

# Gaussian Process

qhy

2017 年 8 月 18 日

## 目录

<b>1 Gaussian Process</b>	<b>2</b>
1.1 数学基础	2
1.1.1 高斯分布	2
1.1.2 高斯变量的贝叶斯定理	2
1.2 随机过程与高斯过程	3
1.3 高斯过程回归 Gaussian Process Regression	3
1.3.1 无噪声GPR	3
1.3.2 带噪声高斯过程回归	4
1.4 高斯过程的图解	5
1.5 超参数的学习	5
1.5.1 均值	5
1.5.2 协方差矩阵 $\mathbf{K}$	5
1.5.3 学习超参数	6
1.6 自动相关性确定 automatic relevance determination ARD	6
1.7 怎样从高斯分布中采样	6
1.8 与神经网络的关系	7
1.9 与岭回归的关系	7
1.10 与贝叶斯线性回归模型的关系	7
1.11 参考资料	8

# 1 Gaussian Process

## 1.1 数学基础

在GP中需要用到的数学知识,这里给出结论,证明见PRML第二章。

### 1.1.1 高斯分布

给定一个联合高斯分布 $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,其中 $\boldsymbol{\Lambda} \equiv \boldsymbol{\Sigma}^{-1}$ ,且

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}, \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad (1)$$

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}, \boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix} \quad (2)$$

有以下两个结论:

条件概率分布:

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Lambda}_{aa}^{-1}) \quad (3)$$

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1}\boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \quad (4)$$

或

$$\boldsymbol{\mu}_{ab} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b) \quad (5)$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba} \quad (6)$$

边缘概率分布:

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa}) \quad (7)$$

这里证明了边缘概率和联合分布之间的关系,当我们在求边缘概率的时候,就可以通过求联合分布得到(这里在计算的时候省略的积分的过程,但实际的证明就是通过积分证明得到的。)

### 1.1.2 高斯变量的贝叶斯定理

给定 $\mathbf{x}$ 的一个边缘高斯分布,以及在给定 $\mathbf{x}$ 的条件下的 $\mathbf{y}$ 的条件高斯分布,形式为:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \quad (8)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{Ax} + \mathbf{b}, \mathbf{L}^{-1}) \quad (9)$$

$\mathbf{y}$ 的边缘分布以及给定 $\mathbf{y}$ 条件下的 $\mathbf{x}$ 的条件分布为:

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{Ax} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T) \quad (10)$$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma}) \quad (11)$$

其中,

$$\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1} \quad (12)$$

## 1.2 随机过程与高斯过程

**随机过程:**许多随机变量的集合

**高斯过程:**许多高斯随机变量的集合

e.g. 给定 $n$ 个随机变量: $\mathbf{y} = \{y_1, y_2, \dots, y_n\}$ , 其中 $y_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

高斯随机过程就是:

$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_n \end{bmatrix} \quad (13)$$

没错,就是这些随机变量的集合。在求这个随机过程的分布的时候,实际上就是求这些随机变量的联合分布。然后,随机过程到底结束。

## 1.3 高斯过程回归 Gaussian Process Regression

### 1.3.1 无噪声GPR

**问题描述:**给定样本 $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$ , 和待预测数据 $x_*$ , 相似性矩阵 $\mathbf{K}$ , 求 $p(y_* | \mathcal{D})$

高斯过程回归:假设目标函数是一条平滑的曲线,在相近的之间,其输出也相近。也就是相近的随机变量(它们采样的值也相近),而这个距离则是通过 $x$ 来确定的,这就保证的相邻的 $x$ 之间的曲线的平滑度(曲线不会突然飙到一个很大或者很小的值)。

下面看高斯过程回归的具体过程:

把每个样本的 $y_i$ 当作随机变量,并且假设这些随机变量满足高斯分布(即,这些随机变量组成高斯过程)。

i.e.  $y_i \sim \mathcal{N}(\mu_i, \Sigma_{ii})$

则这个高斯过程满足以下条件:由前面的数学公式可以得到

$$\mathbf{y} \triangleq \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}) \quad (14)$$

注: $\triangleq$ 表示定义的意思

其中,

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} \quad (15)$$

其中 $\mathbf{K}$ 为相似性矩阵,用来描述随机变量之间的距离。

$$K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m) = \frac{1}{\alpha} \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) \quad (16)$$

根据前面随机变量的定义,有

$$K_{ii} = \Sigma_{ii} \quad (17)$$

给定 $x_*$ ,假设

$$\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu} \\ \mu_* \end{bmatrix}, \begin{bmatrix} \mathbf{K} & \mathbf{K}_* \\ \mathbf{K}_*^T & k_{**} \end{bmatrix}\right) \quad (18)$$

其中,

$$\mathbf{K}_*^T = [k(\mathbf{x}_1, \mathbf{x}_*), \dots, k(\mathbf{x}_n, \mathbf{x}_*)] \quad (19)$$

根据第一部分的数学知识,我们有以下结果:

$$p(y_*|\mathbf{y}) = \mathcal{N}(\tilde{\mu}, \tilde{\Sigma}) \quad (20)$$

其中,

$$\tilde{\mu} = \mu_* + \mathbf{K}_*^T \mathbf{K}^{-1}(\mathbf{y} - \boldsymbol{\mu}) \quad (21)$$

$$\tilde{\Sigma} = k_{**} - \mathbf{K}_*^T \mathbf{K}^{-1} \mathbf{K}_* \quad (22)$$

一般地,我们对均值没有任何的先验知识,根据对称性,我们取 $\boldsymbol{\mu} = \mu_* = \mathbf{0}$

### 1.3.2 带噪声高斯过程回归

考虑目标观测值的噪声,其形式为:

$$t_i = y_i + \epsilon_i \quad (23)$$

其中, $\epsilon_i$ 是一个噪声随机变量,它的值对每个观测 $n$ 是独立的。满足以下形式:

$$\epsilon_i \sim \mathcal{N}(0, \beta^{-1}) \quad (24)$$

因为 $\epsilon_i$ 和 $y_i$ 是独立的,因此,对于 $\epsilon_i$ 来说可以把 $y_i$ 看作常量 我们有以下形式的高斯分布的噪声过程:

$$p(t_i|y_i) = \mathcal{N}(t_i|y_i, \beta^{-1}) \quad (25)$$

其中, $\beta^{-1}$ 是一个超参数,表示噪声的精度。

考虑多个样本,以 $\mathbf{y} = (y_1, \dots, y_n)$ 为条件,以 $\mathbf{t} = (t_1, \dots, t_n)$ 为目标值的联合概率分布。这个分布是一个各向同性的高斯分布,形式为:

$$p(\mathbf{t}|\mathbf{y}) = \mathcal{N}(\mathbf{t}|\mathbf{y}, \beta^{-1} \mathbf{I}_n) \quad (26)$$

根据高斯过程的定义:上一节的内容  $p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K})$

为了找到以输入值 $\mathbf{x}_1, \dots, \mathbf{x}_n$ (即相似性矩阵 $\mathbf{K}$ )为条件的边缘概率分布 $p(\mathbf{t})$ ,我们需要对 $\mathbf{y}$ 积分:

$$p(\mathbf{t}) \triangleq \int p(\mathbf{t}|\mathbf{y})p(\mathbf{y})d\mathbf{y} = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C}) \quad (27)$$

这里计算的时候,并没有进行积分,只是求对应的联合分布,在联合分布的表达中,找到对应的值(这个对应关系的证明是直接积分得到的)

其中协方差矩阵 $\mathbf{C}$ 为:  $\mathbf{C} = \mathbf{K} + \beta^{-1} \mathbf{I}_n$

给定 $x_*$ ,则对应输出 $y_*$ 与样本 $\mathbf{t}$ 的联合概率分布为:

$$\begin{bmatrix} \mathbf{t} \\ y_* \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} \mathbf{C} & \mathbf{K}_* \\ \mathbf{K}_*^T & k_{**} \end{bmatrix}\right) \quad (28)$$

则 $p(y_*|\mathbf{t})$ 为高斯分布,形式为:形式和无噪声的情况是一样的,只不过 $\mathbf{K}$ 变成了 $\mathbf{C}$

$$p(y_*|\mathbf{t}) \sim \mathcal{N}(\mathbf{K}_*^T \mathbf{C}^{-1}(\mathbf{y} - \boldsymbol{\mu}), k_{**} - \mathbf{K}_*^T \mathbf{C}^{-1} \mathbf{K}_*) \quad (29)$$

也有求 $t_*$ 的分布而不是 $y_*$ 的分布,过程类似,只不过对应的方差 $k_{**}$ 变成了 $c = k_{**} + \beta^{-1}$

## 1.4 高斯过程的图解

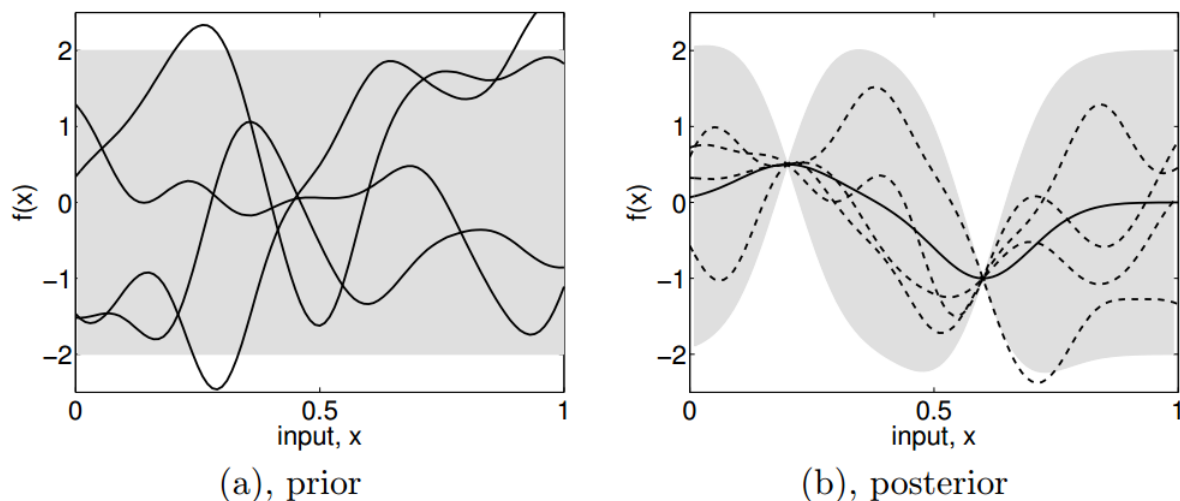


图 1: Panel (a) shows four samples drawn from the prior distribution. Panel (b) shows the situation after two datapoints have been observed. The mean prediction is shown as the solid line and four samples from the posterior are shown as dashed lines. In both plots the shaded region denotes twice the standard deviation at each input value  $x$ .

上图是高斯过程的一个图解: 在先验概率的条件下, 采样空间有无数个多的函数可以选择, 在后验概率的条件下, 则采样的范围就会缩小 而在高斯过程中, 我们进一步把采样的范围缩小到高斯分布的范围内。

这里的函数的采样, 具体地看, 就是对函数的参数  $w$  的采样, 而参数也满足高斯分布, 因此, 参数是从高斯分布采样得到的。(具体见贝叶斯回归与高斯过程的关系)

## 1.5 超参数的学习

在高斯过程中, 我们抛弃参数模型, 直接定义函数上的先验概率分布(具体表现为函数值  $y_i$  的概率分布), 而这个分布则取为高斯分布, 确定一个高斯分布, 只需要找到它的均值和方差即可。

### 1.5.1 均值

前面已经提到, 在大部分应用中, 我们对于函数值  $y_i$  的均值没有任何的先验知识, 因此根据对称性, 我们取均值为 0

### 1.5.2 协方差矩阵 $K$

$K$  矩阵叫相似性矩阵, 用来描述随机变量之间的相似性,

i.e. 对于相似的  $\mathbf{x}_n$  和  $\mathbf{x}_m$ , 对应的函数值  $y(\mathbf{x}_n)$  和  $y(\mathbf{x}_m)$  的相关性要大于不相似的点。这里相似性的概念取决于实际应用。

对于高斯过程回归, 一个广泛使用的和函数的形式为指数项的二次型加上常数项和线性项, i.e.

$$k(\mathbf{x}_n, \mathbf{x}_m) = \theta_0 \exp\left\{-\frac{\theta_1}{2} \|\mathbf{x}_n - \mathbf{x}_m\|\right\} + \theta_2 + \theta_3 \mathbf{x}_n^T \mathbf{x}_m \quad (30)$$

### 1.5.3 学习超参数

学习超参数的方法基于计算似然函数 $p(\mathbf{t}|\boldsymbol{\theta})$ ,其中 $\boldsymbol{\theta}$ 表示高斯过程模型的超参数。最简单的方法是通过最大化似然函数的方法进行 $\boldsymbol{\theta}$ 的点估计。

使用多元高斯分布的标准形式,高斯过程模型的对数似然函数形式为:

$$\ln p(\mathbf{t}|\boldsymbol{\theta}) = -\frac{1}{2} \ln |\mathbf{C}| - \frac{1}{2} \mathbf{t}_T \mathbf{C}^{-1} \mathbf{t} - \frac{n}{2} \ln 2\pi \quad (31)$$

## 1.6 自动相关性确定 automatic relevance determination ARD

上一节用似然函数求参数的方法,这里我们通过为每一个输入变量整合一个单独的参数,在通过最大似然函数方法进行参数的最优化,就能够自动地将不同输入的相对重要性从数据中推断出来。

e.g.考虑二维输入空间 $\mathbf{x} = (x_1, x_2)$ ,有以下形式的核函数:

$$k(\mathbf{x}, \mathbf{x}') = \theta_0 \exp \frac{1}{2} \sum_{i=1}^2 \eta_i (x_i - x'_i)^2 \quad (32)$$

这里的相关性,某个属性的相关性,而不是某个输入样本的相关性

### 1.7 怎样从高斯分布中采样

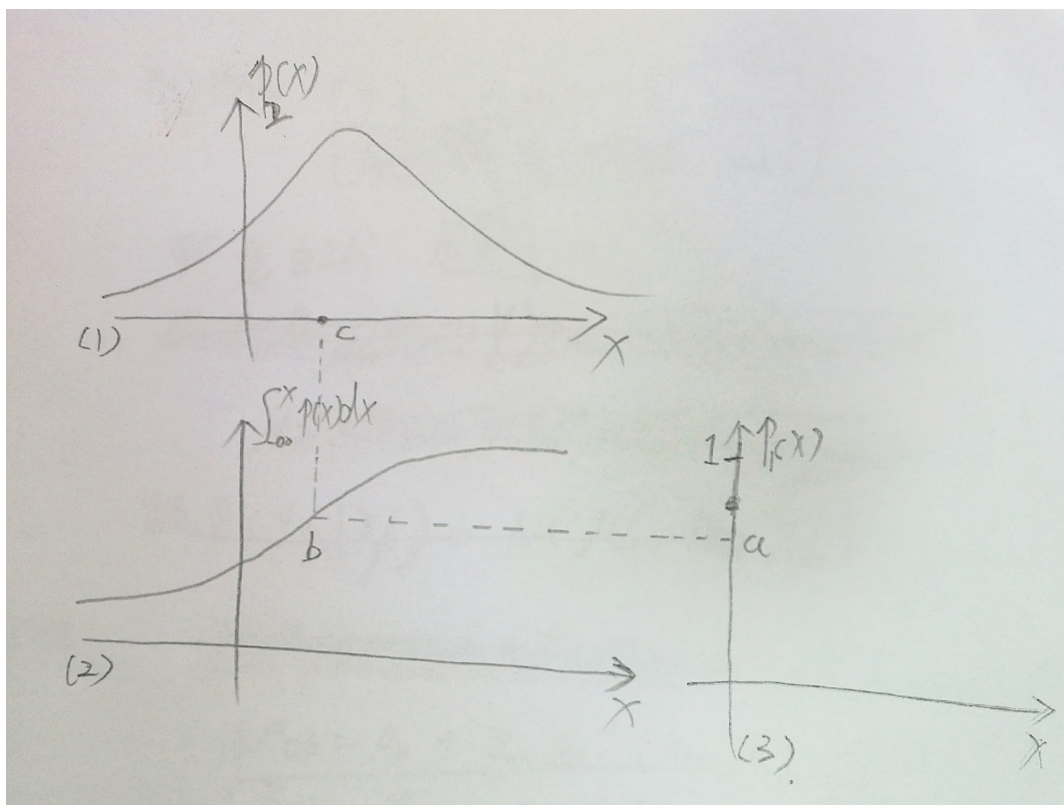


图 2: (1):表示高斯分布,(2)表示从负无穷到 $x$ 曲线下的面积,(3)表示均匀分布

采样过程:从(3)的均匀分布中采样得到点 $a$ ,点 $a$ 对应面积曲线中的点 $b$ ,点 $b$ 对应到高斯分布的点 $c$ ,点 $c$ 就是高斯分布采样的一个结果。

## 1.8 与神经网络的关系

在贝叶斯神经网络中,定义参数 $\mathbf{w}$ 上的先验分布以及网络函数 $f(\mathbf{x}, \mathbf{w})$ 产生了函数 $y(\mathbf{x})$ 上的先验概率分布,其中 $\mathbf{y}$ 表示网络的输出向量。

在 $M \rightarrow \infty$ 的情况下,神经网络产生的函数的分布将会趋于高斯过程。

然而,值得注意的是,在这种期限情况下,神经网络的输出将会变得相互独立。神经网络的优势之一是输出之间共享隐含单元,因此它们可以互相“借统计优势”,即,每个隐含节点关联的权值将被所有的输出变量影响,而不只是被它们中的某一个影响。这个性质在极限的高斯过程中丢失了。

## 1.9 与岭回归的关系

岭回归的优化目标是:

$$\min_{\boldsymbol{\theta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\| + \delta^2 \|\boldsymbol{\theta}\|^2 \quad (33)$$

有解析解:

$$\boldsymbol{\theta} = \mathbf{X}^T \boldsymbol{\alpha} \quad (34)$$

其中,

$$\boldsymbol{\alpha} = (\mathbf{X}\mathbf{X}^T + \delta^2 \mathbf{I}_n)^{-1} \mathbf{y} \quad (35)$$

那么给定 $\mathbf{x}_*$ ,其预测值 $y_*$ 为:

$$y_* = \mathbf{x}_* \boldsymbol{\theta} = \mathbf{x}_* \mathbf{X}^T \boldsymbol{\alpha} = \mathbf{x}_* \mathbf{X}^T (\mathbf{X}\mathbf{X}^T + \delta^2 \mathbf{I}_n)^{-1} \mathbf{y} = \mathbf{K}_*^T \mathbf{C}^{-1} \mathbf{y} \quad (36)$$

即,最后预测的值和高斯过程中, $y_*$ 分布的均值是相等的。

这里体现了高斯过程的一个有点,高斯过程不仅能得到岭回归下的预测结果(均值),还能得到对于每个 $y_*$ 预测的上下边界。

## 1.10 与贝叶斯线性回归模型的关系

贝叶斯线性回归模型:具体见PRML3.3章 引入 $\mathbf{w}$ 的先验概率分布,形式为:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) \quad (37)$$

在输入 $\mathbf{t}$ 下的后验概率分布形式为:

$$p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N) \quad (38)$$

其中,

$$\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \boldsymbol{\Phi}^T \mathbf{t}) \quad (39)$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \boldsymbol{\Phi}_T^T \boldsymbol{\Phi} \quad (40)$$

预测分布:也是求在输入和超参的条件下,预测值的分布

$$p(t | \mathbf{t}, \alpha, \beta) = \int p(t | \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{t}, \alpha, \beta) d\mathbf{w} \quad (41)$$

预测分布的形式为:

$$p(t | \mathbf{t}, \alpha, \beta) = \mathcal{N}(t | \mathbf{m}_N^T \boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x})) \quad (42)$$

其中,

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \boldsymbol{\phi}(\mathbf{x})^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}) \quad (43)$$

因此,贝叶斯线性回归模型是高斯过程的一个具体例子。

### 1.11 参考资料

- PRML 2.3,3.35.7,6.4,8.1
- Gaussian Processes for Machine Learning
- 【机器学习技术】高斯过程初探
- 如何通俗易懂地介绍 Gaussian Process? -知乎
- 贝叶斯线性回归 (Bayesian Linear Regression)
- Machine learning - Introduction to Gaussian processes-UBC