

# AutoML: Gaussian Processes

## Mean Functions for Gaussian Processes

Bernd Bischl   Frank Hutter   Lars Kotthoff  
Marius Lindauer   Joaquin Vanschoren

# The Role of Mean Functions I

- It is common but by no means necessary to consider GPs with zero-mean functions:

$$m(\mathbf{x}) \equiv 0$$

- This is not necessarily a drastic limitation, since the mean of the posterior process is not confined to be zero:

$$\mathbf{f}_* | \mathbf{X}_*, \mathbf{X}, \mathbf{f} \sim \mathcal{N}(\mathbf{K}_*^\top \mathbf{K}^{-1} \mathbf{f}, \mathbf{K}_{**} - \mathbf{K}_*^\top \mathbf{K}^{-1} \mathbf{K}_*).$$

- There are several reasons why one might wish to explicitly model a mean function, including interpretability, convenience of expressing prior informations, etc.

# The Role of Mean Functions II

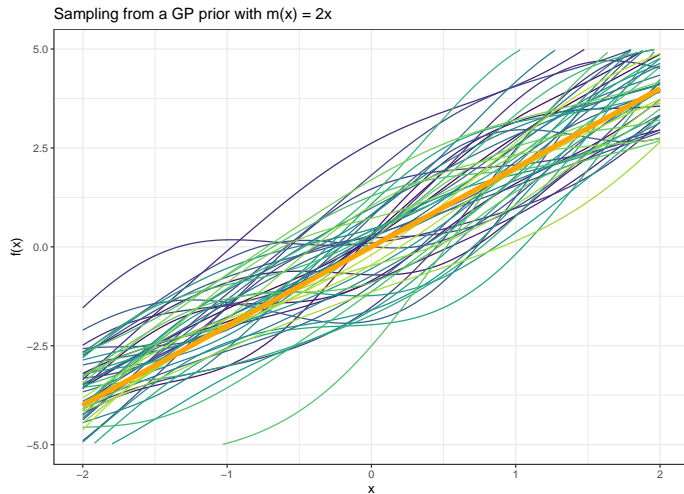
- When assuming a non-zero mean GP prior  $\mathcal{G}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$  with mean  $m(\mathbf{x})$ , the predictive mean becomes:

$$m(\mathbf{X}_*) + \mathbf{K}_* \mathbf{K}_y^{-1} (\mathbf{y} - m(\mathbf{X}))$$

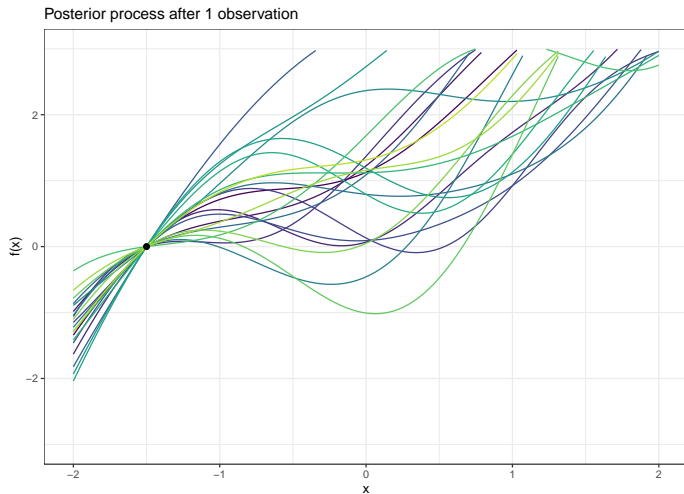
but the predictive variance remains unchanged.

- Gaussian processes with non-zero mean Gaussian process priors are also called **Gaussian processes with trend**.

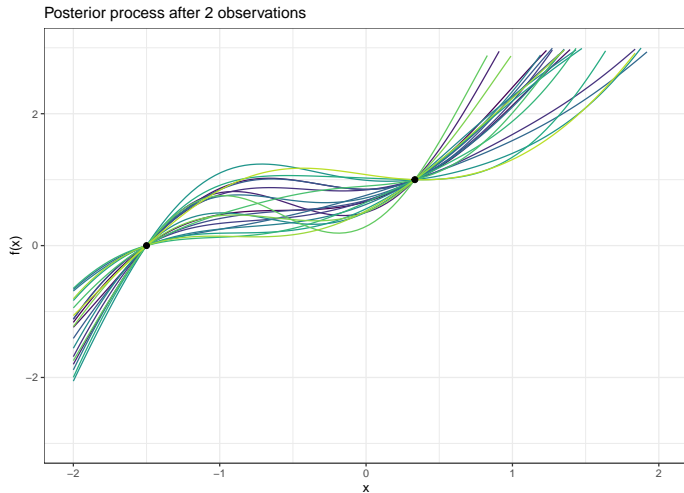
# The Role of Mean Functions III



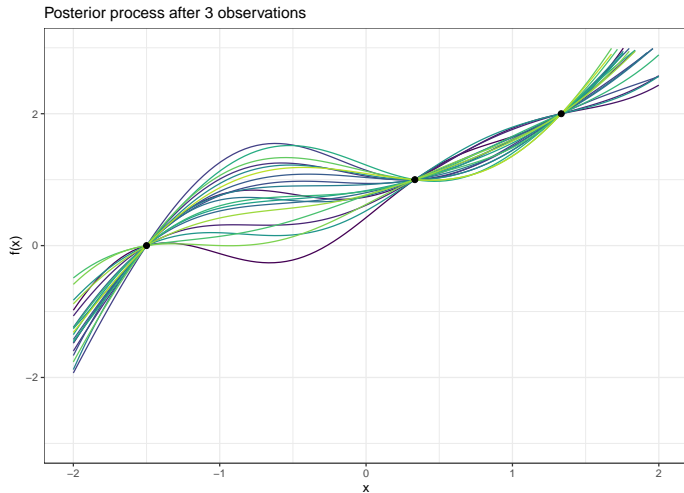
# The Role of Mean Functions IV



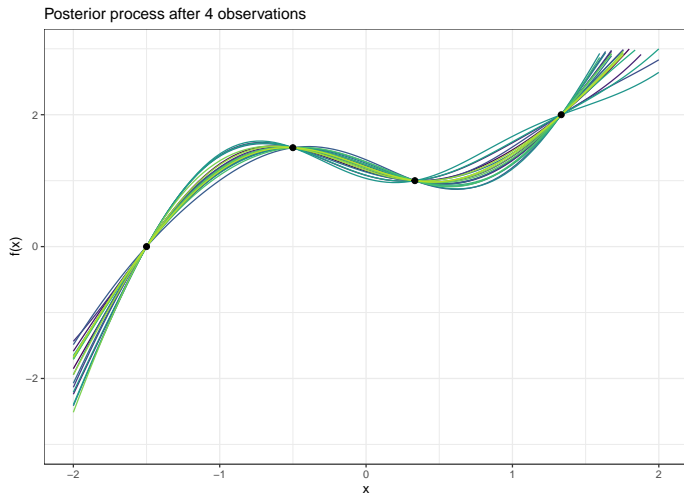
# The Role of Mean Functions $V$



# The Role of Mean Functions VI



# The Role of Mean Functions VII





# The Role of Mean Functions VIII

- In practice it is often difficult to specify a fixed mean function. Instead, it may be more convenient to specify a few fixed basis functions, whose coefficients,  $\beta$ , are to be inferred from the data.

- Consider

$$g(\mathbf{x}) = f(\mathbf{x}) + \mathbf{h}(\mathbf{x})^\top \beta, \text{ with } f(\mathbf{x}) \sim \mathcal{G}(0, k(\mathbf{x}, \mathbf{x}'))$$

where  $f(\mathbf{x})$  is a zero mean GP,  $\mathbf{h}(\mathbf{x})$  are a set of fixed basis functions, and  $\beta$  are additional parameters.

- This formulation expresses that the data is close to a global linear model with the residuals being modelled by a GP.
- By taking the prior on  $\beta$  to be Gaussian,  $\beta \sim \mathcal{N}(\mathbf{b}, B)$ , another GP with an added contribution in the covariance function can be obtained:

$$g(\mathbf{x}) \sim \mathcal{G}\left(\mathbf{h}(\mathbf{x})^\top \mathbf{b}, k(\mathbf{x}, \mathbf{x}') + \mathbf{h}(\mathbf{x})^\top B \mathbf{h}(\mathbf{x}')\right)$$

# The Role of Mean Functions IX

- The mean and covariance functions of  $g(\mathbf{x})$  will be then:

$$\begin{aligned}\bar{\mathbf{g}}(X_*) &= \bar{\mathbf{f}}(X_*) + R^\top \bar{\boldsymbol{\beta}}, \\ \text{cov}(\mathbf{g}_*) &= \text{cov}(\mathbf{f}_*) + R^\top (B^{-1} + HK_y^{-1}H^\top)^{-1}R,\end{aligned}$$

where  $\bar{\boldsymbol{\beta}} = (B^{-1} + HK_y^{-1}H^\top)^{-1}(HK_y^{-1}\mathbf{y} + B^{-1}\mathbf{b})$  and  $R = H_* - HK_y^{-1}K_*$ .

- Notice that  $\bar{\boldsymbol{\beta}}$  is the mean of the global linear model parameters and the  $H$  matrix collects the  $\mathbf{h}(\mathbf{x})$  vectors for all training and test cases.
- Hence, the predictive mean is the sum of the mean linear output and what the GP model predicts from the residuals. The covariance is the usual covariance term plus a non negative contribution.

# The Role of Mean Functions X

- In the limiting case where the prior on the  $\beta$  becomes vague ( $B^{-1} \rightarrow O$ ), the predictive distribution becomes independent of  $\mathbf{b}$ :

$$\begin{aligned}\bar{\mathbf{g}}(X_*) &= \bar{\mathbf{f}}(X_*) + R^\top \bar{\beta}, \\ \text{cov}(\mathbf{g}_*) &= \text{cov}(\mathbf{f}_*) + R^\top (H K_y^{-1} H^\top)^{-1} R,\end{aligned}$$

where the limiting  $\bar{\beta} = (H K_y^{-1} H^\top)^{-1} H K_y^{-1} \mathbf{y}$ .

- To make predictions under the limit  $B^{-1} \rightarrow O$ , it is important to use the above equation. Otherwise, by naively plugging the modified covariance function into the standard prediction equations, the entries of the covariance function tend to infinity making it unsuitable for numerical implementation.

# The Role of Mean Functions XI

- The marginal likelihood for the model with a Gaussian prior  $\beta \sim \mathcal{N}(\mathbf{b}, B)$  can be expressed as what follows (the explicit mean has been included).

$$\begin{aligned}\log p(\mathbf{y}|X, \mathbf{b}, B) = & -\frac{1}{2} \left( H^\top \mathbf{b} - \mathbf{y} \right)^\top \left( K_y + H^\top B H \right)^{-1} \left( H^\top \mathbf{b} - \mathbf{y} \right) \\ & - \frac{1}{2} \log \left| K_y + H^\top B H \right| - \frac{n}{2} \log 2\pi.\end{aligned}$$

- We are interested in exploring the limit where  $B^{-1} \rightarrow O$ , i.e. when the prior is vague. In this limit the mean of the prior is irrelevant, so without loss of generality we assume for now that the mean is zero,  $\mathbf{b} = \mathbf{0}$ .

## The Role of Mean Functions XII

- This assumption gives:

$$\begin{aligned}\log p(\mathbf{y}|X, \mathbf{b} = \mathbf{0}, B) = & -\frac{1}{2}\mathbf{y}^\top K_y^{-1}\mathbf{y} + \frac{1}{2}\mathbf{y}^\top C\mathbf{y} \\ & -\frac{1}{2}\log |K_y| - \frac{1}{2}\log |B| - \frac{1}{2}\log |A| - \frac{n}{2}\log 2\pi,\end{aligned}$$

where  $A = B^{-1} + HK_y^{-1}H^\top$  and  $C = K_y^{-1}H^\top A^{-1}HK_y^{-1}$ .

- The behaviour of the log marginal likelihood in the limiting case can be explored now. The variances of the Gaussian in the directions spanned by columns of  $H^\top$  will become infinite, and it will require special treatment.
- 💡 The log marginal likelihood consists of three terms: a quadratic form in  $\mathbf{y}$ , a log determinant term, and a term involving  $\log 2\pi$ .

## The Role of Mean Functions XIII

- Performing an eigendecomposition of the covariance matrix, the contributions to quadratic form term from the infinite-variance directions will be zero. However, the log determinant term will tend to minus infinity.
- The standard solution in this case is to project  $\mathbf{y}$  onto the directions orthogonal to the span of  $H^\top$  and compute the marginal likelihood in this subspace.
- By taking  $m$  to denote the rank of  $H^\top$ , we can discard the terms  $-\frac{1}{2} \log |B| - \frac{m}{2} \log 2\pi$  from the previous equation to give:

$$\log p(\mathbf{y}|X) = -\frac{1}{2} \mathbf{y}^\top K_y^{-1} \mathbf{y} + \frac{1}{2} \mathbf{y}^\top C \mathbf{y} - \frac{1}{2} \log |K_y| - \frac{1}{2} \log |A| - \frac{n-m}{2} \log 2\pi,$$

where  $A = HK_y^{-1}H^\top$  and  $C = K_y^{-1}H^\top A^{-1}HK_y^{-1}$ .