

$$K(x, x') = \sigma^2 \exp\left(-\frac{|x - x'|^2}{2l^2}\right) \quad (11)$$

$$K(x, x') = \sigma^2 \exp\left(\sum_{i=1}^k -\frac{|x_i - x'_i|^2}{2l_i^2}\right) \quad (12)$$

$$\Sigma(x, x') = K(x, x') + I\sigma_y \quad (13)$$

Equation 12 shows the kernel function extended to multidimensional problem. for k dimensions there are $k+3$ hyper-parameters. l is the scale length or the horizontal scaling, effectively how quickly the correlation between two points decays for each dimension. σ is the vertical scaling. σ_y , shown in equation 13, is a representation of noise in the evaluations, this maintains some uncertainty around evaluation points. In hyper-parameter optimisation and other applications with noisy evaluations, Gaussian noise is added to the covariance matrix to avoid over-fitting. These hyper-parameters have a significant effect of the expression of the model and can be key to the final result, in particular l . Because of this, these settings are normally dealt with automatically rather than hand tuned. One approach used in [9] is to integrate over the hyper-parameters using Monte Carlo estimates, however this can be computationally expensive. Another common approach is to use a marginal likelihood estimates and optimise these settings via gradient descent.

Interesting idea
↓
explain rationale

$$p(f|\theta, D) = \mathcal{GP}(0, K(x, x')) \quad (14)$$

$$p(y|\theta, D) = \mathcal{GP}(0, K(x, x') + I\sigma_y) \quad (15)$$

Equations 14 and 15 show a definition of Gaussian processes over an noiseless and noisy function respectively. For a point of interested, $y(x')$, the Gaussian process can be considered a joint distribution over the $y(x')$ and the query observation pair history D . Using the marginalisation property of gaussians this can be restructured as equation 16. An estimate of y at x' is simply $y(x')$ conditioned on $y(x)$, shown in equation 17

$$p(y(x), y(x')) \sim \mathcal{N}\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} K(x, x) & K(x, x') \\ K(x, x')^T & K(x', x') \end{pmatrix}\right) \quad (16)$$

$$p(y(x')|y(x)) = \frac{p(x, x')}{p(x)} \sim \mathcal{N}(\mu', \sigma') \quad (17)$$

$$\mu' = K(x, x')K(x, x)^{-1}y(x) \quad (18)$$

$$\sigma' = K(x', x') - K(x, x')K(x, x)^{-1}K(x, x')^T \quad (19)$$