

Reading Note 9 for Gaussian Process

Xiang Pan

October 16, 2021

1 Introduction

In short, the ability for a model to learn from data is determined by:

1. The support of the model: what solutions we think are a priori possible.
2. The inductive biases of the model: what solutions we think are a priori likely.

The **capacity (flexibility)** of a model \mathcal{M}_i can be defined as the mutual information between the data \mathbf{y} (at N locations X) and predictions made by the model \mathbf{y}_* (at test locations X_*)

$$I_{i,N} = \sum_{\mathbf{y}, \mathbf{y}_*} p(\mathbf{y}, \mathbf{y}_* | \mathcal{M}_i) \log \frac{p(\mathbf{y}, \mathbf{y}_* | \mathcal{M}_i)}{p(\mathbf{y} | \mathcal{M}_i) p(\mathbf{y}_* | \mathcal{M}_i)} \quad (1)$$

$$I_{i,N} = p(\mathbf{y}) \int p(\mathbf{y}_* | \mathbf{y}) \log \frac{p(\mathbf{y}_* | \mathbf{y})}{p(\mathbf{y}_*)} d\mathbf{y}_* \quad (2)$$

2 GP

We are ultimately more interested in – and have stronger intuitions about – the functions that model data than the weights \mathbf{w} in a parametric model, and we can express those intuitions with a covariance kernel.

$$\begin{aligned} p(\mathbf{y}_* | \mathbf{y}) &= \int p(\mathbf{y}_* | f(x)) p(f(x) | \mathbf{y}) df(x) \\ p(f(x) | \mathbf{y}) &\propto p(\mathbf{y} | f(x)) p(f(x)) \end{aligned} \quad (3)$$

$$p(f_* | \mathbf{y}) = \int p(f_* | \mathbf{f}) p(\mathbf{f} | \mathbf{y}) d\mathbf{f} \quad (4)$$

Dependency: Hyperparameters \rightarrow Parameters \rightarrow Data

$$\log p(\mathbf{y} | \boldsymbol{\theta}, X) = \underbrace{-\frac{1}{2} \mathbf{y}^\top (K_{\boldsymbol{\theta}} + \sigma^2 I)^{-1} \mathbf{y}}_{\text{model fit}} - \underbrace{\frac{1}{2} \log |K_{\boldsymbol{\theta}} + \sigma^2 I|}_{\text{complexity penalty}} - \frac{N}{2} \log(2\pi) \quad (5)$$

Prediction

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

$$p(\mathbf{f}_* | X_*, X, \mathbf{y}) = \int p(\mathbf{f}_* | X_*, X, \mathbf{y}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} \quad (7)$$

$$p(\boldsymbol{\theta} | \mathbf{y}) \propto p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) \quad (8)$$

3 Kernel

A kernel is **stationary** if it is invariant to translations of the inputs.

Covariance function	Expression	Stationary
Constant	a_0	Yes
Linear	$x \cdot x'$	No
Polynomial	$(x \cdot x' + a_0)^p$	No
Squared Exponential	$\exp\left(-\frac{ x-x' ^2}{2l^2}\right)$	Yes
Matérn	$\frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} x-x' }{l}\right)^\nu K_\nu\left(\frac{\sqrt{2\nu} x-x' }{l}\right)$	Yes
Ornstein-Uhlenbeck	$\exp\left(-\frac{ x-x' }{l}\right)$	Yes
Rational Quadratic	$\left(1 + \frac{ x-x' ^2}{2\alpha l^2}\right)^{-\alpha}$	Yes
Periodic	$\exp\left(-\frac{2\sin^2\left(\frac{x-x'}{2}\right)}{l^2}\right)$	Yes
Gibbs	No	No
Spectral Mixture	$\sum_{q=1}^Q w_q \prod_{p=1}^P \exp\left\{-2\pi^2 (x - x')_p^2 v_{qp}\right\} \cos\left(2\pi (x - x')_p \mu_{qp}\right)$	Yes

4 Mean Function

The mean function is also a powerful way to encode assumptions (inductive biases) into a Gaussian process model, the Gaussian process can leverage the assumptions of a parametric model through a mean function and also reflect the belief that the parametric form of that model will not be entirely accurate.

5 Feature Of GP

- Expressive Kernels
- Exact Efficient Inference: Efficiently determine the eigenvalues of a covariance matrix K
- Multi-Output Gaussian Processes
- Sampling Kernel Hyperparameters

References

- [1] Andrew Gordon Wilson. *Covariance kernels for fast automatic pattern discovery and extrapolation with Gaussian processes*. PhD thesis, Citeseer, 2014.