

Machine Learning

Gaussian Processes

Wojciech Czech

January 9, 2024

Institute of Computer Science

Table of contents

1. Stochastic processes
2. Gaussian processes
3. Gaussian processes for regression
4. Gaussian process kernels

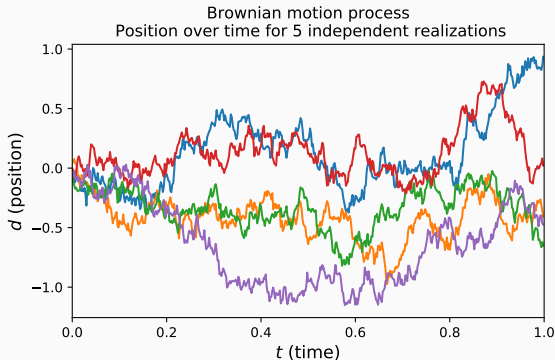
Stochastic processes

Brownian motion

- Random motion of particles suspended in the fluid
- Random walk of particles due to other particles randomly bumping into them

$$d(t + \Delta t) = d(t) + \Delta d \quad (1)$$

$$\Delta d \sim \mathcal{N}(0, \Delta t) \quad (2)$$



Stochastic processes - distributions over functions

- Stochastic processes describes systems randomly changing over time
- Trajectories (paths, realizations) of the stochastic process are different for the same starting point
- Every realization corresponds to a function
- Stochastic process can be understood as a distribution over functions

Gaussian processes

Gaussian processes

Gaussian process

Gaussian process is a Gaussian distribution over functions $f(\mathbf{x})$, defined by mean function $m(\mathbf{x})$ and positive definite covariance function (**kernel function**) $k(\mathbf{x}, \mathbf{x}')$,

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (3)$$

such that the set of values of $f(\mathbf{x})$ evaluated at an arbitrary set of points $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ jointly have Gaussian distribution:

$$\mathbf{y} = f(\mathbf{X}) \sim \mathcal{N}(m(\mathbf{X}), k(\mathbf{X}, \mathbf{X})) \quad (4)$$

with mean vector $m(\mathbf{X})$ and covariance matrix $k(\mathbf{X}, \mathbf{X})$.

Positive definite matrix

Positive definite matrix

A symmetric matrix \mathbf{A} whose eigenvalues are positive is called *positive definite* ($\mathbf{z}^T \mathbf{A} \mathbf{z} > 0$ for $\mathbf{z} \in \mathbb{R}^{n \times 1}$ and $\mathbf{z} \neq \mathbf{0}$).

Kernel function as prior

- Sampling functions from Gaussian process requires defining mean and covariance (kernel) functions
- Covariance function models the joint variability of Gaussian process random variables and returns modeled covariance for each pair of inputs $k(\mathbf{x}, \mathbf{x}')$
- By choosing a specific kernel function k it is possible to set prior information on the distribution of $f(\mathbf{x})$

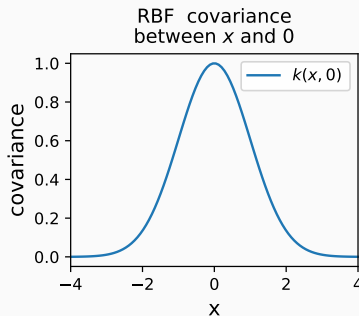
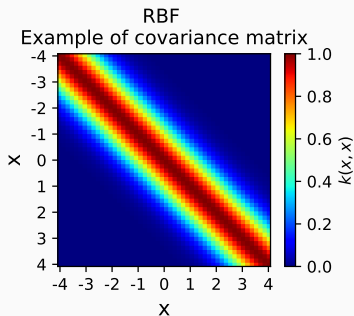
1. Sample function evaluations \mathbf{y} from function $f(\mathbf{x})$ drawn from Gaussian process at arbitrary set of points \mathbf{X}

$$\mathbf{y} = f(\mathbf{X}) \tag{5}$$

2. Finite dimensional subset of the Gaussian process distribution results in Gaussian distribution $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\mu} = m(\mathbf{X})$ and $\boldsymbol{\Sigma} = k(\mathbf{X}, \mathbf{X})$

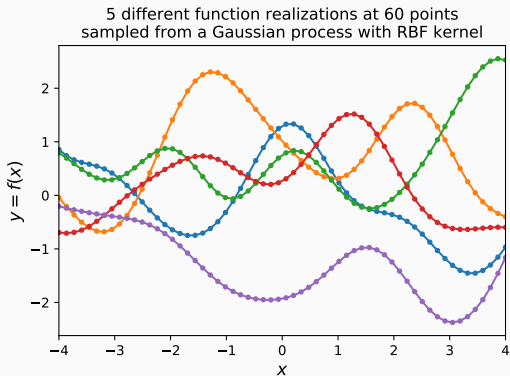
Example 1: RBF kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}'\|^2\right) \quad (6)$$

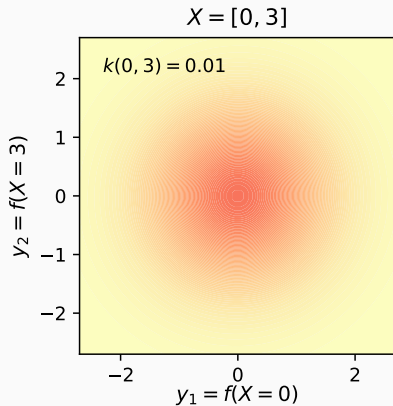
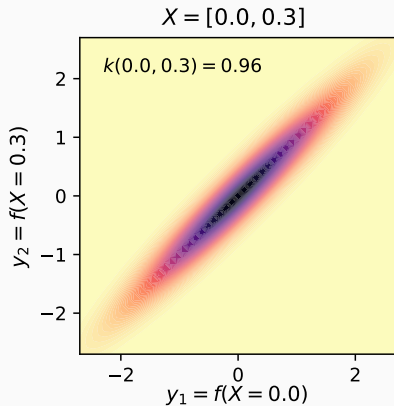


Example 1: Sampling from prior

1. Generate 60 equidistant samples from $[-4, 4]$: $\mathbf{X} = \{x_1, \dots, x_{60}\}$
2. Calculate covariance matrix $\mathbf{K} = k(\mathbf{X}, \mathbf{X})$
3. Generate 5 samples (realizations) from 60-dimensional Gaussian distribution with covariance matrix $\mathbf{\Sigma}$ and zero mean $\mathbf{0}$

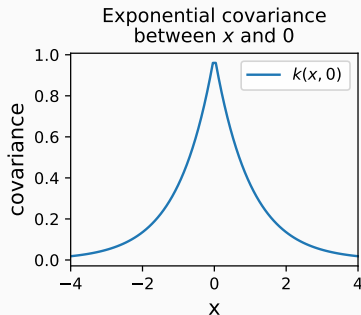
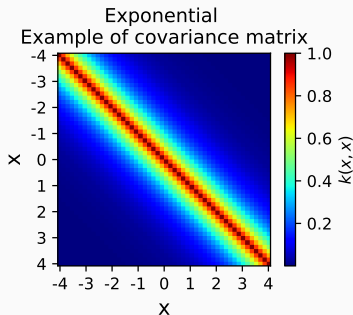


Example 1: 2D marginal distributions



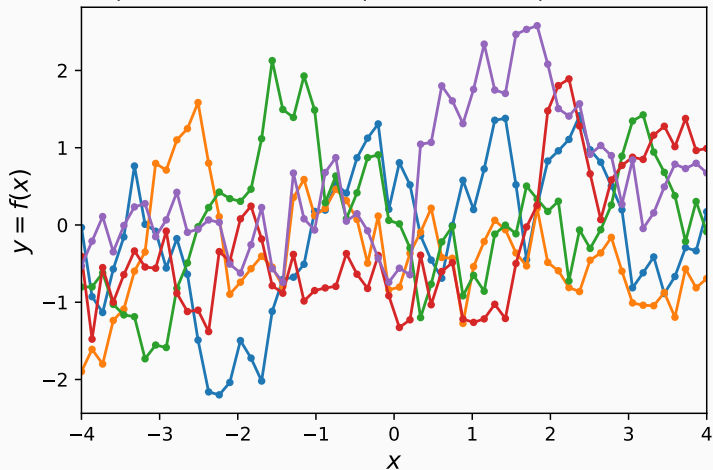
Example 2: Exponential kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\theta \|\mathbf{x} - \mathbf{x}'\|) \quad (7)$$

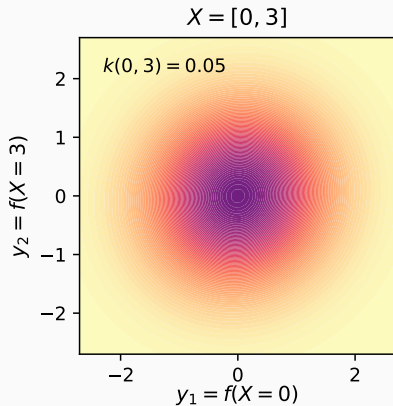
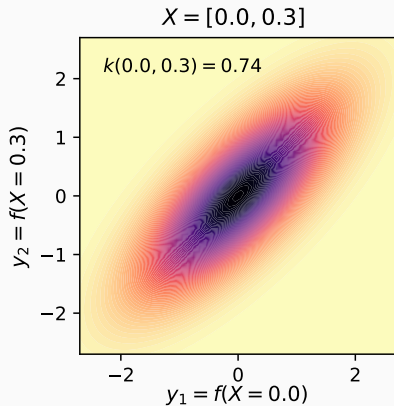


Example 2: Sampling from prior

5 different function realizations at 60 points
sampled from a Gaussian process with exponential kernel



Example 2: 2D marginal distributions



Gaussian processes for regression

Linear regression revisited

- Linear combination of fixed basis functions

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) \quad (8)$$

- Prior distribution over \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I}) \quad (9)$$

- Probability distribution over \mathbf{w} induces probability distribution over functions $y(\mathbf{w})$
- Joint probability of the function values: $\mathbf{y} = [y(\mathbf{x}_1), \dots, y(\mathbf{x}_n)]^T$

$$\mathbf{y} = \Phi \mathbf{w} \quad (10)$$

where $[\Phi]_{i,j} = \phi_j(\mathbf{x}_i)$ - design matrix

Joint distribution over y and Gaussian process

$$\mathbb{E}[\mathbf{y}] = \Phi \mathbb{E}[\mathbf{w}] = \mathbf{0} \quad (11)$$

$$\text{cov}[\mathbf{y}] = \mathbb{E}[\mathbf{y}\mathbf{y}^T] = \Phi \mathbb{E}[\mathbf{w}\mathbf{w}^T] \Phi^T = \frac{1}{\alpha} \Phi \Phi^T = \mathbf{K} \quad (12)$$

$$[\mathbf{K}]_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{\alpha} \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \quad (13)$$

\mathbf{K} - Gram matrix

$k(\mathbf{x}_i, \mathbf{x}_j)$ - kernel function vs. selecting basis

Prior distribution over \mathbf{w} vs. prior distribution over functions $y(\mathbf{x})$

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{K}) \quad (14)$$

Gaussian process for regression

1. Use Gaussian process defined by kernel function as prior
2. Generate posterior distribution given training data
3. Use posterior distribution to predict the expected value and probability of the output variable y given input variables

Regression problem with noise

- Training data
 - a) Target values: $\mathbf{t}_n = [t_1, \dots, t_n]^T$
 - b) Input values: $\mathbf{X}_n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$
- Noise on observed target values

$$t_i = y_i + \epsilon_i \quad (15)$$

where $y_i = f(\mathbf{x}_i)$, $\mathbf{y}_n = [y_1, \dots, y_n]^T$

- Noise has Gaussian distribution

$$p(t_i|y_i) = \mathcal{N}(t_i|y_i, \beta^{-1}) \quad (16)$$

- Joint distribution on data points (see [1])

$$p(\mathbf{t}_n|\mathbf{y}_n) = \mathcal{N}(\mathbf{t}_n|\mathbf{y}_n, \beta^{-1}\mathbf{I}_n) \quad (17)$$

- Model of the joint distribution over sets of data points

$$p(\mathbf{t}) = \int p(\mathbf{t}|\mathbf{y})p(\mathbf{y})d\mathbf{y} = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C}) \quad (18)$$

$$[\mathbf{C}]_{i,j} = C(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j) + \beta^{-1}\delta_{ij} \quad (19)$$

- Prediction problem: provide target value t_{n+1} for a new input vector \mathbf{x}_{n+1}
- Find predictive distribution

$$p(t_{n+1}|\mathbf{t}_n, \mathbf{X}_n, \mathbf{x}_{n+1}) = p(t_{n+1}|\mathbf{t}_n, \mathbf{x}_1, \dots, \mathbf{x}_{n+1})$$

Step 1: Joint distribution $p(\mathbf{t}_{n+1})$

$$p(\mathbf{t}_{n+1}) = p(t_1, \dots, t_{n+1}) = \mathcal{N}(\mathbf{t}_{n+1} | \mathbf{0}, \mathbf{C}_{n+1}) \quad (20)$$

$\mathbf{C}_{n+1} \in \mathbb{R}^{(n+1) \times (n+1)}$ - covariance matrix

$$C(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j) + \beta^{-1} \delta_{ij} \quad (21)$$

$$\mathbf{C}_{n+1} = \begin{bmatrix} \mathbf{C}_n & \mathbf{k} \\ \mathbf{k}^T & c \end{bmatrix} \quad (22)$$

$$c = k(\mathbf{x}_{n+1}, \mathbf{x}_{n+1}) + \beta^{-1} \quad (23)$$

$$\mathbf{k} = [k(\mathbf{x}_1, \mathbf{x}_{n+1}), \dots, k(\mathbf{x}_n, \mathbf{x}_{n+1})]^T \quad (24)$$

Step 2: Conditional distribution $p(\mathbf{t}_{n+1}|\mathbf{t}_n)$

$$p(t_{n+1}|\mathbf{t}_n, \mathbf{X}_n, \mathbf{x}_{n+1}) = \mathcal{N}(t_{n+1}|m(\mathbf{x}_{n+1}), \sigma^2(\mathbf{x}_{n+1})) \quad (25)$$

Partitioned Gaussian distributions

Two sets of variables

If two sets of variables are jointly Gaussian, then conditional distribution of one set conditioned on the other is also Gaussian. Marginal distribution of either set is also Gaussian.

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}, \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix} \quad (26)$$

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{bmatrix}, \boldsymbol{\Sigma} = \boldsymbol{\Sigma}^T, \boldsymbol{\Sigma}_{ab} = \boldsymbol{\Sigma}_{ba}^T \quad (27)$$

Conditional distribution

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b}) \quad (28)$$

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b) \quad (29)$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba} \quad (30)$$

$$p(t_{n+1}|\mathbf{t}_n, \mathbf{X}_n, \mathbf{x}_{n+1}) = \mathcal{N}(t_{n+1}|m(\mathbf{x}_{n+1}), \sigma^2(\mathbf{x}_{n+1})) \quad (31)$$

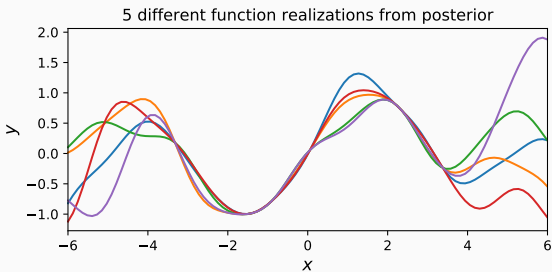
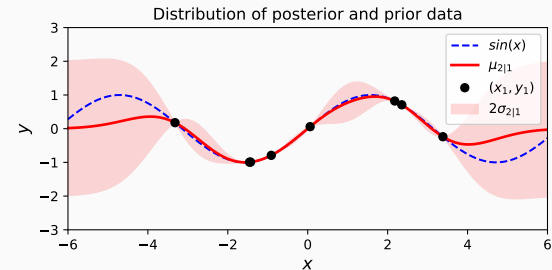
$$m(\mathbf{x}_{n+1}) = \mathbf{k}^T \mathbf{C}_n^{-1} \mathbf{t}_n \quad (32)$$

$$\sigma^2(\mathbf{x}_{n+1}) = c - \mathbf{k}^T \mathbf{C}_n^{-1} \mathbf{k}_n \quad (33)$$

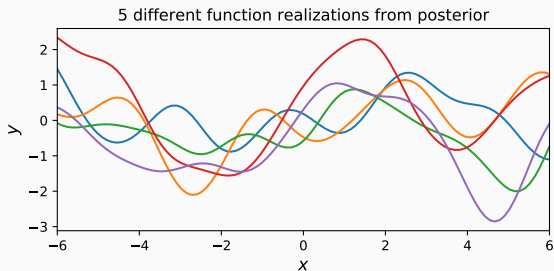
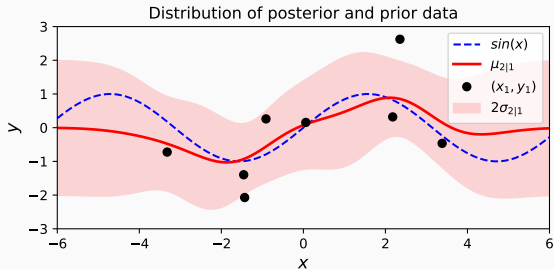
Gaussian process regression - algorithm

1. Collect training points
2. Define testing points in posterior (e.g. uniformly spaced to capture function)
3. Compute posterior mean and variance using Equation 32 and Equation 33
4. Compute standard deviations at tests points
5. Plot posterior distribution and samples from posterior (see [2])

Example 3: Gaussian process regression - noiseless distribution



Example 4: Gaussian process regression - noisy observations



Gaussian process kernels

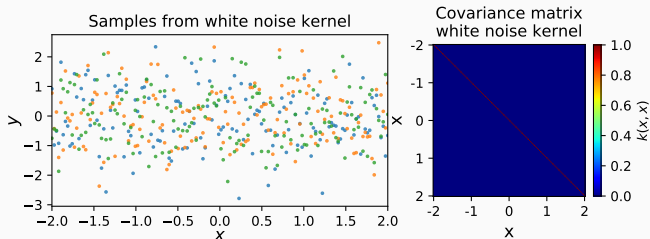
$$\Sigma = k(\mathbf{X}, \mathbf{X}) \quad (34)$$

- Positive definite $\mathbf{x}^T \Sigma \mathbf{x} > 0$ (\mathbf{x} is non-zero vector)
- Σ is symmetric
- Σ is invertible

White noise kernel

$$k(\mathbf{x}_a, \mathbf{x}_b) = \beta^{-1} \mathbf{I}_n \quad (35)$$

- β - precision of the noise, $\beta^{-1} = \sigma^2$
- \mathbf{I}_n - identity matrix
- Diagonal covariance matrix (noise is uncorrelated)



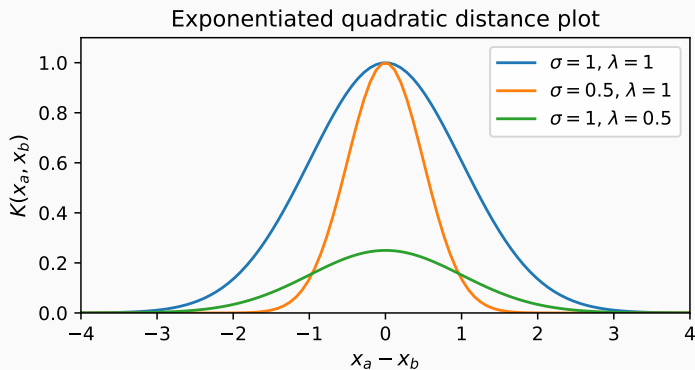
Exponential quadratic kernel

$$k(\mathbf{x}_a, \mathbf{x}_b) = \lambda^2 \exp \left(-\frac{\|\mathbf{x}_a - \mathbf{x}_b\|^2}{2\sigma^2} \right) \quad (36)$$

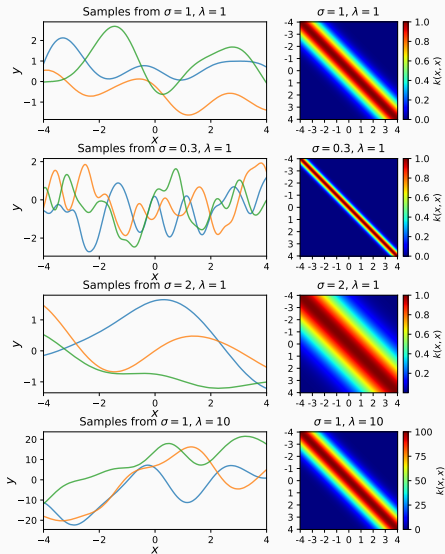
- λ^2 - overall variance (affects maximum value of the covariance)
- σ - length scale (the spread of the covariance)
- Smooth prior on functions sampled from Gaussian process

Exponential quadratic kernel - distance

$$x_a = 0$$



Exponential quadratic kernel - samples



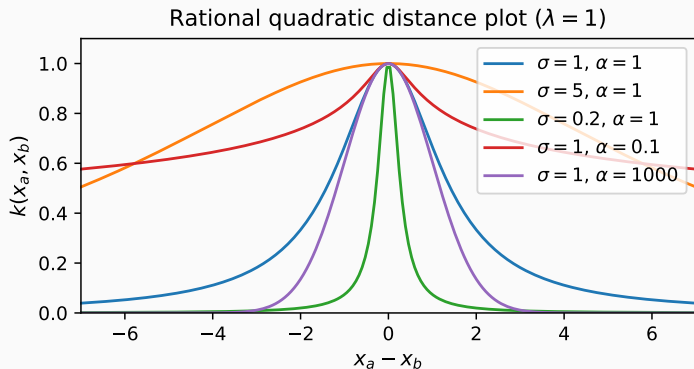
Rational quadratic kernel

$$k(\mathbf{x}_a, \mathbf{x}_b) = \lambda^2 \left(1 + \frac{\|\mathbf{x}_a - \mathbf{x}_b\|^2}{2\alpha\sigma^2} \right)^{-\alpha} \quad (37)$$

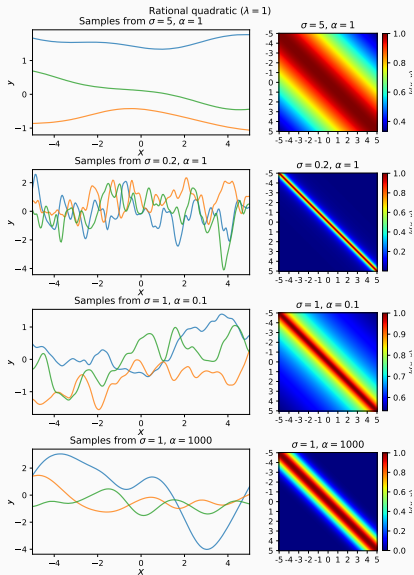
- λ^2 - overall variance, σ - length scale
- α - scale mixture ($\alpha > 0$)
- Infinite sum of different exponentiated quadratic kernels with different length scales and with α determining the weighting between different length scales

Rational quadratic kernel - distance

$$x_a = 0$$



Rational quadratic kernel - samples

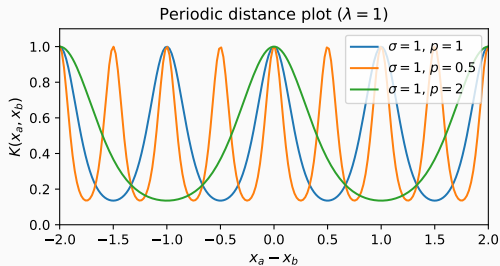
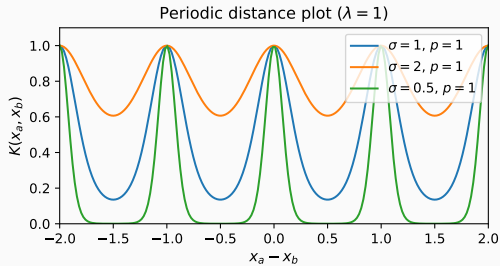


Periodic kernel

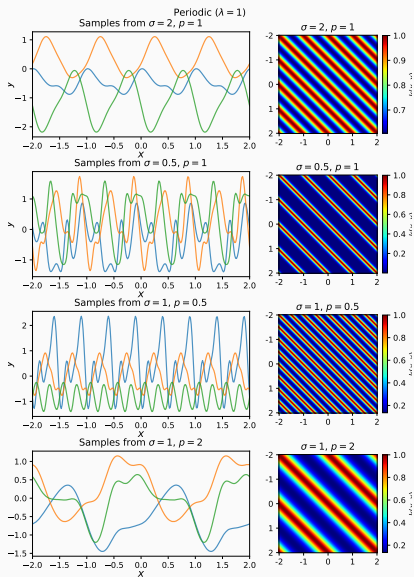
$$k(\mathbf{x}_a, \mathbf{x}_b) = \lambda^2 \exp \left(-\frac{2}{\sigma^2} \sin^2 \left(\pi \frac{\|\mathbf{x}_a - \mathbf{x}_b\|}{p} \right) \right) \quad (38)$$

- λ^2 - overall variance (λ is amplitude), σ - the length scale
- p - the period (distance between repetitions)
- Allows for modeling periodic functions

Periodic kernel - distance



Periodic kernel - samples



Combining kernels by multiplication

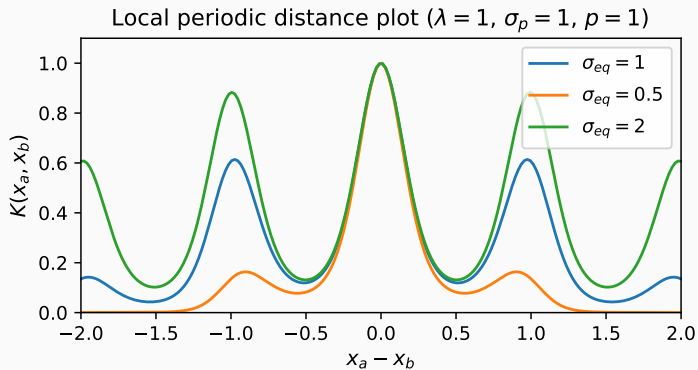
- Multiplying kernels is an element wise multiplication of their corresponding covariance matrices (AND operator)
- Allows for modeling functions that are only locally periodic - the shape of the repeating part of the function can change over time

Local periodic kernel

$$k(\mathbf{x}_a, \mathbf{x}_b) = \lambda^2 \exp \left(-\frac{2}{\sigma_p^2} \sin^2 \left(\pi \frac{\|\mathbf{x}_a - \mathbf{x}_b\|}{p} \right) \right) \exp \left(-\frac{\|\mathbf{x}_a - \mathbf{x}_b\|^2}{2\sigma_{eq}^2} \right) \quad (39)$$

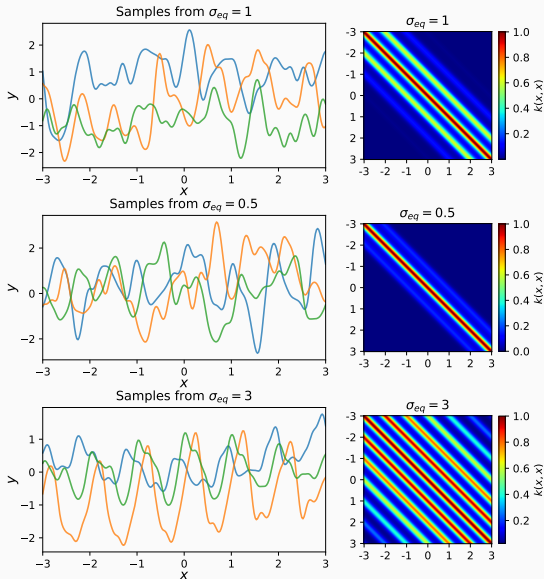
- λ^2 - overall variance (λ is amplitude), σ_p - the length scale of the periodic function, σ_{eq} - the length scale of the exponentiated quadratic function
- p - the period (distance between repetitions)
- Allows the periods to vary over longer distances

Local periodic kernel - distance



Local periodic kernel - samples

Local periodic ($\lambda = 1$, $\sigma_p = 1$, $\rho = 1$)



Combining kernels by addition

- Kernels can be combined by adding them together
- Adding kernels is an element wise addition of their corresponding covariance matrices (OR operator)

Fitting Gaussian process kernel

Create kernel combined by addition to reflect different characteristics of the data:

- Long term smooth change
- Seasonality
- Short term irregularities
- Medium term irregularities
- White noise



C. M. Bishop.

Pattern Recognition and Machine Learning.

Springer, 2006.



P. Roelants.

Understanding Gaussian processes, Accessed November, 2020.

Available at <https://peterroelants.github.io/posts/gaussian-process-tutorial/>.