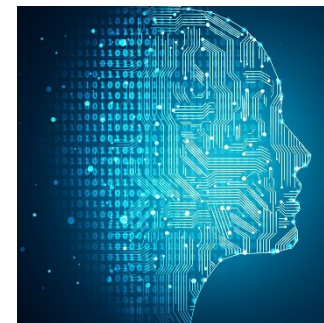


# Machine Learning

# Gaussian Processes



Kevin Moon (kevin.moon@usu.edu)  
STAT/CS 5810/6655



# Motivation



- Apply Gaussian random processes + Bayesian estimation theory to regression
- **Result:** probabilistic version of kernel ridge regression
  - Enables uncertainty quantification
- Let's cover Bayesian estimation first

# Bayesian Estimation



- Suppose we have parameters we want to estimate
  - E.g. regression coefficients or neural network parameters
- We observe  $\mathbf{Z} \sim p(\mathbf{z}; \boldsymbol{\theta})$
- **Objective:** estimate  $\boldsymbol{\theta}$  that best explains the observations
- In Bayesian estimation, we assume  $\boldsymbol{\theta}$  is random with prior distribution  $\boldsymbol{\theta} \sim p(\boldsymbol{\theta})$
- The posterior distribution (by Bayes rule) is:

$$p(\boldsymbol{\theta}|\mathbf{z}) = \frac{p(\mathbf{z}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{z})}$$

# Bayesian estimation



- The posterior distribution (by Bayes rule) is:

$$p(\boldsymbol{\theta}|\mathbf{z}) = \frac{p(\mathbf{z}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{z})}$$

- Bayesian parameter estimates are based on the posterior
  - Posterior mode:  $\hat{\boldsymbol{\theta}}(\mathbf{z}) = \arg \min_{\boldsymbol{\theta}} p(\boldsymbol{\theta}|\mathbf{z})$
  - Posterior mean:  $\hat{\boldsymbol{\theta}}(\mathbf{z}) = \mathbb{E}[\boldsymbol{\theta}|\mathbf{Z} = \mathbf{z}] = \int \boldsymbol{\theta} p(\boldsymbol{\theta}|\mathbf{z}) d\boldsymbol{\theta}$

# Random Processes



- In most Bayesian estimation problems,  $\theta$  is finite dimensional
- In our case, the “parameter” is the unknown regression function  $f$  in

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

- $\mathbf{x} \in \mathbb{R}^d, y \in \mathbb{R}, \epsilon \in \mathbb{R}$
- We’ll use a *random process* as a prior for  $f$
- Random process = a family of random variables indexed by a potentially uncountably infinite variable
  - $\mathbb{R}^d$  is the index in our case

# Gaussian Processes



- Gaussian process = a collection of variables, where any finite number of them has a multivariate Gaussian distribution
- Gaussian process completely specified by:
  1. Mean function  $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$
  2. Covariance function
$$k(\mathbf{x}, \mathbf{x}') = \text{cov}(f(\mathbf{x}), f(\mathbf{x}'))$$
$$= \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$
- Notation:  $f(\mathbf{x}) \sim GP(m, k)$
- How could you plot a random function like this?

# Gaussian Processes



- Restrictions on  $k(\mathbf{x}, \mathbf{x}')$ :
  1. Symmetric:  $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$
  2. Positive semi-definite. I.e. for any finite collection of points  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ , the matrix formed by  $k(\mathbf{x}_i, \mathbf{x}_j)$  must be PSD
- No restrictions on  $m(\mathbf{x})$
- We'll need the *Gaussian conditioning property*:
  - If  $\mathbf{a} \in \mathbb{R}^p$  and  $\mathbf{b} \in \mathbb{R}^q$  and
$$\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{m}_a \\ \mathbf{m}_b \end{bmatrix}, \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix} \right),$$
$$\Rightarrow \mathbf{a} | \mathbf{b} \sim \mathcal{N}(\mathbf{m}_a + \Sigma_{ab} \Sigma_{bb}^{-1} (\mathbf{b} - \mathbf{m}_b), \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba})$$

# Noise-free case



## Given

- Training inputs  $\mathbf{x}_1, \dots, \mathbf{x}_m$
- Training outputs  $y_i = f(\mathbf{x}_i)$
- Test inputs  $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$
- Mean function  $m(\mathbf{x})$
- Covariance function  $k(\mathbf{x}, \mathbf{x}')$

## Goal

- Predict the test outputs  $f(\tilde{\mathbf{x}}_i)$



# Noise-free case



## Given

- Training inputs  $\mathbf{x}_1, \dots, \mathbf{x}_m$
- Training outputs  $y_i = f(\mathbf{x}_i)$
- Test inputs  $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$
- Mean function  $m(\mathbf{x})$  and covariance function  $k(\mathbf{x}, \mathbf{x}')$

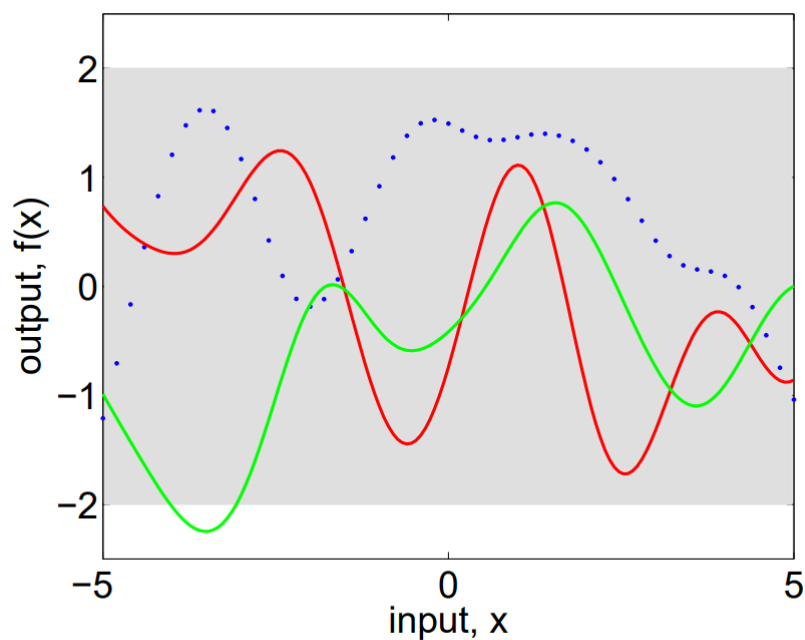
$$\mathbf{f} = [f(\mathbf{x}_1) \quad \dots \quad f(\mathbf{x}_m)]^T$$
$$\tilde{\mathbf{f}} = [f(\tilde{\mathbf{x}}_1) \quad \dots \quad f(\tilde{\mathbf{x}}_n)]^T$$

- $\begin{bmatrix} \mathbf{f} \\ \tilde{\mathbf{f}} \end{bmatrix}$  is multivariate Gaussian distributed
- Compute the posterior distribution of  $\tilde{\mathbf{f}} | \mathbf{f}$  using the Gaussian conditioning lemma

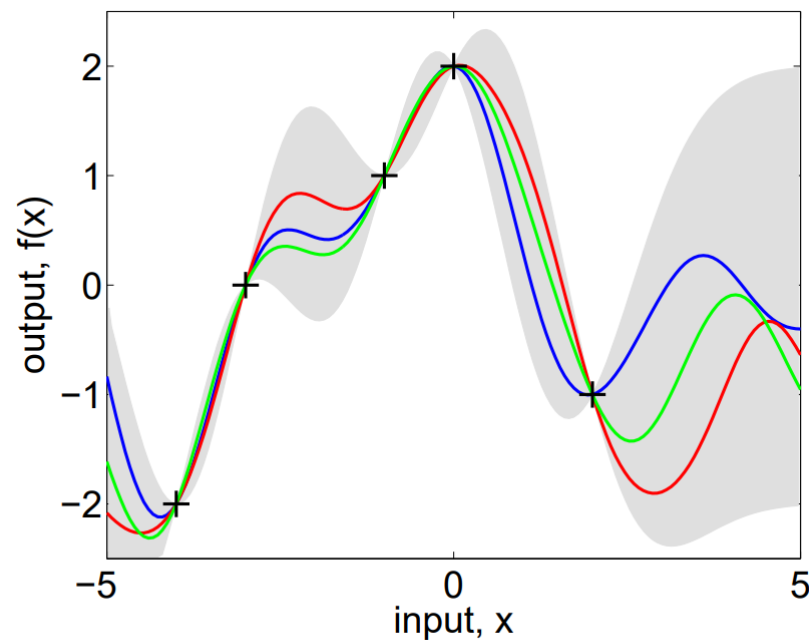
# Noise free case: Example



- $m(x) = 0, k(x, x') = \sigma_g^2 \exp\left(-\frac{\|x - x'\|^2}{2\sigma_s^2}\right)$



(a), prior



(b), posterior

- 3 samples from the prior and posterior with 95% CI
- What effect does  $\sigma_s^2$  have on the functions?

# Noisy case



## Given

- Training inputs  $\mathbf{x}_1, \dots, \mathbf{x}_m$
- Noisy training outputs  $y_i = f(\mathbf{x}_i) + \epsilon_i$
- Test inputs  $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$
- $m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'), \sigma_\epsilon^2$

## Assume

- $\epsilon_i \sim \mathcal{N}(0, \sigma_\epsilon^2)$

## Goal

- Predict  $\tilde{f}_i = f(\tilde{\mathbf{x}}_i)$

# Noisy case



## Given

- Training inputs  $\mathbf{x}_1, \dots, \mathbf{x}_m$
- Noisy training outputs  $y_i = f(\mathbf{x}_i) + \epsilon_i$
- Test inputs  $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$
- $m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'), \sigma_\epsilon^2$
- $\begin{bmatrix} \mathbf{y} \\ \tilde{\mathbf{f}} \end{bmatrix}$  is multivariate Gaussian distributed
- Compute the posterior distribution of  $\tilde{\mathbf{f}} | \mathbf{y}$  using the Gaussian conditioning lemma

# Notation



- $X = [\mathbf{x}_1 \quad \dots \quad \mathbf{x}_m] \in \mathbb{R}^{d \times m}$
- $\tilde{X} = [\tilde{\mathbf{x}}_1 \quad \dots \quad \tilde{\mathbf{x}}_n] \in \mathbb{R}^{d \times n}$
- $\mathbf{f} = [f(\mathbf{x}_1) \quad \dots \quad f(\mathbf{x}_m)]^T$
- $\tilde{\mathbf{f}} = [f(\tilde{\mathbf{x}}_1) \quad \dots \quad f(\tilde{\mathbf{x}}_n)]^T$
- $\mathbf{y} = [y_1 \quad \dots \quad y_m]^T$
- $\mathbf{m} = [m(\mathbf{x}_1) \quad \dots \quad m(\mathbf{x}_m)]^T$
- $\tilde{\mathbf{m}} = [m(\tilde{\mathbf{x}}_1) \quad \dots \quad m(\tilde{\mathbf{x}}_n)]^T$
- $K(X, X) = [k(\mathbf{x}_i, \mathbf{x}_j)]_{i,j=1}^m \in \mathbb{R}^{m \times m}$
- $K(X, \tilde{X}) = [k(\mathbf{x}_i, \tilde{\mathbf{x}}_j)]_{i,j=1}^{m,n} \in \mathbb{R}^{m \times n}$
- $K(\tilde{X}, \tilde{X}) = [k(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j)]_{i,j=1}^n \in \mathbb{R}^{n \times n}$

# Noisy case



$$\begin{bmatrix} \tilde{\mathbf{f}} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \tilde{\mathbf{m}} \\ \mathbf{m} \end{bmatrix}, \begin{bmatrix} K(\tilde{X}, \tilde{X}) & K(\tilde{X}, X) \\ K(X, \tilde{X}) & K(X, X) + \sigma_\epsilon^2 I \end{bmatrix} \right)$$

By the Gaussian conditioning lemma,  $\tilde{\mathbf{f}} | \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$

- $\boldsymbol{\mu} = \tilde{\mathbf{m}} + K(X, \tilde{X})^T [K(X, X) + \sigma_\epsilon^2 I]^{-1} (\mathbf{y} - \mathbf{m})$
- $\Sigma = K(\tilde{X}, \tilde{X}) - K(X, \tilde{X})^T [K(X, X) + \sigma_\epsilon^2 I]^{-1} K(X, \tilde{X})$
- Due to symmetry, the estimate of  $\tilde{\mathbf{f}}$  is:

$$\hat{\tilde{\mathbf{f}}} = \tilde{\mathbf{m}} + K(X, \tilde{X})^T [K(X, X) + \sigma_\epsilon^2 I]^{-1} (\mathbf{y} - \mathbf{m})$$

- This predictor is linear in  $\mathbf{y}$  but nonlinear in  $\mathbf{x}_i$  and  $\tilde{\mathbf{x}}_j$

# Connection to KRR



- Assume from now on that  $m(\mathbf{x}) = 0 \Rightarrow \mathbf{m}, \tilde{\mathbf{m}} = \mathbf{0}$
- Denote  $\boldsymbol{\alpha} = [K(X, X) + \sigma_\epsilon^2 I]^{-1} \mathbf{y}$
- Consider a single test point  $\mathbf{x}$
- The GP prediction is

$$\hat{f}(\mathbf{x}) = K(X, \mathbf{x})\boldsymbol{\alpha} = \sum_{i=1}^m \alpha_i k(\mathbf{x}_i, \mathbf{x})$$

- This is identical to the KRR estimate with  $\lambda = \sigma_\epsilon^2$
- Have we gained anything?

# Confidence bands



- GP predictor comes with a variance:

$$\begin{aligned} & \text{var} \left( \hat{f}(\mathbf{x}) \right) \\ &= K(\mathbf{x}, \mathbf{x}) - K(X, \mathbf{x})^T [K(X, X) + \sigma_\epsilon^2 I]^{-1} K(X, \mathbf{x}) \end{aligned}$$

- Can place a confidence band around the estimate



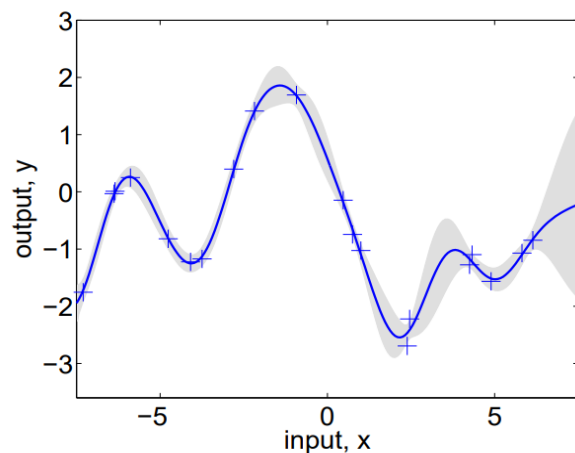
# Noisy case: Example



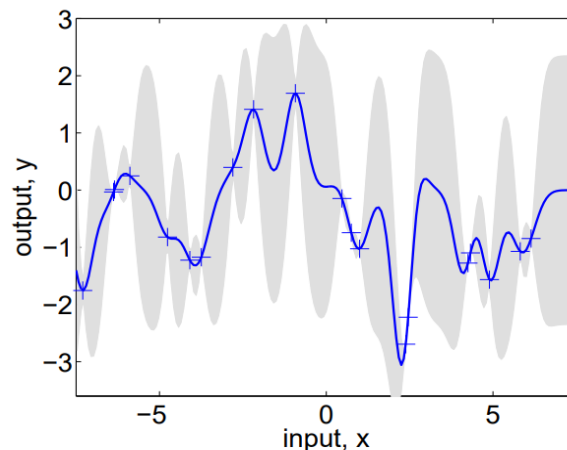
a)  $\sigma_{\epsilon}^2 = 0.1, \sigma_s^2 = 1, \sigma_g^2 = 1$

b)  $\sigma_{\epsilon}^2 = 0.00005, \sigma_s^2 = 0.3, \sigma_g^2 = 1.08$

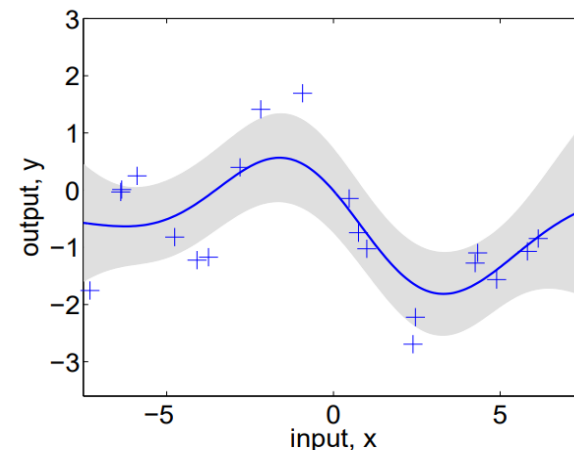
c)  $\sigma_{\epsilon}^2 = 0.89, \sigma_s^2 = 3, \sigma_g^2 = 1.16$



(a),  $\ell = 1$



(b),  $\ell = 0.3$



(c),  $\ell = 3$

# Conclusion



- Gaussian processes can be applied to obtain a Bayesian version of kernel ridge regression
- Can derive it via an alternative way: start with Bayesian linear regression, and then apply the kernel trick
  - See [https://web.eecs.umich.edu/~cscott/past\\_courses/eecs545f16/32\\_gaussian\\_processes.pdf](https://web.eecs.umich.edu/~cscott/past_courses/eecs545f16/32_gaussian_processes.pdf)
- Gaussian processes can also be applied to classification

# Further reading



- Rasmussen and Williams, Gaussian Processes for Machine Learning, 2006