

Machine Learning and Data Analytics

ME 5013- Fall 2019

Lecture 07

- Gaussian Process



The University of Texas at San Antonio™

Adel Alaeddini, PhD

Associate Professor of Mechanical Engineering

Advanced Data Engineering Lab

adel.alaeddini@utsa.edu

- Here are some data points! What function did they come from?
 - I have *no idea*.
- Oh. Okay. Uh, you think this point is likely in the function too?
 - I have *no idea*.

- You can't get anywhere without making some assumptions
- GPs are a nice way of expressing this 'prior on functions' idea.
- Can do a bunch of cool stuff
 - Regression
 - Classification
 - Optimization

- Having some observed input-output pairs (\mathbf{x}_i, y_i) where y_i might be corrupted by some noise ε_i

$$y_i = f(\mathbf{x}_i) + \varepsilon_i \text{ for } i = 1, \dots, n$$

Scalar Vector, i.e. $\mathbf{x}_i = (x_i^{(1)}, \dots, x_i^{(p)})$

- ε_i is the additive independent identically distributed Gaussian noise with variance σ_n^2
- The underlying function f is not known (black-box function)
- The functional evaluation at the test (new) point $\mathbf{x} \in X$ is denoted as f_* (or y_* or $f(\mathbf{x}_*)$)

Matrix, i.e. $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)$

Gaussian process is a collection of random variables, a finite number of which have a joint Gaussian distribution

Input:

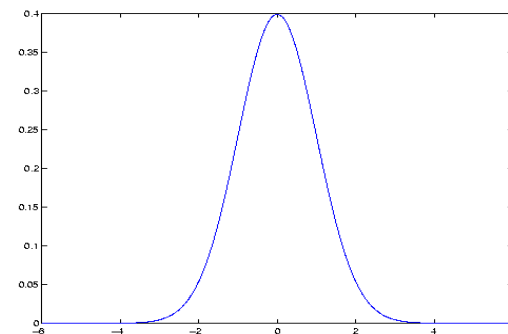
- ✓ Training set : $\{(\mathbf{x}_i, y_i), i = 1, 2, \dots, n\}$ from $y_i = f(\mathbf{x}_i) + \varepsilon_i$
- ✓ Test set X_* - Get predictions

Can be more than a single point $\mathbf{x}_{*1}, \dots, \mathbf{x}_{*k}$

- According to the joint distribution of training outputs and test outputs we have,

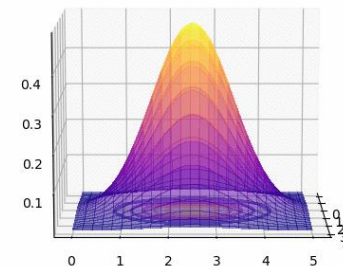
$$\begin{bmatrix} y \\ f_* \end{bmatrix} \sim \mathcal{N} \left(0, \begin{bmatrix} K(X, X) + \sigma_n^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix} \right)$$

Kernel function between pair of points



- Squared Exponential Kernel is the most popular and highly used kernel

$$K(x, x') = \sigma_f^2 \exp\left(-\frac{1}{2l^2} (x - x')^2\right)$$



- By conditional distribution we get,

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

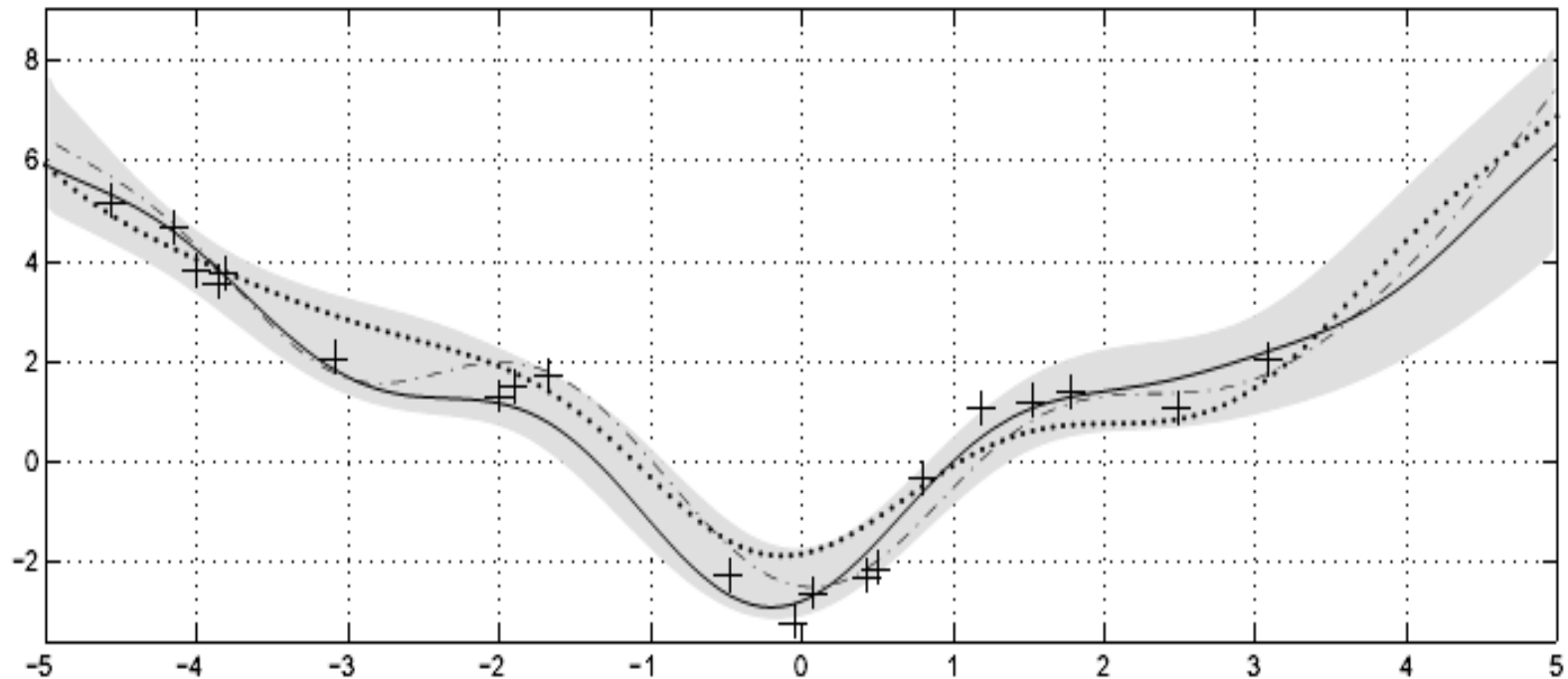
Procedure:

- The maximum likelihood estimate

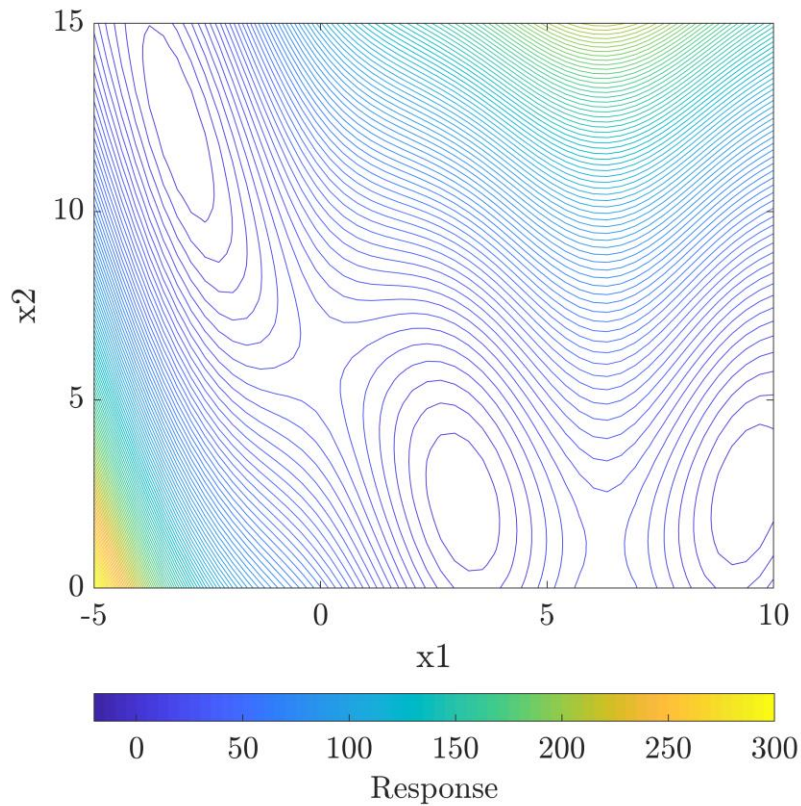
$$J(a) = \frac{1}{2} a^T K a + \frac{1}{2\sigma_n^2} (\mathbf{y} - K a)^2$$

- First term is the ridge penalty term to prevent overfitting
- Second term is the standard loss function
- Minimizing above function, $\frac{\partial J}{\partial a} = 0$, gives $a = (K + \sigma_n^2 I)^{-1} \mathbf{y}$

- 20 Training data
- GP posterior
- noise level $\sigma_n^2=0$ and $\sigma_n^2=7$



True Contour



GP Estimated

