

# Gaussian Process Cheatsheet

Shuai Guo

January 26, 2021

## 1 Gaussian Process (GP) Definition

In its basic form, a Gaussian Process  $f(\cdot)$  is fully characterized by a mean  $\mu$ , a process variance  $\sigma^2$ , and a **kernel function**  $K(x, x^*)$ , such that a finite collection of  $\mathbf{f} = [f(x_1), f(x_2), \dots, f(x_n)]$  follows a multivariate Gaussian distribution, i.e.,

$$\mathbf{f} \sim \mathcal{N}(\mathbf{1}\mu, \sigma^2 \mathbf{K}), \quad (1)$$

where  $\mathbf{1}$  is a vector with  $n$  ones, and  $\mathbf{K}$  is the correlation matrix, with its element  $\mathbf{K}_{i,j} = K(x_i, x_j)$ .

## 2 Gaussian Kernel Function

A one-dimensional Gaussian kernel  $K(x_i, x_j)$  is expressed as:

$$K(x_i, x_j) = e^{-\theta(x_i - x_j)^2}, \quad (2)$$

where  $\theta$  is a kernel parameter that controls the correlation strength. Similarly, a  $m$ -dimensional Gaussian kernel is expressed as:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp \left[ - \sum_{k=1}^m \theta_k (x_i^k - x_j^k)^2 \right], \quad (3)$$

which is simply a series of multiplication of the one-dimensional Gaussian kernel for each feature. Here, we have the kernel parameters  $\boldsymbol{\theta} = [\theta_1, \theta_2, \dots, \theta_m]$ .

## 3 GP Model Training

Maximum likelihood estimation is used to derive  $\mu, \sigma^2$ , and  $\boldsymbol{\theta}$ . The likelihood  $L$  of observing the labels  $(y_1, y_2, \dots, y_n)$  of the training instances  $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$  is expressed as:

$$L(\mathbf{y}|\mu, \sigma^2, \boldsymbol{\theta}) = \frac{1}{\sqrt{(2\pi\sigma^2)^n |\mathbf{K}|}} \exp \left[ - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{1}\mu)^T \mathbf{K}^{-1} (\mathbf{y} - \mathbf{1}\mu) \right], \quad (4)$$

where  $\mathbf{y} = [y_1, y_2, \dots, y_n]$  and  $\mathbf{K}$  is the correlation matrix of the training instances.

In practice, the logarithm of the likelihood  $L$  is maximized to avoid round-off error:

$$\begin{aligned} \ln(L) = & -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2} \ln(|\mathbf{K}|) \\ & - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{1}\mu)^T \mathbf{K}^{-1} (\mathbf{y} - \mathbf{1}\mu) \end{aligned} \quad (5)$$

By setting the derivatives of  $\ln(L)$  with respect to  $\mu$  and  $\sigma^2$  to zero, we can derive the analytical expressions for their optimum values:

$$\mu = (\mathbf{1}^T \mathbf{K}^{-1} \mathbf{1})^{-1} \mathbf{1}^T \mathbf{K}^{-1} \mathbf{y} \quad (6)$$

$$\sigma^2 = \frac{1}{n} (\mathbf{y} - \mathbf{1}\mu)^T \mathbf{K}^{-1} (\mathbf{y} - \mathbf{1}\mu) \quad (7)$$

For  $\theta$ , its estimation requires solving an auxiliary optimization problem:

$$\theta = \operatorname{argmax}_{\theta} \left[ -\frac{n}{2} \ln(\sigma^2) - \frac{1}{2} \ln(|\mathbf{K}|) \right]. \quad (8)$$

Eq. (8) is obtained via substituting Eq. (7) into Eq. (5) and removing the constant term  $-n/2 \ln(2\pi)$ .

## 4 GP Model Prediction

To predict  $f^*$  at  $\mathbf{x}^*$  with a trained GP model, first of all, we write out the joint distribution of  $f^*$  and  $\mathbf{y}$  (i.e., the observed labels of the training instances):

$$\begin{pmatrix} \mathbf{y} \\ f^* \end{pmatrix} \sim \mathcal{N} \left( \mu, \sigma^2 \begin{pmatrix} \mathbf{K} & \mathbf{k}^* \\ \mathbf{k}^{*T} & 1 \end{pmatrix} \right), \quad (9)$$

where  $\mathbf{k}^*$  is a correlation vector between the testing and training instances, with its  $i$ th element being  $k_i^* = K(x^*, x_i)$ .

In a second step, we derive the distribution of  $f^*$  conditioned on  $\mathbf{y}$  from their joint distribution. This conditional distribution of  $f^*$  is written as  $f^* | \mathbf{y} \sim \mathcal{N}(\mu^*, \Sigma^*)$ , with

$$\mu^* = \mu + \mathbf{k}^{*T} \mathbf{K}^{-1} (\mathbf{y} - \mathbf{1}\mu) \quad (10)$$

$$\Sigma^* = \sigma^2 (1 - \mathbf{k}^{*T} \mathbf{K}^{-1} \mathbf{k}^*). \quad (11)$$

$f^* | \mathbf{y} \sim \mathcal{N}(\mu^*, \Sigma^*)$  fully characterizes the GP prediction at  $\mathbf{x}^*$ .