

# 1 GP

Gaussian processes are mathematically equivalent to many well known models, including Bayesian linear models, spline models, large neural networks (under suitable conditions).

## 1.1 weight-space view

Projecting the inputs into a high-dimensional feature space and applying the linear model there.

## 1.2 function-space view

Defining a distribution over functions, and inference taking place directly in the space of functions.

A Gaussian process is a collection of random variables, any finite number of which have a joint Gaussian distribution.

$$\begin{aligned} m(\mathbf{x}) &= \mathbb{E}[f(\mathbf{x})] \\ k(\mathbf{x}, \mathbf{x}') &= \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))] \end{aligned} \quad (1)$$

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (2)$$

### Prediction with Noise-free Observations

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(X, X) & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right) \quad (3)$$

### Prediction using Noisy Observations

$$\text{cov}(y_p, y_q) = k(\mathbf{x}_p, \mathbf{x}_q) + \sigma_n^2 \delta_{pq} \text{ or } \text{cov}(\mathbf{y}) = K(X, X) + \sigma_n^2 I \quad (4)$$

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(X, X) + \sigma_n^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix}\right) \quad (5)$$

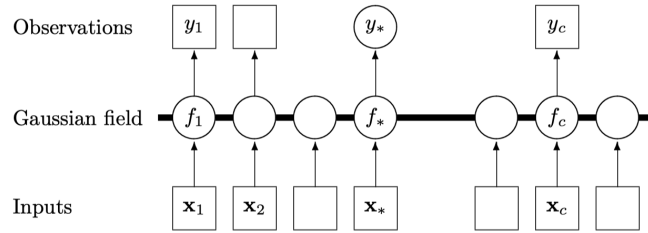


Figure 1: Graphical Model for GP

$$\mathbf{f}_* | X, \mathbf{y}, X_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)) \quad (6)$$

$$\bar{\mathbf{f}}_* \triangleq \mathbb{E}[\mathbf{f}_* | X, \mathbf{y}, X_*] = K(X_*, X) [K(X, X) + \sigma_n^2 I]^{-1} \mathbf{y} \quad (7)$$

$$\text{cov}(\mathbf{f}_*) = K(X_*, X_*) - K(X_*, X) [K(X, X) + \sigma_n^2 I]^{-1} K(X, X_*) \quad (8)$$

$$\begin{aligned} \bar{f}_* &= \mathbf{k}_*^\top (K + \sigma_n^2 I)^{-1} \mathbf{y} \\ \mathbb{V}[f_*] &= k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{k}_*^\top (K + \sigma_n^2 I)^{-1} \mathbf{k}_* \end{aligned} \quad (9)$$

## 2 Varying the Hyperparameters

Of course we can take the position of a quickly-varying signal with low noise, or a slowly-varying signal with high noise to extremes; the former would give rise to a white-noise process model for the signal, while the latter would give rise to a constant signal with added white noise.

## 3 Decision Theory for Regression

In general the value of  $y_{\text{guess}}$  that minimizes the risk for the loss function  $|y_{\text{guess}} - y_*|$  is the median of  $p(y_* | \mathbf{x}_*, \mathcal{D})$ , while for the squared loss  $(y_{\text{guess}} - y_*)^2$  it is the mean of this distribution. When the predictive distribution is Gaussian, the mean and the median coincide.

## References

- [1] Christopher K Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*, volume 2. MIT press Cambridge, MA, 2006.