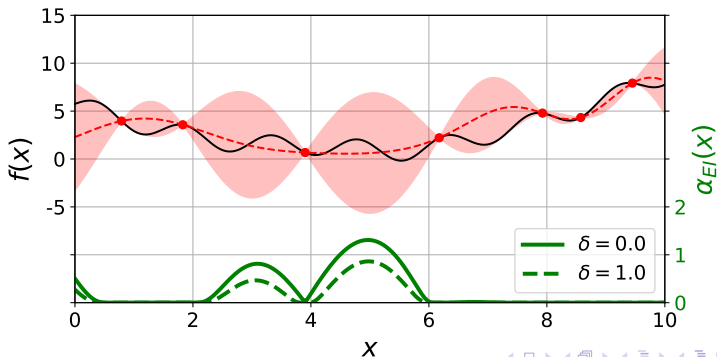


Acquisition Function: Expected Improvement

- Not all improvements are equally helpful.
- Expected improvement compared to threshold τ is

$$\begin{aligned}\alpha_{EI}(\mathbf{x}; \mathcal{D}) &= \mathbb{E} [\max(0, \tau - \mathcal{N} [\mu(\mathbf{x}), \sigma^2(\mathbf{x})])] \\ &= (\tau - \mu(\mathbf{x})) \Phi \left(\frac{\tau - \mu(\mathbf{x})}{\sigma(\mathbf{x})} \right) + \sigma(\mathbf{x}) \phi \left(\frac{\tau - \mu(\mathbf{x})}{\sigma(\mathbf{x})} \right).\end{aligned}$$

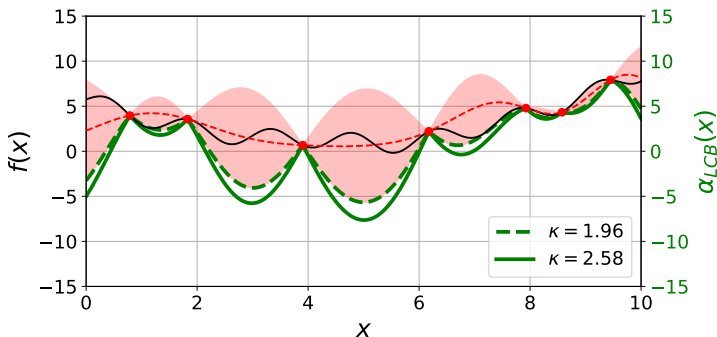
- Less prone than PI to become stuck in a local minimum.



Acquisition Function: Lower Confidence Bound

- Which inputs have superior best-case outputs?
- Usually expressed as a loss to be minimized:

$$\alpha_{LCB}(\mathbf{x}; \mathcal{D}) = \mu(\mathbf{x}) - \kappa\sigma(\mathbf{x}).$$



Computational Cost

- Inverting $K(X, X) \in \mathbb{R}^{n \times n}$ requires $O(n^3)$ operations.
- Cholesky decomposition must be updated each iteration.
- Approximation techniques exchange accuracy for speed.
- Using a sparse kernel for the GPR may help.

Pre-processing Data

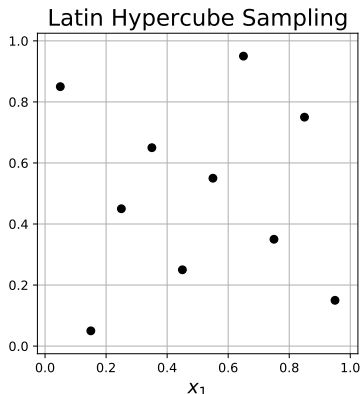
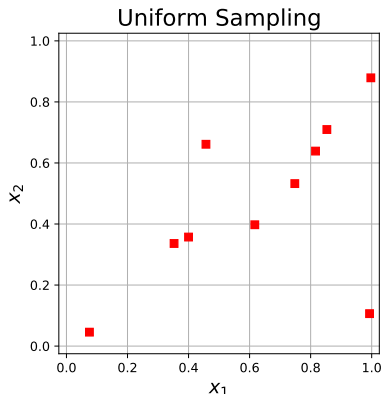
- Transforming sample data can facilitate GPR performance.
 - ▶ Normalize observations:

$$\tilde{\mathbf{x}}_i = \frac{\mathbf{x}_i - \bar{\mathbf{x}}}{s_x} \qquad \tilde{y}_i = \frac{y_i - \bar{y}}{s_y}$$

- ▶ If outputs are necessarily positive: $\tilde{y}_i = \log(y_i)$.
 - ▶ Dimensionality reduction: focus on important directions.
- Apply knowledge of problem when deciding how to process data.
- Must reverse transformations to recover interpretable quantities.

Generating a Useful Initial Sample

- Samples drawn uniformly may not fill the domain.
- Latin hypercube sampling offers better samples for regression.



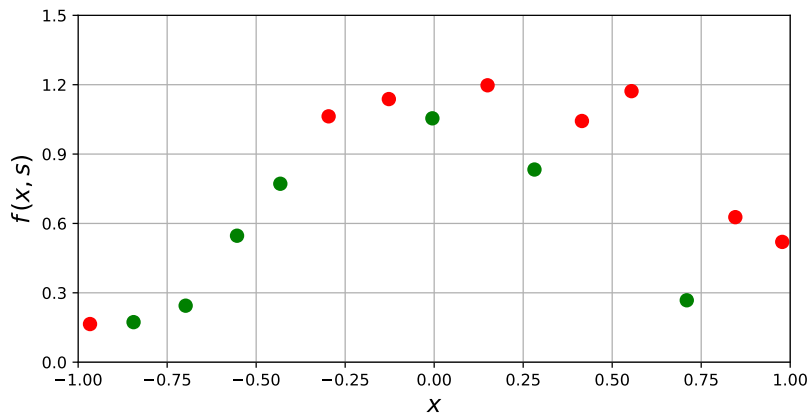
GPR with Qualitative Inputs

- Model a system as $f : \mathbb{R}^d \times S \rightarrow \mathbb{R}$ with $S = \{1, \dots, n_s\}$.
- Can we do better than n_s separate fits?
- If the n_s surfaces are correlated, we can use expanded kernel

$$k((\mathbf{x}, s), (\mathbf{x}', s')) = C_{s,s'} \tilde{k}(\mathbf{x}, \mathbf{x}'),$$

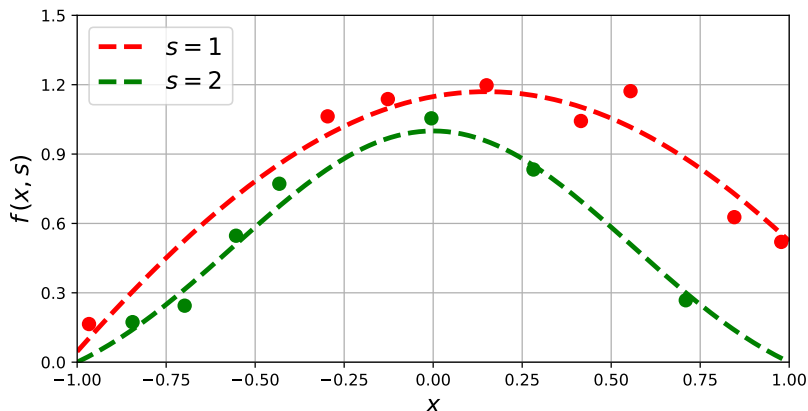
where $C \in \mathbb{R}^{n_s \times n_s}$ is a cross-correlation matrix.

A Mixed-Model Example



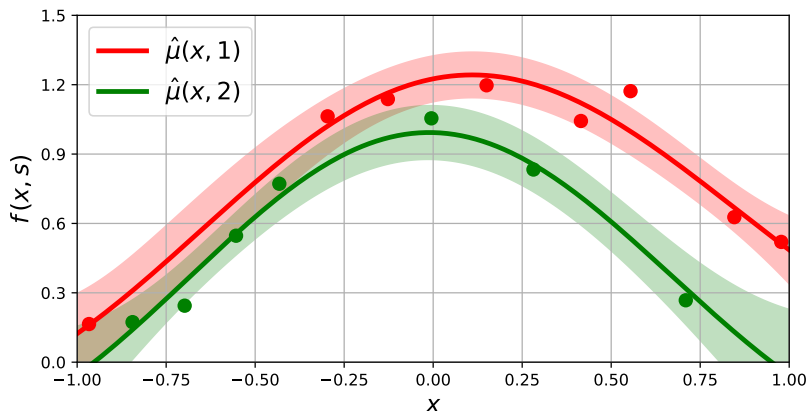
A Mixed-Model Example

$$f : [-1, 1] \times \{1, 2\} : (x, s) \mapsto \begin{cases} 1.9 \cos(x - 0.15) - 0.73 & : s = 1 \\ (1 - x^2)e^{-x^2} & : s = 2 \end{cases}$$



A Mixed-Model Example

$$f : [-1, 1] \times \{1, 2\} : (x, s) \mapsto \begin{cases} 1.9 \cos(x - 0.15) - 0.73 & : s = 1 \\ (1 - x^2)e^{-x^2} & : s = 2 \end{cases}$$

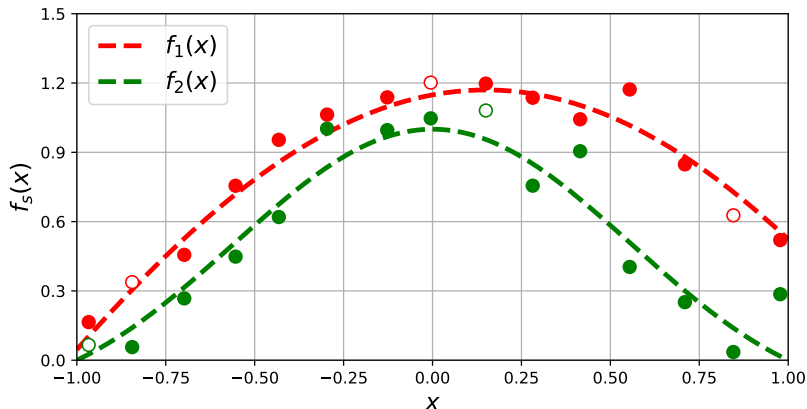


Multiple Regression

- The qualitative input scheme also works for multiple outputs:

$$\mathbf{y} = \mathbf{f}(\mathbf{x}) \in \mathbb{R}^{d'} \iff y_s = f(\mathbf{x}, s)$$

- Such an approach can accommodate missing output data.



References

- Frazier, P.I. “A Tutorial on Bayesian Optimization,” *arXiv Preprint*, 22 pp. (2018). [arXiv:1807.02811](https://arxiv.org/abs/1807.02811)
- Luong, P.; Gupta, S.; Nguyen, D.; Rana, S.; Venkatesh, S. “Bayesian Optimization with Discrete Variables,” in *AI 2019: Advances in Artificial Intelligence*, 473–84 (2019).
- Moćkus, J. “On Bayesian Methods for Seeking the Extremum,” in *Optimization Techniques: IFIP Technical Conference*, 400–04 (1975).
- Rasmussen, C.E.; Williams, C.K.I. *Gaussian Processes for Machine Learning*, MIT Press: Cambridge, MA (2006).
- Santner, T.J.; Williams, B.J.; Notz, W.I. *The Design and Analysis of Computer Experiments*, 2nd ed., Springer: New York, NY (2018).
- Shahriari, B.; Swersky, K.; Wang, Z.; Adams, R.P.; de Freitas, N. “Taking the Human Out of the Loop: A Review of Bayesian Optimization,” *Proc. IEEE*, 104(1), 148–75 (2015).