

Basic Theory of Gaussian Processes

EN.540.782: Statistical Uncertainty Quantification

N. Wichrowski

Johns Hopkins University

March 16, 2021

Outline

- 1 Motivation
- 2 Definition and Sampling
- 3 Covariance Functions
- 4 Regression and Prediction
- 5 Conclusion and Appendix

Prelude: Linear Regression

- Consider data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$.
- We want the weight vector $\mathbf{w} \in \mathbb{R}^d$ that yields the best linear fit

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}.$$

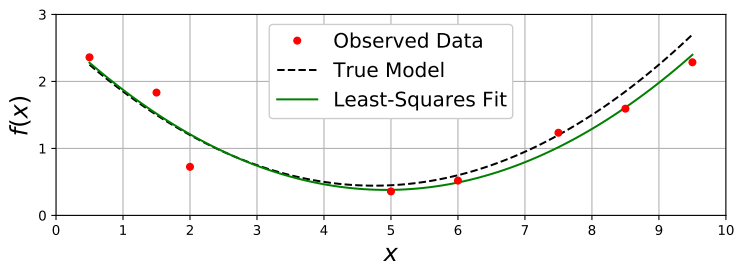
- Let $X \in \mathbb{R}^{n \times d}$ be the design matrix, $\mathbf{y} \in \mathbb{R}^n$ the response vector.
- Then $\hat{\mathbf{w}} = X^\dagger \mathbf{y} = (X^\top X)^{-1} X^\top \mathbf{y}$ minimizes the squared error:

$$X^\dagger \mathbf{y} \in \underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \|X\mathbf{w} - \mathbf{y}\|_2^2.$$

Linear Regression with Basis Functions

- Consider data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$.
- For weight vector $\mathbf{w} \in \mathbb{R}^N$ and basis functions $\phi_j : \mathbb{R}^d \rightarrow \mathbb{R}$,

$$f(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}) = \sum_{j=1}^N w_j \phi_j(\mathbf{x}).$$



Bayesian Linear Regression

- The standard linear model for “frequentist” regression is given by

$$\begin{aligned}f(\mathbf{x}) &= \mathbf{w}^\top \mathbf{x}, \\y_i &= f(\mathbf{x}_i) + \varepsilon_i, \\ \varepsilon_i &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_n^2).\end{aligned}$$

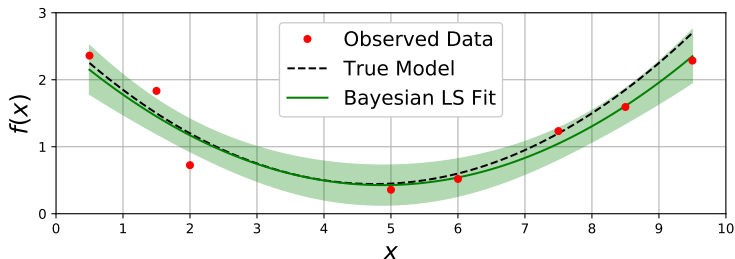
- In a Bayesian framework, we have $\mathbf{y} \mid X, \mathbf{w} \sim \mathcal{N}(X\mathbf{w}, \sigma_n^2 I)$.
- With prior $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p)$, the resulting posterior is

$$\begin{aligned}\mathbf{w} \mid X, \mathbf{y} &\sim \mathcal{N}(\bar{\mathbf{w}}, C), \\C^{-1} &= \sigma_n^{-2} X^\top X + \Sigma_p^{-1}, \\ \bar{\mathbf{w}} &= \sigma_n^{-2} C X^\top \mathbf{y}.\end{aligned}$$

Bayesian Linear Regression: Prediction

- We can use the posterior on \mathbf{w} to predict new observations.
- Given new input $\mathbf{x}_\star \in \mathbb{R}^d$,

$$y_\star | \mathbf{x}_\star, X, \mathbf{y} \sim \mathcal{N}(\bar{\mathbf{w}}^\top \mathbf{x}_\star, \mathbf{x}_\star^\top C \mathbf{x}_\star).$$



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.

- The index set \mathcal{X} is usually an interval $T \subseteq \mathbb{R}$.
- A GP is completely specified by:
 - ▶ mean function: $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$
 - ▶ covariance kernel: $k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$

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.

- The index set \mathcal{X} is usually an interval $T \subseteq \mathbb{R}$.
- A GP is completely specified by:
 - ▶ mean function: $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$
 - ▶ covariance kernel: $k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$

Example

Brownian motion on \mathbb{R} :

$$m(t) = 0$$

$$k(t, t') = \min(t, t')$$

Ornstein-Uhlenbeck Process:

$$m(t) = 0$$

$$k(t, t') = e^{-|t-t'|}$$

Function-Space View of GPs

- The mean and covariance functions of a GP define a distribution over real-valued functions defined on the index set \mathcal{X} .
- How do we sample $f \sim \mathcal{GP}[m(\cdot), k(\cdot, \cdot)]$?
 - ▶ Choose $\mathcal{D}_X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$.
 - ▶ Compute $m_i = m(\mathbf{x}_i)$, $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$.
 - ▶ Draw $\mathbf{y} \sim \mathcal{N}(\mathbf{m}, K)$ and take $f(\mathbf{x}_i) = y_i$.
- Use conditional distributions to sample $y_\star = f(\mathbf{x}_\star)$ for $\mathbf{x}_\star \notin \mathcal{D}_X$:

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

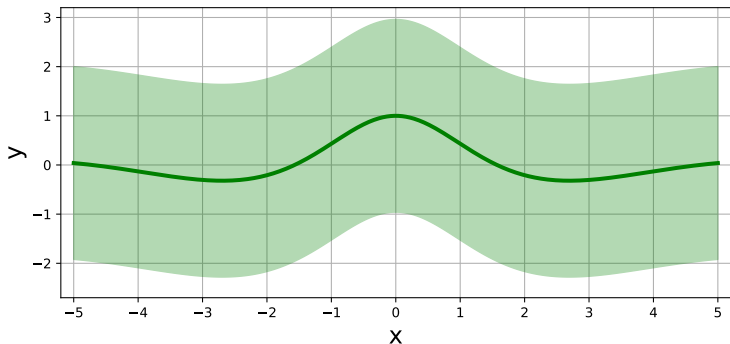
$$\text{Cov}[\mathbf{y}_\star | X_\star, X, \mathbf{y}] = K(X_\star, X_\star) - K(X_\star, X)K(X, X)^{-1}K(X, X_\star)$$

Sampling from a GP

Consider the following prior distribution on univariate functions:

$$m(x) = \frac{\cos(x)}{1 + 0.25x^2}$$

$$k(x, x') = \exp \left[-\frac{1}{2} (x - x')^2 \right]$$

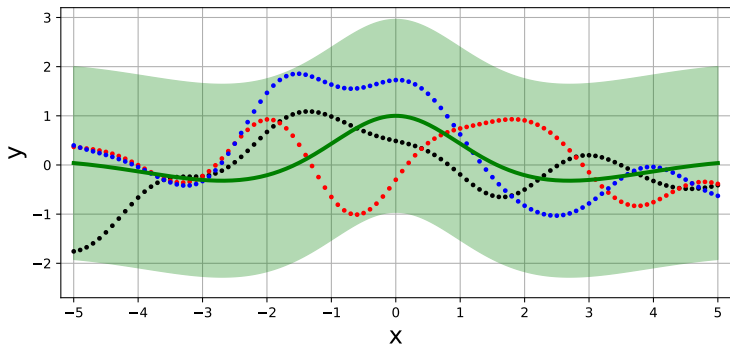


Sampling from a GP

Consider the following prior distribution on univariate functions:

$$m(x) = \frac{\cos(x)}{1 + 0.25x^2}$$

$$k(x, x') = \exp \left[-\frac{1}{2} (x - x')^2 \right]$$

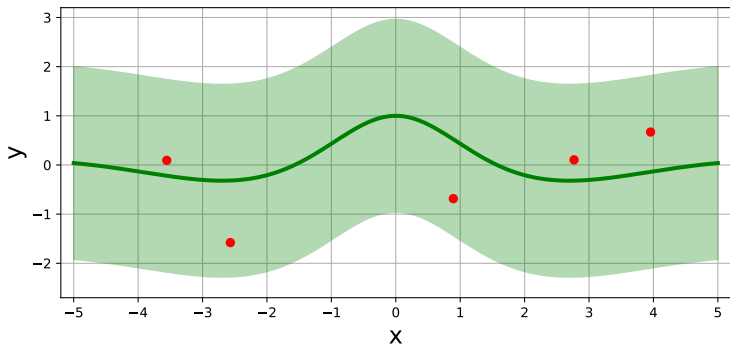


Sampling from a GP

Consider the following prior distribution on univariate functions:

$$m(x) = \frac{\cos(x)}{1 + 0.25x^2}$$

$$k(x, x') = \exp \left[-\frac{1}{2} (x - x')^2 \right]$$

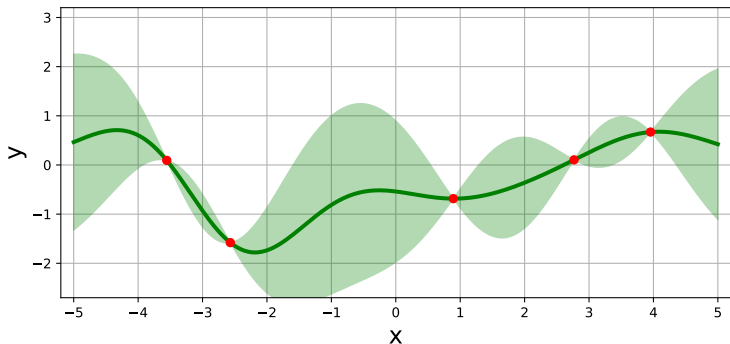


Sampling from a GP

Consider the following prior distribution on univariate functions:

$$m(x) = \frac{\cos(x)}{1 + 0.25x^2}$$

$$k(x, x') = \exp\left[-\frac{1}{2}(x - x')^2\right]$$

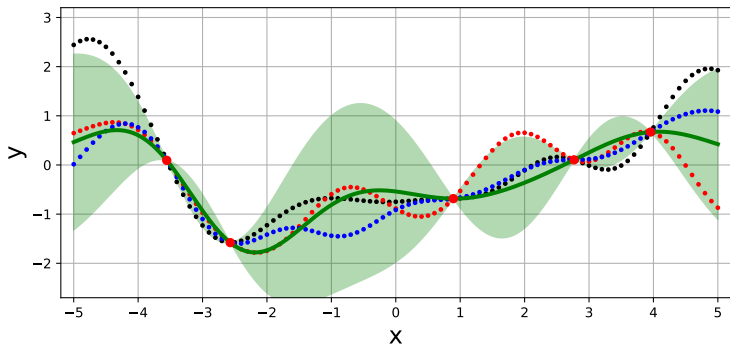


Sampling from a GP

Consider the following prior distribution on univariate functions:

$$m(x) = \frac{\cos(x)}{1 + 0.25x^2}$$

$$k(x, x') = \exp \left[-\frac{1}{2} (x - x')^2 \right]$$



Covariance Functions: Concepts

- A kernel must be symmetric and positive semidefinite.

- ▶ For all $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$.
- ▶ Given measure μ , for all $f \in L^2(\mathcal{X}, \mu)$,

$$\iint_{\mathcal{X} \times \mathcal{X}} k(\mathbf{x}, \mathbf{x}') f(\mathbf{x}) f(\mathbf{x}') d\mu(\mathbf{x}) d\mu(\mathbf{x}') \geq 0.$$

- We say k is stationary if $k(\mathbf{x}, \mathbf{x}') = g(\mathbf{x} - \mathbf{x}')$.
- We say k is isotropic if $k(\mathbf{x}, \mathbf{x}') = g(\|\mathbf{x} - \mathbf{x}'\|)$.
- Kernels can be added, multiplied, and scaled:

$$k(\mathbf{x}, \mathbf{x}') = c_1^2 k_1(\mathbf{x}, \mathbf{x}') + c_2^2 k_2(\mathbf{x}, \mathbf{x}') k_3(\mathbf{x}, \mathbf{x}').$$

Continuity and Differentiability

Definition

Let $f \sim \mathcal{GP}[m(\cdot), k(\cdot, \cdot)]$ be a Gaussian process on $\mathcal{X} \subseteq \mathbb{R}^d$. Then f is continuous in mean square (CMS) at $\mathbf{x}_\star \in \mathcal{X}$ if

$$\lim_{\mathbf{x} \rightarrow \mathbf{x}_\star} \mathbb{E}[|f(\mathbf{x}) - f(\mathbf{x}_\star)|^2] = 0.$$

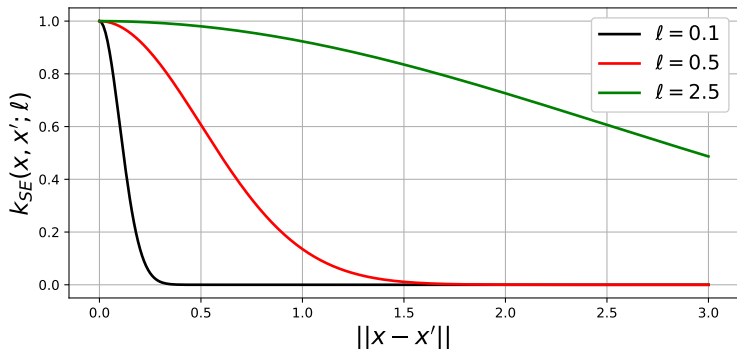
We say that f is mean-square differentiable (MSD) at \mathbf{x}_\star with partial derivatives $\partial f(\mathbf{x}_\star)/\partial x_i$ if, for $i \in \{1, \dots, d\}$,

$$\lim_{h \rightarrow 0} \mathbb{E} \left[\left(\frac{f(\mathbf{x}_\star + h\mathbf{e}_i) - f(\mathbf{x}_\star)}{h} - \frac{\partial f(\mathbf{x}_\star)}{\partial x_i} \right)^2 \right] = 0.$$

- GP with kernel k is CMS at $\mathbf{x}_\star \in \mathcal{X}$ iff k is continuous at $(\mathbf{x}_\star, \mathbf{x}_\star)$.
- A $2p$ -order derivative of $k(\mathbf{x}_\star, \mathbf{x}_\star)$ ensures f is MSD p times at \mathbf{x}_\star .

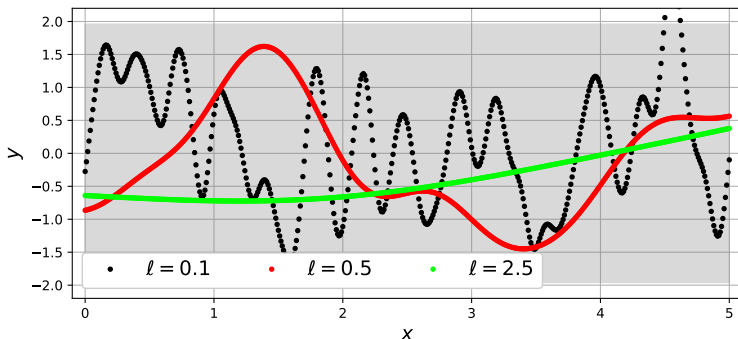
Covariance Functions: Squared Exponential

$$k(\mathbf{x}, \mathbf{x}') = \exp \left[-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2} \right] \quad \ell > 0$$



Covariance Functions: Squared Exponential

$$k(\mathbf{x}, \mathbf{x}') = \exp \left[-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2} \right] \quad \ell > 0$$

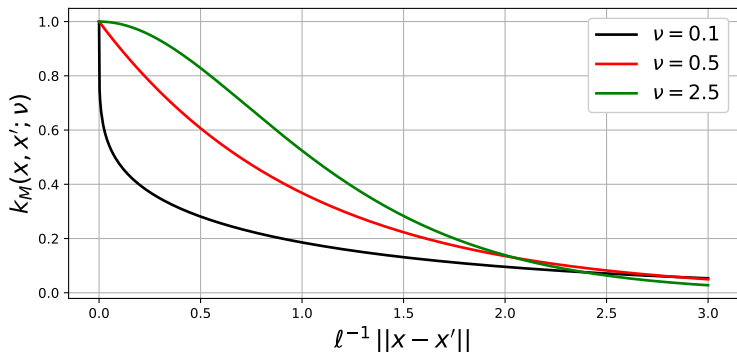


Covariance Functions: Matérn

$$k(\mathbf{x}, \mathbf{x}') = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{r\sqrt{2\nu}}{\ell} \right)^\nu K_\nu \left(\frac{r\sqrt{2\nu}}{\ell} \right) \quad \ell, \nu > 0$$

$K_\nu(\cdot)$ is a modified Bessel function

Can be expressed more simply if $\nu + \frac{1}{2} \in \mathbb{N}$ (see [Appendix](#)).

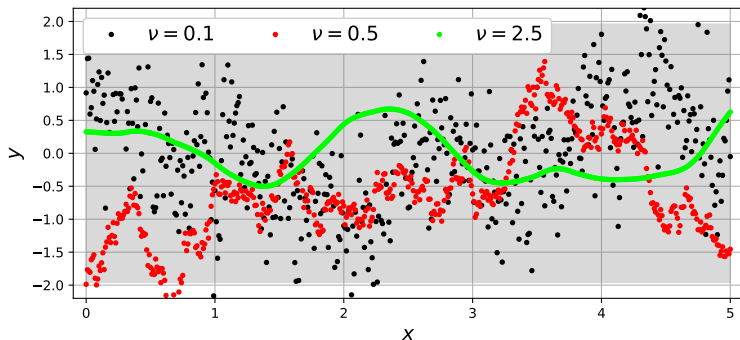


Covariance Functions: Matérn

$$k(\mathbf{x}, \mathbf{x}') = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{r\sqrt{2\nu}}{\ell} \right)^\nu K_\nu \left(\frac{r\sqrt{2\nu}}{\ell} \right) \quad \ell, \nu > 0$$

$K_\nu(\cdot)$ is a modified Bessel function

Can be expressed more simply if $\nu + \frac{1}{2} \in \mathbb{N}$ (see [Appendix](#)).



Fitting GPs to Data

- Consider data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$.
- Recall the prediction for $y_\star = f(\mathbf{x}_\star)$ at $\mathbf{x}_\star \notin \mathcal{D}_X$:

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

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

- If noisy measurements, then replace $K(X, X)$ by $K(X, X) + \sigma_n^2 I$.
- Given \mathcal{D} , the GPR model is specified by our choice of kernel.

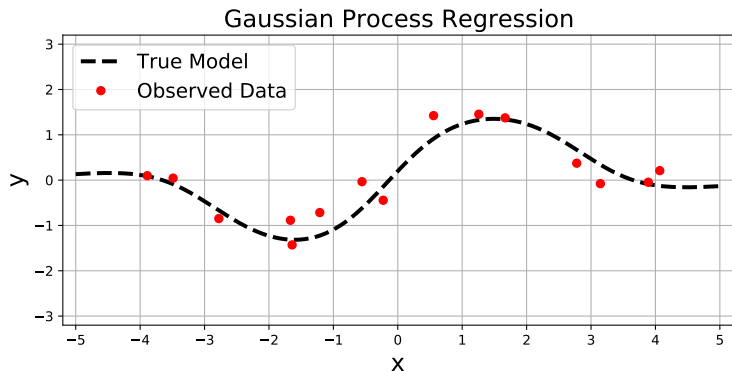
Model Selection

- GP Regression requires:
 - ▶ Selecting a kernel function (model).
 - ▶ Tuning the hyperparameters.
- “Minimize the generalization error.”
- Log marginal likelihood:

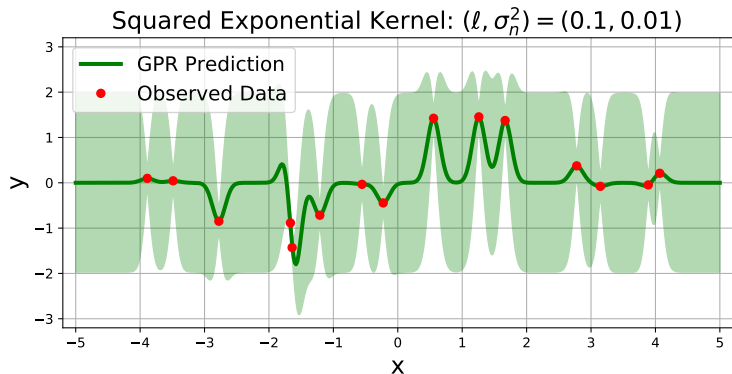
$$\log p(\mathbf{y} | X, \boldsymbol{\theta}) = -\frac{1}{2} \mathbf{y}^\top K_y \mathbf{y} - \frac{1}{2} \log |K_y| - \frac{n}{2} \log(2\pi)$$
$$K_y = K(X, X; \boldsymbol{\theta}) + \sigma_n^2 I$$

- Cross-validation

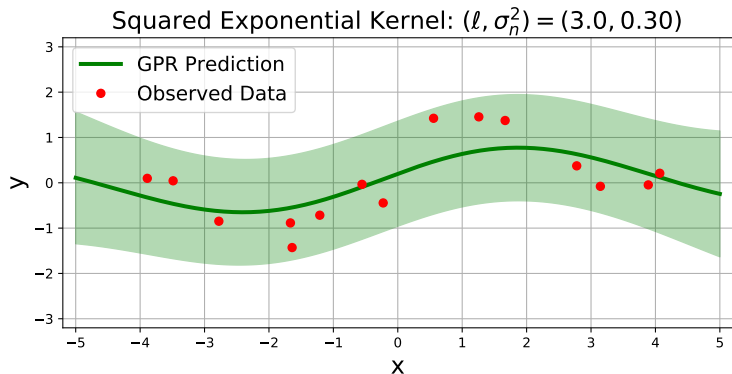
Hyperparameter Tuning: An Illustration



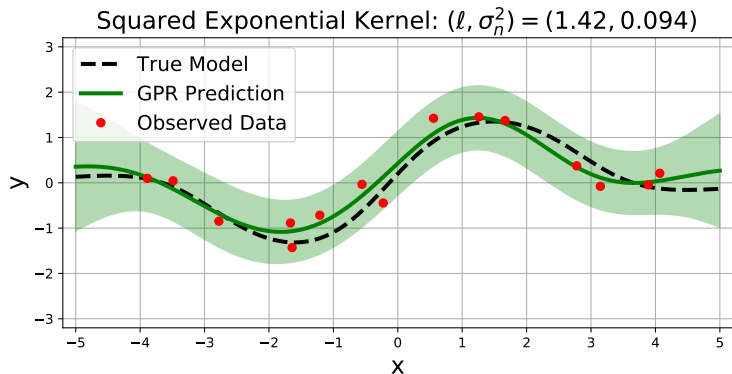
Hyperparameter Tuning: An Illustration



Hyperparameter Tuning: An Illustration



Hyperparameter Tuning: An Illustration



Final Thoughts

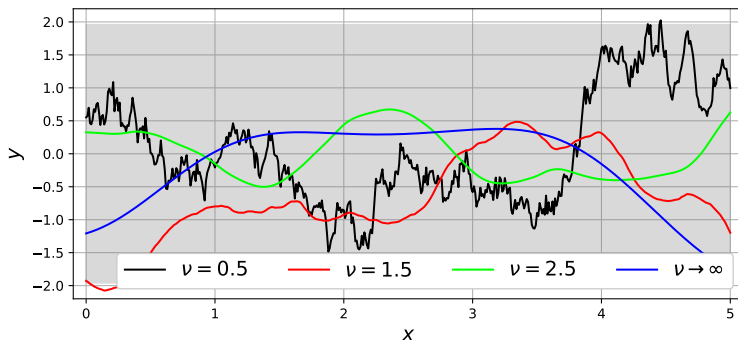
- GPs offer benefits for both regression and fitting.
- Precise functional forms are not needed.
- Bayesian framework provides a natural measure of uncertainty.

References

- Duvenaud, D.K. *Automatic Model Construction with Gaussian Processes* [Ph.D. dissertation], Pembroke College: Cambridge, England (2014).
- Goldberg, P.W.; Williams, C.K.I.; Bishop, C.M. “Regression with Input-Dependent Noise: A Gaussian Process Treatment,” *Adv Neural Inf Process Syst*, 10, 493–499 (1997).
- Jun, M. *Spatial Statistics* [[Lecture Slides](#)], Texas A&M University: College Station, TX (2009).
- Lifshits, M. *Lectures on Gaussian Processes*, Springer: New York, NY (2012).
- Rasmussen, C.E.; Williams, C.K.I. *Gaussian Processes for Machine Learning*, MIT Press: Cambridge, MA (2006).

Matérn Kernel with $\nu \in \{\frac{1}{2}, \frac{3}{2}, \frac{5}{2}\}$

$$k_{1/2}(r) = \exp\left[-\frac{r}{\ell}\right] \quad k_{3/2}(r) = \left(1 + \frac{r\sqrt{3}}{\ell}\right) \exp\left[-\frac{r\sqrt{3}}{\ell}\right]$$
$$k_{5/2}(r) = \left(1 + \frac{r\sqrt{5}}{\ell} + \frac{5r^2}{3\ell^2}\right) \exp\left[-\frac{r\sqrt{5}}{\ell}\right]$$

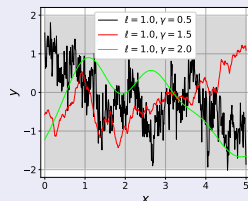


Additional (Isotropic) Families of Covariance Kernels

Gamma-Exponential

$$k(r) = \exp \left[- \left(\frac{r}{\ell} \right)^\gamma \right] \quad 0 < \gamma \leq 2, \ell > 0$$

- Reduces to Squared Exponential when $\gamma = 2$.
- Not mean-square differentiable for any $\gamma < 2$.



Rational Quadratic

$$k(r) = \left(1 + \frac{r^2}{2\alpha\ell^2} \right)^{-\alpha} \quad \alpha, \ell > 0$$

- Reduces to Squared Exponential as $\alpha \rightarrow \infty$.
- Mean-square differentiable for all $\alpha > 0$.

