

# Gaussian Processes for Regression

# Outline

- Review of Bayesian estimation
- Review of Bayesian linear regression
- Gaussian process regression
- Bayesian optimization

# Bayesian Estimation

user-specified

- Maximum likelihood assumes that  $\theta$  is nonrandom
- In Bayesian estimation,  $\theta$  is viewed as **random**
- Our uncertainty about  $\theta$  depends on whether we have observed data or not
- Let  $\underline{Z} = (Z_1, \dots, Z_n)$  denote all observed quantities
- Let  $p(\theta)$  denote the distribution of  $\theta$  before observing a realization  $\underline{z}$  of  $\underline{Z}$ . **prior distribution**
- Let  $p(\theta | \underline{z})$  denote the distribution of  $\theta$  after observing a realization  $\underline{z}$  of  $\underline{Z}$ . **posterior distribution**
- These are related by Bayes rule:

$$p(\theta | \underline{z}) = \frac{p(\underline{z} | \theta) p(\theta)}{p(\underline{z})} = \int p(\underline{z} | \theta) p(\theta) d\theta$$

$(\underline{z}, \theta)$  jointly distributed

# Bayesian Estimation

- A parameter estimate can be obtained from the posterior in multiple ways:

- Posterior mean:  $\hat{\theta}(z) = \mathbb{E}[\theta | z] = \int \theta p(\theta | z) d\theta$

- Maximum a posteriori (MAP) estimation:

$$\begin{aligned}\hat{\theta}(z) &= \arg \max_{\theta} \log p(\theta | z) \\ &= \arg \max_{\theta} \log p(z | \theta) + \log p(\theta)\end{aligned}$$

- In addition to parameter estimates, the main advantage of Bayesian methods over frequentist methods is they also provide natural

confidence intervals.

# Bayesian Linear Regression

- Training data:  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ . View  $\mathbf{x}_i$  as fixed.
- No offset,  $\theta = \mathbf{w}$
- Gaussian likelihood:

$$y_i = \mathbf{w}^T \mathbf{x}_i + \epsilon_i, \quad \iff \quad y_i \sim N(\mathbf{w}^T \mathbf{x}_i, \sigma_e^2)$$

where  $\epsilon_i \sim N(0, \sigma_e^2)$

- Prior:

$$\mathbf{w} \sim N(\mathbf{0}, \sigma_w^2 \mathbf{I}) \iff p(\mathbf{w}) = (2\pi\sigma_w^2)^{-\frac{d}{2}} \exp\left\{-\frac{\|\mathbf{w}\|^2}{2\sigma_w^2}\right\}$$

$$\iff \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}, \quad w_i \stackrel{iid}{\sim} N(0, \sigma_w^2)$$

- $\sigma_e^2, \sigma_w^2$  are hyperparameters

# Bayesian Linear Regression

- Posterior:

$$w|y \sim N(\mu(y), \Sigma(y))$$

where

$$\hat{\omega} = \mu(y) = \left( \mathbf{X} \mathbf{X}^T + \frac{\sigma_e^2}{\sigma_w^2} \mathbf{I} \right)^{-1} \mathbf{X} y$$

$$\Sigma(y) = \left( \frac{1}{\sigma_e^2} \mathbf{X} \mathbf{X}^T + \frac{1}{\sigma_w^2} \mathbf{I} \right)^{-1}$$

and

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n \end{bmatrix} \in \mathbb{R}^{d \times n}$$

# Bayesian Linear Regression

- Linear transformation of a MVG is another MVG
- In particular, if  $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , and  $\mathbf{A}$  is a matrix and  $\mathbf{b}$  a column vector of appropriate dimensions, then

$$\mathbf{AZ} + \mathbf{b} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$$

- If  $\mathbf{x}$  is a test point, then the predicted output

$$\hat{f}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \mathbf{x}^T \mathbf{w}$$

- Thus, the predicted output is a random variable with distribution

$$\hat{f}(\mathbf{x}) | \mathbf{y} \sim \mathcal{N}(\mathbf{x}^T \boldsymbol{\mu}(\mathbf{y}), \mathbf{x}^T \boldsymbol{\Sigma}(\mathbf{y}) \mathbf{x})$$

$$\sim \mathcal{N}\left(\underbrace{\mathbf{x}^T (\mathbf{X}\mathbf{X}^T + \frac{\sigma_e^2}{\sigma_w^2} \mathbf{I})^{-1} \mathbf{X}_y}_{M_x(y)}, \underbrace{\mathbf{x}^T (\mathbf{X}\mathbf{X}^T + \frac{\sigma_e^2}{\sigma_w^2} \mathbf{I})^{-1} \mathbf{x}}_{\sigma_x^2(y)}\right)$$

# Bayesian Linear Regression

- Using matrix identities (e.g., the matrix inversion lemma), it can be shown that

$$\hat{f}(\mathbf{x}) \mid \mathbf{y} \sim \mathcal{N}(\mu_{\mathbf{x}}(\mathbf{y}), \sigma_{\mathbf{x}}^2(\mathbf{y}))$$

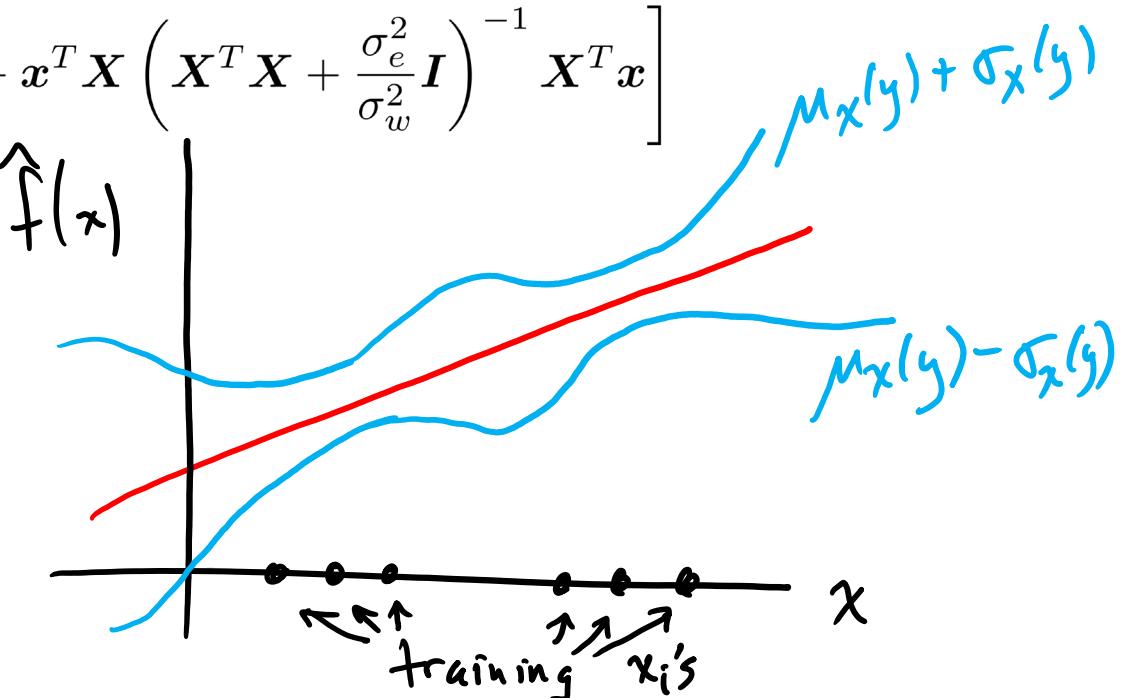
where

$$\mu_{\mathbf{x}}(\mathbf{y}) = \sigma_w^2 \mathbf{x}^T \mathbf{X} \left( \mathbf{X}^T \mathbf{X} + \frac{\sigma_e^2}{\sigma_w^2} \mathbf{I} \right)^{-1} \mathbf{y}$$

$$\sigma_{\mathbf{x}}^2(\mathbf{y}) = \sigma_w^2 \left[ \mathbf{x}^T \mathbf{x} - \mathbf{x}^T \mathbf{X} \left( \mathbf{X}^T \mathbf{X} + \frac{\sigma_e^2}{\sigma_w^2} \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{x} \right]$$

- Conclusion:

You can  
kernelize  
this  
method



# Kernel Bayesian Linear Regression

- Let  $k$  be a symmetric, positive definite kernel
- Then Bayesian linear regression can be kernelized, leading to the prediction

$$\hat{f}(\mathbf{x}) \mid \mathbf{y} \sim \mathcal{N}(\mu_{\mathbf{x}}(\mathbf{y}), \sigma_{\mathbf{x}}^2(\mathbf{y}))$$

where

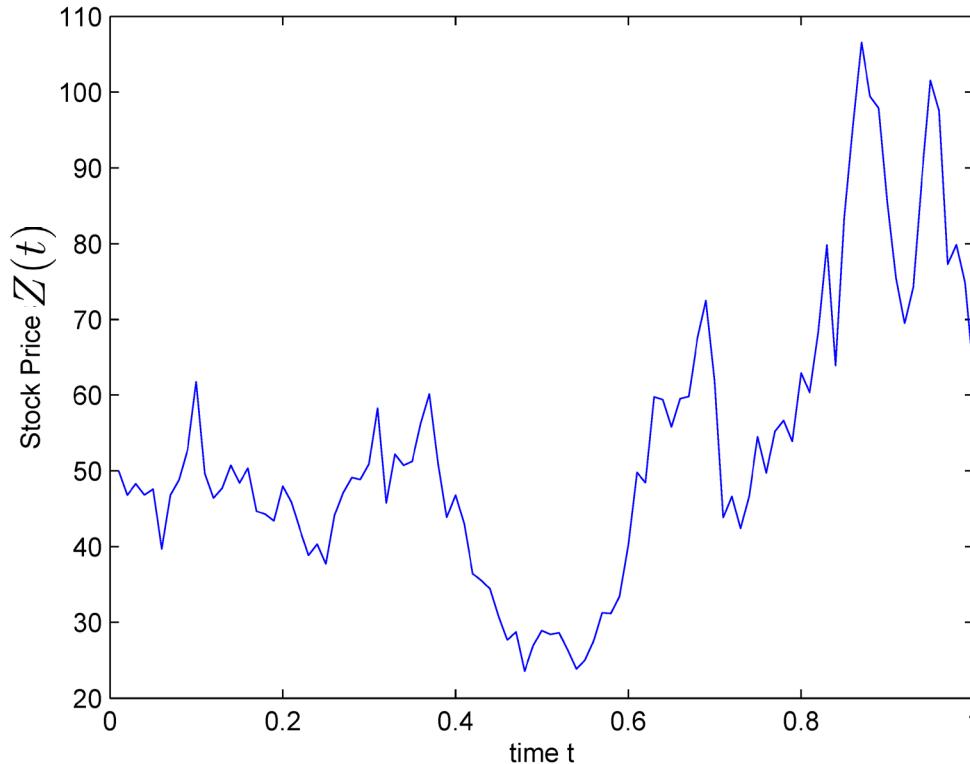
$$\mu_{\mathbf{x}}(\mathbf{y}) = \sigma_w^2 \mathbf{k}(\mathbf{x})^T \left( \mathbf{K} + \frac{\sigma_e^2}{\sigma_w^2} \mathbf{I} \right)^{-1} \mathbf{y}$$

$$\sigma_{\mathbf{x}}^2(\mathbf{y}) = \sigma_w^2 \left[ k(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x})^T \left( \mathbf{K} + \frac{\sigma_e^2}{\sigma_w^2} \mathbf{I} \right)^{-1} \mathbf{k}(\mathbf{x}) \right]$$

$$\mathbf{k}(\mathbf{x}) = \begin{bmatrix} k(x, x_1) \\ \vdots \\ k(x, x_n) \end{bmatrix}, \quad \mathbf{K} = \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix}$$

# Random Processes

- A *random process* is a family of (scalar) random variables  $\{Z(t)\}$ ,  $t \in T$ , where  $T$  is an index set (e.g.,  $T = \mathbb{R}$ )
- A random process is *zero-mean* if  $\mathbb{E}[Z(t)] = 0$  for all  $t$ .



# Gaussian Processes

$\{Z(x)\}_{x \in \mathbb{R}^d} \leftrightarrow \tilde{Z}(x) = \text{random function on } \mathbb{R}^d$

- Switch notation:  $t \in T$  to  $x \in \mathbb{R}^d$
- A *Gaussian process* is a random process, indexed by  $\mathbb{R}^d$ , such that for any finite set of indices  $x_1, \dots, x_N \in \mathbb{R}^d$ , the values  $Z(x_1), \dots, Z(x_N) \in \mathbb{R}$  are jointly Gaussian.
- A zero-mean Gaussian process is completely characterized by its *covariance function*

$$k(x, x') = \mathbb{E}[Z(x)Z(x')]$$

- Thus, for any  $x_1, \dots, x_N$ ,

$$\begin{bmatrix} Z(x_1) \\ \vdots \\ Z(x_N) \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_N) \\ \vdots & \ddots & \vdots \\ k(x_N, x_1) & \cdots & k(x_N, x_N) \end{bmatrix} \right)$$

# Covariance Function

- To be a valid covariance function,  $k$  must satisfy
  - $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$  for all  $\mathbf{x}, \mathbf{x}'$
  - For all  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ , the matrix

$$\begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

must be positive semi-definite.

- In other words,  $k$  must be a *symmetric, positive definite kernel*
- Examples:

Gaussian:  $k(\mathbf{x}, \mathbf{x}') = \sigma_k^2 \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$

Laplacian:  $k(\mathbf{x}, \mathbf{x}') = \sigma_k^2 \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|)$

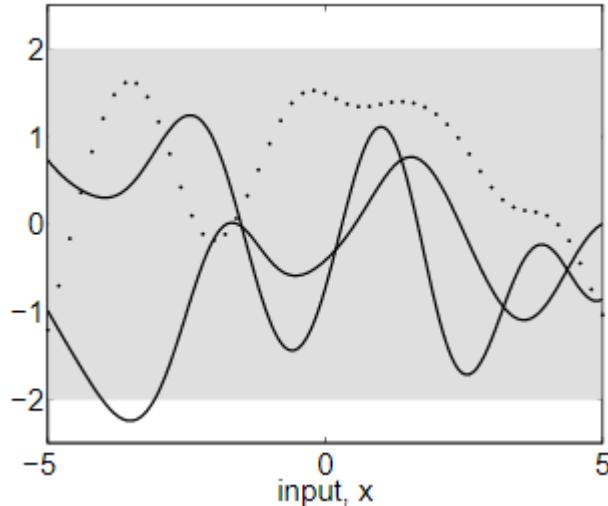
# Smooth Realizations

- A realization of a Gaussian process, i.e., a random function  $Z(x)$  defined on a domain  $\mathcal{X} \subseteq \mathbb{R}^d$ , is a continuous function provided

$\mathcal{X}$  is a compact set (e.g., a closed rectangle)

$k(x, x')$  is Lipschitz continuous (e.g., if  $k$  is differentiable)

- To plot a realization, choose a grid of points  $x_1, \dots, x_n$  at which to plot the function, and simulate a realization of



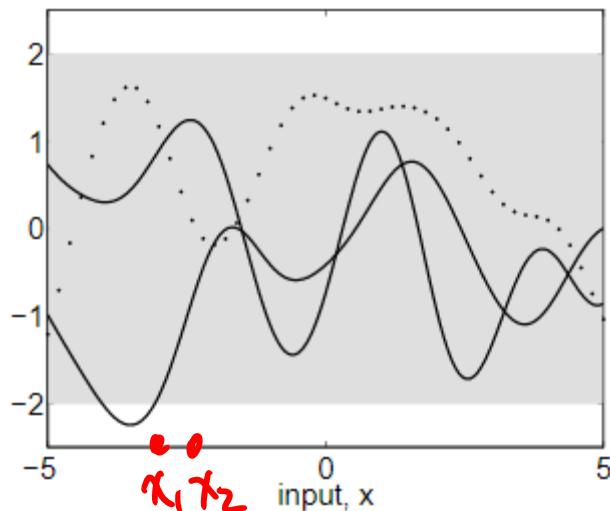
$$\begin{bmatrix} Z(x_1) \\ \vdots \\ Z(x_n) \end{bmatrix} = \mathcal{N} \left( \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} k(x_1, x_1) & \cdots & \\ & \ddots & \\ & & k(x_n, x_n) \end{bmatrix} \right)$$

# Poll

$$\text{cov}(z(x), z(x')) \\ //$$

Consider the Gaussian covariance  $k(x, x') = \exp(-\gamma \|x - x'\|^2)$ . As  $\gamma$  increases, realizations of the associated GP become

- (A) Smoother (less wiggly)
- (B) Less smooth (more wiggly) ✓



$$\begin{bmatrix} z(x_1) \\ z(x_2) \end{bmatrix} \sim N \left( 0, \begin{bmatrix} 1 & e^{-\gamma \|x_1 - x_2\|^2} \\ e^{-\gamma \|x_1 - x_2\|^2} & 1 \end{bmatrix} \right)$$

# GPs for Regression

- Function to be estimated:  $f(\mathbf{x})$
- Assumption:  $f$  is a realization of a Gaussian process (can be viewed as a prior distribution on  $f$ )
- This is a *nonparametric* Bayesian setup

- Training data:  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

• Gaussian likelihood:  $y_i = f(\mathbf{x}_i) + \varepsilon_i, \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2_e)$

- Fixed test point:  $\mathbf{x}_{m+1}, \dots, \mathbf{x}_{m+n}$

- Predicted values at test points: Given  $y_1, \dots, y_n$ , estimate

$$f(\mathbf{x}_{m+i}), i=1, \dots, n.$$

- Goal: Determine posterior distribution of

$$\mathbf{f} = \begin{bmatrix} f(\mathbf{x}_{m+1}) \\ \vdots \\ f(\mathbf{x}_{m+n}) \end{bmatrix} \quad \left| \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \right.$$

# GPs for Regression

- The posterior is just the conditional distribution of  $\mathbf{f}$  given  $\mathbf{y}$ .
- We will use the Gaussian conditioning property:

If  $\mathbf{A} \in \mathbb{R}^p$  and  $\mathbf{B} \in \mathbb{R}^q$  are such that

$$\begin{bmatrix} \mathbf{A} \\ \mathbf{B} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \boldsymbol{\mu}_A \\ \boldsymbol{\mu}_B \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{AA} & \boldsymbol{\Sigma}_{AB} \\ \boldsymbol{\Sigma}_{BA} & \boldsymbol{\Sigma}_{BB} \end{bmatrix} \right),$$

then

$$\mathbf{A} | \mathbf{B} = \mathbf{b} \sim \mathcal{N} \left( \boldsymbol{\mu}_A + \boldsymbol{\Sigma}_{AB} \boldsymbol{\Sigma}_{BB}^{-1} (\mathbf{b} - \boldsymbol{\mu}_B), \boldsymbol{\Sigma}_{AA} - \boldsymbol{\Sigma}_{AB} \boldsymbol{\Sigma}_{BB}^{-1} \boldsymbol{\Sigma}_{BA} \right).$$

$$y_j = f(x_j) + \varepsilon_j$$

# GPs for Regression

- Introduce the notation

$K$  = train-train kernel matrix ( $m \times m$ )

$K_*$  = train-test kernel matrix ( $m \times n$ )

$K_{**}$  = test-test kernel matrix ( $n \times n$ )

( $i,j$ ) entry  
 $\downarrow$   
 $k(x_i, x_j)$   
 $k(x_i, x_{m+j})$   
 $k(x_{m+i}, x_{m+j})$

- Covariances

- $\text{Cov}(f(x_{m+i}), f(x_{m+j})) = k(x_{m+i}, x_{m+j})$

- $\text{Cov}(f(x_{m+i}), y_j) = \text{Cov}(f(x_{m+i}), f(x_j) + \varepsilon_j) = \text{Cov}(f(x_{m+i}), f(x_j))$

- $\text{Cov}(y_i, y_j) = \text{Cov}(f(x_i) + \varepsilon_i, f(x_j) + \varepsilon_j) = k(x_i, x_j) + \sigma_e^2 \mathbf{1}_{\{i=j\}}$

- The joint distribution of  $f$  and  $y$  is

$$\begin{bmatrix} f \\ y \end{bmatrix} \sim N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} K_{**} & K_*^\top \\ K_* & K + \sigma_e^2 I \end{bmatrix} \right)$$

# GPs for Regression

- Therefore the posterior distribution of the predicted test outputs is

$$\mathbf{f} \mid \mathbf{y} \sim \mathcal{N}(\mathbf{K}_*^T(\mathbf{K} + \sigma_e^2 \mathbf{I})^{-1} \mathbf{y}, \mathbf{K}_{**} - \mathbf{K}_*^T(\mathbf{K} + \sigma_e^2 \mathbf{I})^{-1} \mathbf{K}_*).$$

- Special case of  $n = 1$     *test point:  $\mathbf{x}$*

$$f(\mathbf{x}) \mid \mathbf{y} \sim \mathcal{N}\left(k(\mathbf{x})^T (\mathbf{K} + \sigma_e^2 \mathbf{I})^{-1} \mathbf{y},\right.$$

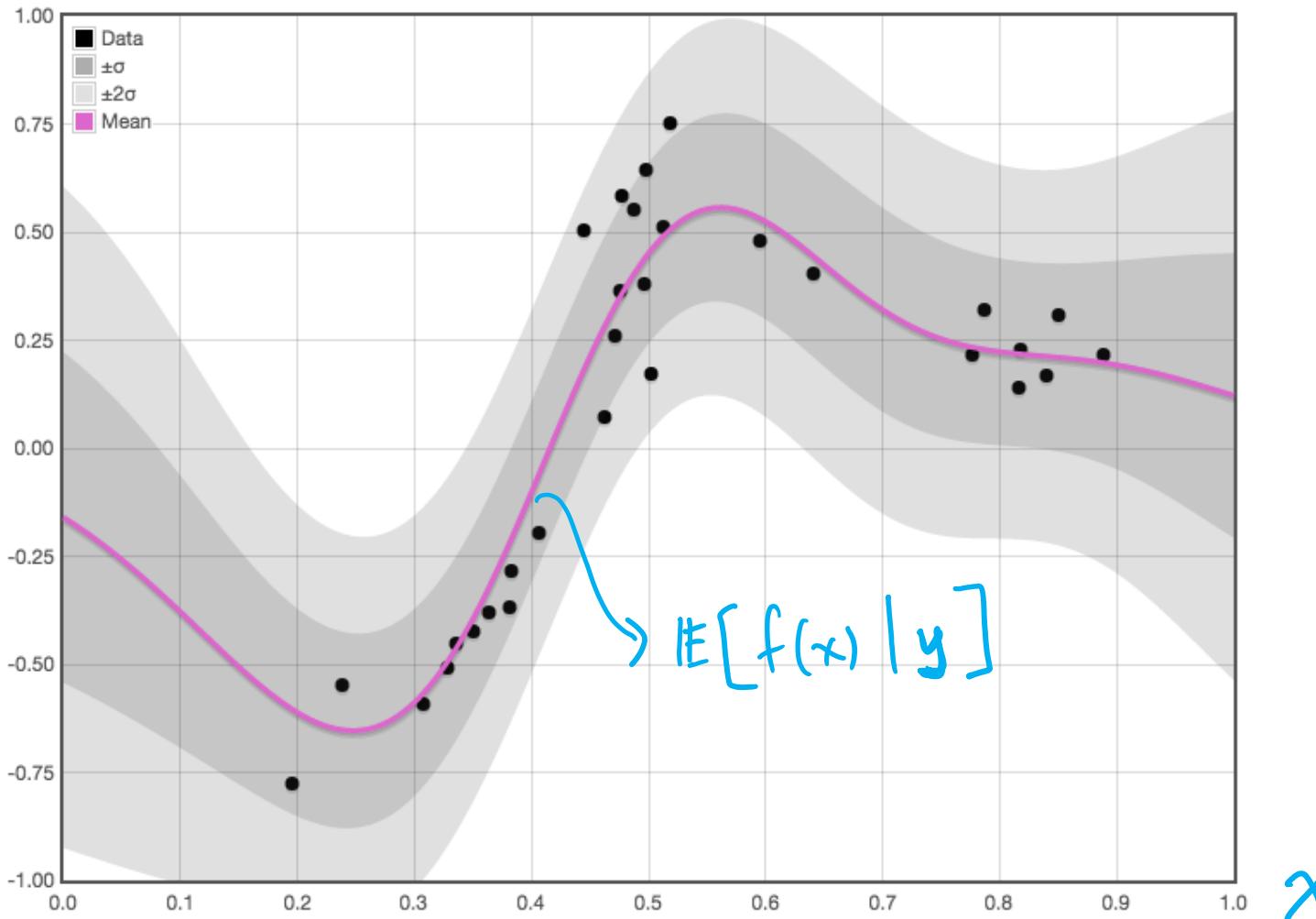
$$\left. k(\mathbf{x}, \mathbf{x}) - k(\mathbf{x})^T (\mathbf{K} + \sigma_e^2 \mathbf{I})^{-1} k(\mathbf{x})\right)$$

where

$$k(\mathbf{x}) = \begin{bmatrix} k(\mathbf{x}, \mathbf{x}_1) \\ \vdots \\ k(\mathbf{x}, \mathbf{x}_m) \end{bmatrix}$$

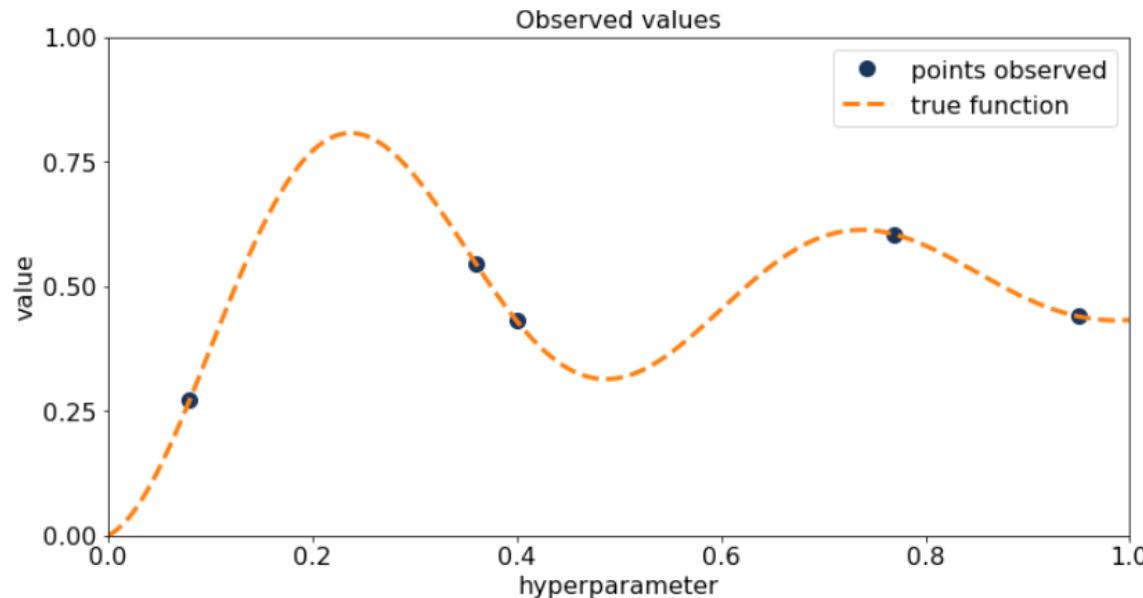
- This agrees with our kernelized version of Bayesian linear regression
- While the estimator (posterior mean/mode) is a linear function of  $\mathbf{y}$ , it is a nonlinear function of  $\mathbf{x}_1, \dots, \mathbf{x}_{m+n}$ .

# GPs for Regression



# Bayesian Optimization

- Method for *black box optimization*, when objective function is hard to evaluate and can't calculate derivatives
- Need to determine next point to sample
- Application: model selection in machine learning (e.g., objective function is cross-validation accuracy estimate of model performance)

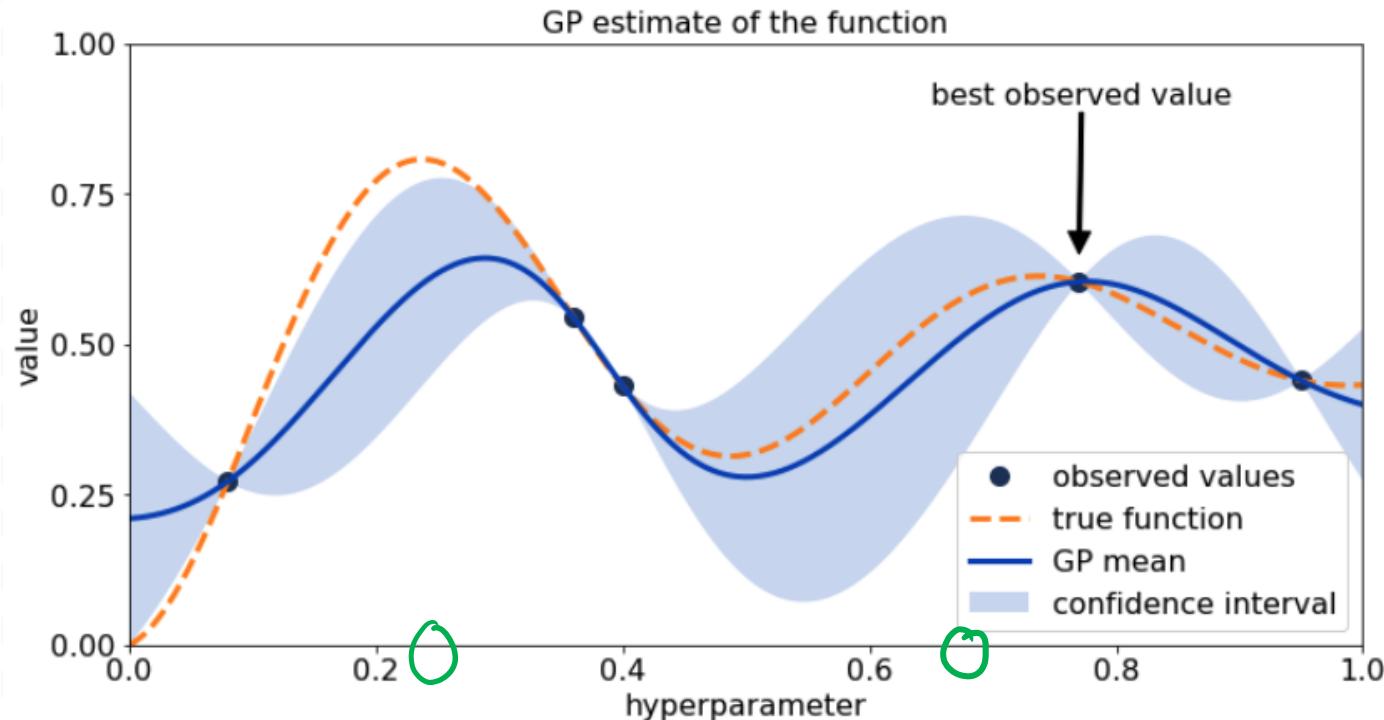


<https://sigopt.com/blog/bayesian-optimization-101/>

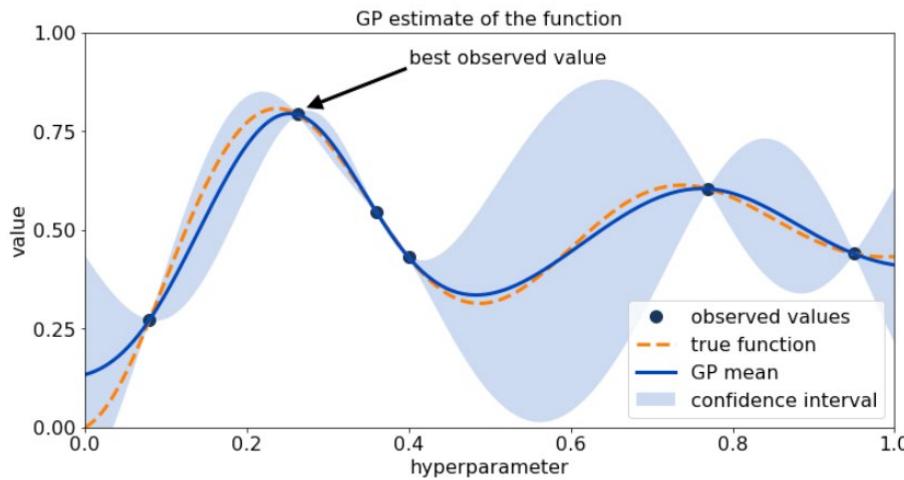
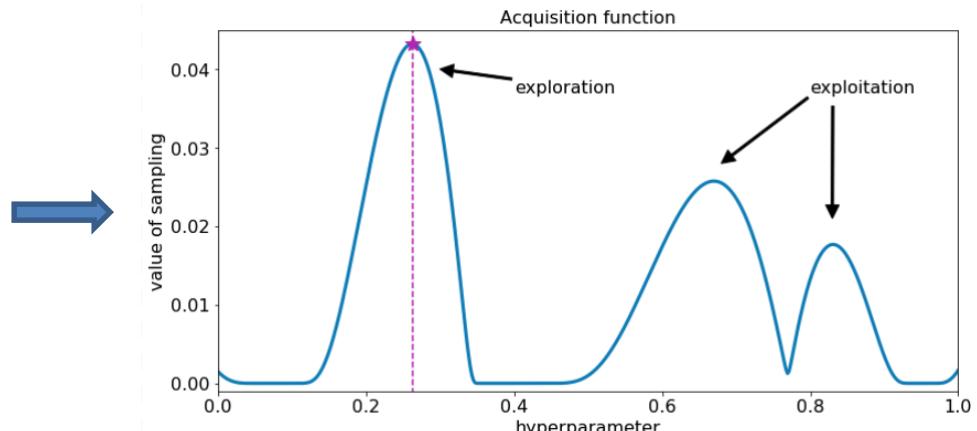
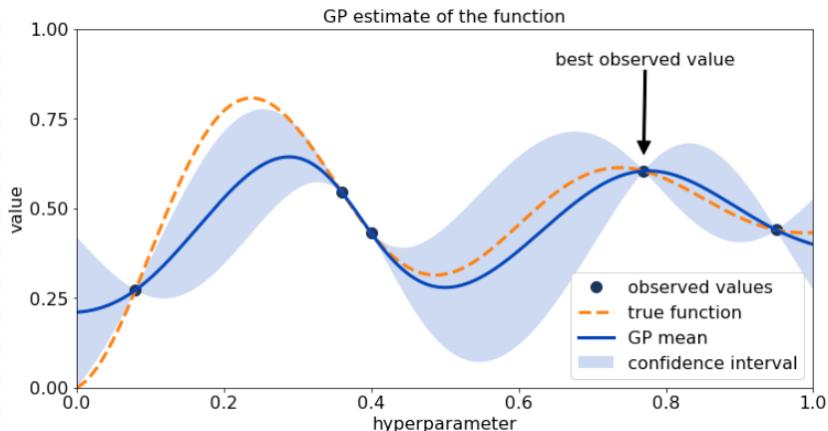
# Bayesian Optimization

- Assume a GP prior on the unknown function
- Use posterior to characterize uncertainty in the unknown function, and determine next point to evaluate

$$\sigma_e^2 = 0$$



# Bayesian Optimization



# Final Thoughts on GPs

- Free parameters can be tuned by maximum likelihood estimation – much faster than cross-validation (although possibly not as accurate)
- In regression, we used a GP prior and a Gaussian likelihood. The math works out nicely because the Gaussian prior is *conjugate* to the Gaussian likelihood.
- GPs can be used with the logistic likelihood as a nonlinear model for classification.
- This gives a nonparametric Bayesian version of kernel logistic regression
- Math is more challenging because Gaussian prior is not conjugate to the logistic likelihood
- Bayesian optimization: black-box optimization without grid search