

Supervised Learning :: CHEAT SHEET

Bayesian Linear Model

Bayesian Linear Model:

$$y^{(i)} = f(\mathbf{x}^{(i)}) + \epsilon^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} + \epsilon^{(i)}, \quad \text{for } i \in \{1, \dots, n\}$$

where $\epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$.

Parameter vector $\boldsymbol{\theta}$ is stochastic and follows a distribution.

Gaussian variant:

- Prior distribution: $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I}_p)$
- Posterior distribution: $\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\sigma^{-2} \mathbf{K}^{-1} \mathbf{X}^\top \mathbf{y}, \mathbf{K}^{-1})$ with $\mathbf{K} := \sigma^{-2} \mathbf{X}^\top \mathbf{X} + \frac{1}{\tau^2} \mathbf{I}_p$
- Predictive distribution of $y_* = \boldsymbol{\theta}^\top \mathbf{x}_*$ for a new observations \mathbf{x}_* :
$$y_* \mid \mathbf{X}, \mathbf{y}, \mathbf{x}_* \sim \mathcal{N}(\sigma^{-2} \mathbf{y}^\top \mathbf{X} \mathbf{A}^{-1} \mathbf{x}_*, \mathbf{x}_*^\top \mathbf{A}^{-1} \mathbf{x}_*)$$

Gaussian Processes

Weight-Space View	Function-Space View
Parameterize functions	Work on functions directly
Example: $f(\mathbf{x} \mid \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{x}$	
Define distributions on $\boldsymbol{\theta}$	Define distributions on f
Inference in parameter space Θ	Inference in function space \mathcal{H}

Gaussian Processes: A function $f(\mathbf{x})$ is generated by a GP $\mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$ if for **any finite** set of inputs $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$, the associated vector of function values $\mathbf{f} = (f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(n)}))$ has a Gaussian distribution

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

with

$$\mathbf{m} := \left(m(\mathbf{x}^{(i)}) \right)_i, \quad \mathbf{K} := \left(k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \right)_{ij},$$

where $m(\mathbf{x})$ is the mean function and $k(\mathbf{x}, \mathbf{x}')$ is the covariance function.

Types of **covariance functions**:

- $k(., .)$ is stationary if it is as a function of $\mathbf{d} = \mathbf{x} - \mathbf{x}'$, $\rightsquigarrow k(\mathbf{d})$
- $k(., .)$ is isotropic if it is a function of $r = \|\mathbf{x} - \mathbf{x}'\|$, $\rightsquigarrow k(r)$
- $k(., .)$ is a dot product covariance function if k is a function of $\mathbf{x}^T \mathbf{x}'$

Commonly used covariance functions:

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

Gaussian Processes Prediction

Posterior Process

Assuming a zero-mean GP prior $\mathcal{GP}(\mathbf{0}, k(\mathbf{x}, \mathbf{x}'))$. For $f_* = f(\mathbf{x}_*)$ on single unobserved test point \mathbf{x}_*

$$f_* \mid \mathbf{x}_*, \mathbf{X}, \mathbf{f} \sim \mathcal{N}(\mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{f}, \mathbf{k}_{**} - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_*),$$

where, $\mathbf{K} = (k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}))_{ij}$, $\mathbf{k}_* = [k(\mathbf{x}_*, \mathbf{x}^{(1)}), \dots, k(\mathbf{x}_*, \mathbf{x}^{(n)})]$ and $\mathbf{k}_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$.

For multiple unobserved test points $\mathbf{f}_* = [f(\mathbf{x}_*^{(1)}), \dots, f(\mathbf{x}_*^{(m)})]$:

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

with $\mathbf{K}_* = (k(\mathbf{x}^{(i)}, \mathbf{x}_*^{(j)}))_{ij}$, $\mathbf{K}_{**} = (k(\mathbf{x}_*^{(i)}, \mathbf{x}_*^{(j)}))_{ij}$.

Predictive mean when assuming a non-zero mean GP prior $\mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$ with mean $m(\mathbf{x})$:

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

Predictive variance remains unchanged.

Noisy Posterior Process

Assuming a zero-mean GP prior $\mathcal{GP}(\mathbf{0}, k(\mathbf{x}, \mathbf{x}'))$:

$$\mathbf{f}_* \mid \mathbf{X}_*, \mathbf{X}, \mathbf{y} \sim \mathcal{N}(\mathbf{m}_{\text{post}}, \mathbf{K}_{\text{post}}).$$

with nugget σ^2 and

$$\mathbf{m}_{\text{post}} = \mathbf{K}_*^T (\mathbf{K} + \sigma^2 \cdot \mathbf{I})^{-1} \mathbf{y}$$
$$\mathbf{K}_{\text{post}} = \mathbf{K}_{**} - \mathbf{K}_*^T (\mathbf{K} + \sigma^2 \cdot \mathbf{I})^{-1} \mathbf{K}_*,$$

Predictive mean when assuming a non-zero mean GP prior $\mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$ with mean $m(\mathbf{x})$:

$$m(\mathbf{X}_*) + \mathbf{K}_* (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} (\mathbf{y} - m(\mathbf{X}))$$

Predictive variance remains unchanged.

Train a Gaussian Processes

We can learn the numerical hyperparameters of a selected covariance function directly during GP training.

Let us assume

$$y = f(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2),$$

where $f(\mathbf{x}) \sim \mathcal{GP}(\mathbf{0}, k(\mathbf{x}, \mathbf{x}' \mid \boldsymbol{\theta}))$.

Observing $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I})$, the marginal log-likelihood (or evidence) is

$$\log p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) = \log \left[(2\pi)^{-n/2} |\mathbf{K}_y|^{-1/2} \exp \left(-\frac{1}{2} \mathbf{y}^\top \mathbf{K}_y^{-1} \mathbf{y} \right) \right]$$
$$= -\frac{1}{2} \mathbf{y}^\top \mathbf{K}_y^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}_y| - \frac{n}{2} \log 2\pi.$$

with $\mathbf{K}_y := \mathbf{K} + \sigma^2 \mathbf{I}$ and $\boldsymbol{\theta}$ denoting the hyperparameters (the parameters of the covariance function).

The three terms of the marginal likelihood have interpretable roles, considering that the model becomes less flexible as the length-scale increases:

- the data fit $-\frac{1}{2} \mathbf{y}^\top \mathbf{K}_y^{-1} \mathbf{y}$, which tends to decrease if the length scale increases
- the complexity penalty $-\frac{1}{2} \log |\mathbf{K}_y|$, which depends on the covariance function only and which increases with the length-scale, because the model gets less complex with growing length-scale
- a normalization constant $-\frac{n}{2} \log 2\pi$