Introduction to Machine Learning

Gaussian Processes Covariance functions for GPs

× 0 0 × 0 × × ×



Learning goals

- Covariance functions encode key assumptions about the GP
- Know common covariance functions like squared exponential and Matérn

COVARIANCE FUNCTION OF A GP

The marginalization property of the Gaussian process implies that for any finite set of input values, the corresponding vector of function values is Gaussian:

$$\boldsymbol{f} = \left[f\left(\boldsymbol{x}^{(1)} \right), ..., f\left(\boldsymbol{x}^{(n)} \right) \right] \sim \mathcal{N}\left(\boldsymbol{m}, \boldsymbol{K} \right),$$

- The covariance matrix *K* is constructed based on the chosen inputs {x⁽¹⁾, ..., x⁽ⁿ⁾}.
- Entry \boldsymbol{K}_{ij} is computed by $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$.
- Technically, for every choice of inputs {x⁽¹⁾, ..., x⁽ⁿ⁾}, K needs to be positive semi-definite in order to be a valid covariance matrix.
- A function k(.,.) satisfying this property is called **positive definite**.

× × ×

COVARIANCE FUNCTION OF A GP / 2

• Recall, the purpose of the covariance function is to control to which degree the following is fulfilled:

If two points $\mathbf{x}^{(i)}, \mathbf{x}^{(j)}$ are close in \mathcal{X} -space, their function values $f(\mathbf{x}^{(i)}), f(\mathbf{x}^{(j)})$ should be close (**correlated**!) in \mathcal{Y} -space.

Closeness of two points x⁽ⁱ⁾, x^(j) in input space X is measured in terms of d = x⁽ⁱ⁾ - x^(j):

$$k(\mathbf{x}^{(i)},\mathbf{x}^{(j)})=k(\mathbf{d})$$



COVARIANCE FUNCTION OF A GP: EXAMPLE

- Let $f(\mathbf{x})$ be a GP with $k(\mathbf{x}, \mathbf{x}') = \exp(-\frac{1}{2} \|\mathbf{d}\|^2)$ with $\mathbf{d} = \mathbf{x} \mathbf{x}'$.
- Consider two points $\mathbf{x}^{(1)} = 3$ and $\mathbf{x}^{(2)} = 2.5$.
- If you want to know how correlated their function values are, compute their correlation!



× × × ×

COVARIANCE FUNCTION OF A GP: EXAMPLE

• Assume we observed a value $y^{(1)} = -0.8$, the value of $y^{(2)}$ should be close under the assumption of the above Gaussian process.



× 0 0 × 0 × ×

COVARIANCE FUNCTION OF A GP: EXAMPLE

- Let us compare another point **x**⁽³⁾ to the point **x**⁽¹⁾
- We again compute their correlation
- Their function values are not very much correlated; y⁽¹⁾ and y⁽³⁾ might be far away from each other





COVARIANCE FUNCTIONS

There are three types of commonly used covariance functions:

k(.,.) is called stationary if it is as a function of *d* = *x* − *x*[′], we write *k*(*d*).

Stationarity is invariance to translations in the input space:

- $k(\pmb{x}, \pmb{x} + \pmb{d}) = k(\pmb{0}, \pmb{d})$
- k(.,.) is called isotropic if it is a function of r = ||x x'||, we write k(r).
 Isotropy is invariance to rotations of the input space and implies

stationarity.

• k(.,.) is a dot product covariance function if k is a function of $\mathbf{x}^T \mathbf{x}'$

× 0 0 × × ×

COMMONLY USED COVARIANCE FUNCTIONS

Name	$k(\mathbf{x}, \mathbf{x}')$
constant	σ_0^2
linear	$\sigma_0^2 + \pmb{x}^{\intercal} \pmb{x}'$
polynomial	$(\sigma_0^2 + \boldsymbol{x}^T \boldsymbol{x}')^p$
squared exponential	$\exp(-\frac{\ \boldsymbol{x}-\boldsymbol{x}'\ ^2}{2\ell^2})$
Matérn	$\frac{1}{2^{\nu}\Gamma(\nu)}\left(\frac{\sqrt{2\nu}}{\ell}\ \boldsymbol{x}-\boldsymbol{x}'\ \right)^{\nu}\mathcal{K}_{\nu}\left(\frac{\sqrt{2\nu}}{\ell}\ \boldsymbol{x}-\boldsymbol{x}'\ \right)$
exponential	$\exp\left(-\frac{\ \mathbf{x}-\mathbf{x}'\ }{\ell}\right)$

× × 0 × × ×

 $K_{\nu}(\cdot)$ is the modified Bessel function of the second kind.

COMMONLY USED COVARIANCE FUNCTIONS / 2



× 0 0 × 0 × ×

- Random functions drawn from Gaussian processes with a Squared Exponential Kernel (left), Polynomial Kernel (middle), and a Matérn Kernel (right, $\ell = 1$).
- The length-scale hyperparameter determines the "wiggliness" of the function.
- For Matérn, the ν parameter determines how differentiable the process is.

SQUARED EXPONENTIAL COVARIANCE FUNCTION

The squared exponential function is one of the most commonly used covariance functions.

$$k(\mathbf{x},\mathbf{x}') = \exp\left(-\frac{\|\mathbf{x}-\mathbf{x}'\|^2}{2\ell^2}\right).$$

× 0 0 × 0 × ×

Properties:

- It depends merely on the distance *r* = ||**x** − **x**'|| → isotropic and stationary.
- Infinitely differentiable → sometimes deemed unrealistic for modeling most of the physical processes.

CHARACTERISTIC LENGTH-SCALE

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

 ℓ is called **characteristic length-scale**. Loosely speaking, the characteristic length-scale describes how far you need to move in input space for the function values to become uncorrelated. Higher ℓ induces smoother functions, lower ℓ induces more wiggly functions.



× 0 0 × × ×

CHARACTERISTIC LENGTH-SCALE / 2

For $p \ge 2$ dimensions, the squared exponential can be parameterized:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}\left(\mathbf{x} - \mathbf{x}'\right)^{\top} \mathbf{M}\left(\mathbf{x} - \mathbf{x}'\right)\right)$$

Possible choices for the matrix **M** include

$$M_1 = \ell^{-2} I$$
 $M_2 = \operatorname{diag}(\ell)^{-2}$ $M_3 = \Gamma \Gamma^{\top} + \operatorname{diag}(\ell)^{-2}$

where ℓ is a *p*-vector of positive values and Γ is a $p \times k$ matrix.

The 2nd (and most important) case can also be written as

$$k(\mathbf{d}) = \exp\left(-\frac{1}{2}\sum_{j=1}^{p}\frac{d_{j}^{2}}{l_{j}^{2}}\right)$$

× 0 0 × × ×

CHARACTERISTIC LENGTH-SCALE / 3

What is the benefit of having an individual hyperparameter ℓ_i for each dimension?

- The ℓ_1, \ldots, ℓ_p hyperparameters play the role of **characteristic** length-scales.
- Loosely speaking, l_i describes how far you need to move along axis *i* in input space for the function values to be uncorrelated.
- Such a covariance function implements automatic relevance determination (ARD), since the inverse of the length-scale ℓ_i determines the relevancy of input feature *i* to the regression.
- If ℓ_i is very large, the covariance will become almost independent of that input, effectively removing it from inference.
- If the features are on different scales, the data can be automatically rescaled by estimating *l*₁,..., *l*_p

× × 0 × × ×

CHARACTERISTIC LENGTH-SCALE / 4



× 0 0 × × ×

For the first plot, we have chosen $\mathbf{M} = \mathbf{I}$: the function varies the same in all directions. The second plot is for $\mathbf{M} = \text{diag}(\ell)^{-2}$ and $\ell = (1,3)$: The function varies less rapidly as a function of x_2 than x_1 as the length-scale for x_1 is less. In the third plot $\mathbf{M} = \Gamma\Gamma^T + \text{diag}(\ell)^{-2}$ for $\Gamma = (1, -1)^T$ and $\ell = (6, 6)^T$. Here Γ gives the direction of the most rapid variation. (Image from Rasmussen & Williams, 2006)