

Sparse Variational Gaussian Processes

目 录

1	Gaussian Distribution	1
1.1	联合概率	1
1.2	边缘概率	1
1.3	条件概率	1
2	Gaussian Processes	2
2.1	先验分布	2
2.2	核函数	2
2.3	似然分布	2
2.4	边缘分布	3
2.5	后验分布	3
2.6	预测应用	4
3	Variational GPs	5
3.1	问题引出	5
3.1.1	先验	5
3.1.2	似然	5
3.1.3	后验	6
3.2	变分近似	6
3.3	ELBO 解析式	7
3.3.1	似然解析式	7
3.3.2	KL 距离解析式	9
3.4	预测应用	10
4	Sparse VGPs	12
4.1	问题引出	12
4.1.1	稀疏先验	12
4.1.2	似然函数	12
4.1.3	后验分布	13
4.2	变分近似	13
4.3	ELBO 解析式	14
4.3.1	似然解析式	14
4.3.2	KL 距离解析式	16
4.4	预测应用	16

第 1 章 Gaussian Distribution

本章主要回顾多元高斯分布及相应边缘概率和条件概率形式

1.1 联合概率

定义 $\mathbf{x} \in \mathbb{R}^d$ 为随机向量，其对应的均值为 $\boldsymbol{\mu}$ ，对应的协方差矩阵为 Σ ，协方差矩阵衡量了随机变量两两间的相关程度。则维度为 d 的随机向量相应的高斯分布概率密度函数如下：

$$p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (1.1)$$

1.2 边缘概率

定义随机向量 \mathbf{x} 和 \mathbf{y} 服从如下联合高斯分布：

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix}\right) \quad (1.2)$$

则随机向量 \mathbf{x} 和 \mathbf{y} 分别服从如下边缘分布：

$$p(\mathbf{x}) \sim \mathcal{N}(\boldsymbol{\mu}_x, \mathbf{A}) \quad (1.3)$$

$$p(\mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}_y, \mathbf{B}) \quad (1.4)$$

1.3 条件概率

定义随机向量 \mathbf{x} 和 \mathbf{y} 服从如下联合高斯分布：

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix}\right) \quad (1.5)$$

则给定随机向量 \mathbf{y} 取值时，随机向量 \mathbf{x} 对应的条件概率服从如下高斯分布：

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_{x|\mathbf{y}}, \Sigma_{x|\mathbf{y}}) \quad (1.6)$$

其中：

$$\boldsymbol{\mu}_{x|\mathbf{y}} = \boldsymbol{\mu}_x + \mathbf{C}\mathbf{B}^{-1}(\mathbf{y} - \boldsymbol{\mu}_y) \quad (1.7)$$

$$\Sigma_{x|\mathbf{y}} = \mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top \quad (1.8)$$

第2章 Gaussian Processes

本章涉及高斯过程先验分布、似然分布及后验分布的推理，最后给出高斯过程如何应用于预测

2.1 先验分布

高斯过程是关于函数值随机变量 $f(\mathbf{x})$ 的分布，分布的定义由均值函数 $m(\mathbf{x})$ 和核函数 $k(\mathbf{x}, \mathbf{x}')$ 给定，其中 \mathbf{x} 和 \mathbf{x}' 为随机变量 $f(\mathbf{x})$ 及 $f(\mathbf{x}')$ 相应的输入。因此对于随机变量 $\mathbf{x} \in \mathbb{R}^d$ 的任意子集构成的矩阵 $\mathbf{X} \in \mathbb{R}^{n \times d}$ ，其对应的随机向量 $f(\mathbf{X})$ 服从如下多元高斯分布：

$$f(\mathbf{X}) \sim \mathcal{N}(m(\mathbf{X}), k(\mathbf{X}, \mathbf{X})) \quad (2.1)$$

其中均值向量为 $\boldsymbol{\mu} = m(\mathbf{X})$ ，协方差矩阵 $\Sigma = k(\mathbf{X}, \mathbf{X})$ 。公式2.1为高斯过程先验 (GP Prior)，应用中一般令均值函数 $m(\mathbf{x}) = 0$ ，协方差矩阵 Σ 由核函数 $k(\mathbf{x}, \mathbf{x}')$ 给定。在高斯过程中称 $f(\mathbf{X})$ 服从的多元高斯分布为高斯过程先验，高斯过程先验即为我们假设随机变量 $f(\mathbf{X})$ 服从的联合分布。

2.2 核函数

核函数有不同的形式，分别对应的不同的随机变量 $f(\mathbf{x})$ 服从的先验分布。这里仅介绍在 EE 问题应用中涉及的径向基核函数 (RBF)，函数定义如下：

$$k_{\text{RBF}}(\mathbf{x}_{i,:}, \mathbf{x}_{j,:}) = \sigma^2 \exp\left(-\frac{1}{2\ell^2} \sum_{q=1}^d (\mathbf{x}_{i,q} - \mathbf{x}_{j,q})^2\right) \quad (2.2)$$

其中 σ^2 为方差， ℓ 为长度缩放系数。

2.3 似然分布

高斯过程先验仅涉及随机变量 \mathbf{x} ，通过似然 (Likelihood) 可以将 GP Prior 中的随机变量 \mathbf{x} 和相应的观察值 y 联系在一起，给定先验和似然后，即可计算后验，有了后验后，即可用于预测。似然为给定随机变量 $f(\mathbf{x})$ 后，观察到 y 的概率。假定给定随机变量 $\mathbf{x} \in \mathbb{R}^d$ 的任意子集构成的矩阵 $\mathbf{X} \in \mathbb{R}^{n \times d}$ 后，随机向量 $y(\mathbf{X})$ 在给定 $f(\mathbf{X})$ 时，服从如下多元高斯分布：

$$p(y(\mathbf{X})|f(\mathbf{X})) = \mathcal{N}(f(\mathbf{X}), \eta^2 \mathbf{I}_n) \quad (2.3)$$

上述形式即为似然函数，其均值为 $f(\mathbf{X})$ ，协方差为 $\eta^2 \mathbf{I}_n$ ，其中 η^2 为模型参数，称为高斯噪声， \mathbf{I}_n 为 $n \times n$ 的单位矩阵。

2.4 边缘分布

似然函数2.3等价定义为如下高斯线性变换：

$$y(\mathbf{X}) = \mathbf{I}_n f(\mathbf{X}) + \epsilon \quad \text{其中: } \epsilon \sim \mathcal{N}(\mathbf{0}, \eta^2 \mathbf{I}_n) \quad (2.4)$$

通过上述形式，可以应用高斯线性变换规则，推理出随机变量 $y(\mathbf{X})$ 的边缘分布。给定随机向量 \mathbf{a} 服从如下多元高斯分布

$$\mathbf{a} \sim \mathcal{N}(\mu_a, \Sigma_a) \quad (2.5)$$

随机向量 \mathbf{b} 为随机向量 \mathbf{a} 的线性变换：

$$\mathbf{b} = \mathbf{A}\mathbf{a} + \eta \quad \text{其中: } \eta \sim \mathcal{N}(0, \Sigma_\eta) \quad (2.6)$$

则随机向量 \mathbf{b} 服从如下多元高斯分布：

$$\mathbf{b} \sim \mathcal{N}(\mathbf{A}\mu_a, \mathbf{A}\Sigma_a\mathbf{A}^\top + \Sigma_\eta) \quad (2.7)$$

给定上述线性变换规则后，由于 $f(\mathbf{X})$ 服从分布2.1，则 $y(\mathbf{X})$ 服从如下分布：

$$y(\mathbf{X}) \sim \mathcal{N}(m(\mathbf{X}), k(\mathbf{X}, \mathbf{X}) + \eta^2 \mathbf{I}_n) \quad (2.8)$$

2.5 后验分布

给定高斯过程先验分布和边缘分布后，可以求得任意输入 \mathbf{x}' 对应的随机变量 $f(\mathbf{x}')$ 的后验分布 $p(f(\mathbf{x}')|y(\mathbf{X}))$ 。为了解该分布，由于随机变量 $f(\mathbf{x}')$ 和随机向量 $y(\mathbf{X})$ 均服从高斯分布，由高斯分布性质可知，其联合分布依然为高斯分布，形式如下：

$$\begin{bmatrix} f(\mathbf{x}') \\ y(\mathbf{X}) \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} m(\mathbf{x}') \\ m(\mathbf{X}) \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}', \mathbf{x}') & \text{cov}(f(\mathbf{x}'), y(\mathbf{X})) \\ \text{cov}(f(\mathbf{x}'), y(\mathbf{X}))^\top & k(\mathbf{X}, \mathbf{X}) + \eta^2 \mathbf{I}_n \end{bmatrix}\right)$$

其中

$$\begin{aligned} \text{cov}(f(\mathbf{x}'), y(\mathbf{X})) &= \mathbb{E}[(f(\mathbf{x}') - m(\mathbf{x}'))(y(\mathbf{X}) - m(\mathbf{X}))] \\ &= \mathbb{E}[(f(\mathbf{x}') - m(\mathbf{x}'))(f(\mathbf{X}) + \epsilon - m(\mathbf{X}))] \\ &= \mathbb{E}[(f(\mathbf{x}') - m(\mathbf{x}'))(f(\mathbf{X}) - m(\mathbf{X})) + (f(\mathbf{x}') - m(\mathbf{x}'))\epsilon] \\ &= \text{cov}(f(\mathbf{x}'), f(\mathbf{X})) + \mathbb{E}[(f(\mathbf{x}') - m(\mathbf{x}'))\epsilon] \\ &= k(\mathbf{x}', \mathbf{X}) + 0 \end{aligned}$$

因此有：

$$\begin{bmatrix} f(\mathbf{x}') \\ y(\mathbf{X}) \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} m(\mathbf{x}') \\ m(\mathbf{X}) \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}', \mathbf{x}') & k(\mathbf{x}', \mathbf{X}) \\ k(\mathbf{x}', \mathbf{X})^\top & k(\mathbf{X}, \mathbf{X}) + \eta^2 \mathbf{I}_n \end{bmatrix} \right)$$

结合上一章的条件概率公式可知

$$p(f(\mathbf{x}')|y(\mathbf{X})) = \mathcal{N}(\mu', \sigma^{2'})$$

其中：

$$\mu' = m(\mathbf{x}') + k(\mathbf{x}', \mathbf{X})(k(\mathbf{X}, \mathbf{X}) + \eta^2 \mathbf{I}_n)^{-1}(y(\mathbf{X}) - m(\mathbf{X})) \quad (2.9)$$

$$\sigma^{2'} = k(\mathbf{x}', \mathbf{x}') - k(\mathbf{x}', \mathbf{X})(k(\mathbf{X}, \mathbf{X}) + \eta^2 \mathbf{I}_n)^{-1}k(\mathbf{x}', \mathbf{X})^\top \quad (2.10)$$

2.6 预测应用

后验分布 $p(f(\mathbf{x}')|y(\mathbf{X}))$ 为给定 $y(\mathbf{X})$ 条件下，隐藏随机变量 $f(\mathbf{x}')$ 的概率分布，由于观测随机变量来自于引入高斯噪声后的隐藏随机变量。因此给定 $y(\mathbf{X})$ 条件时，观测随机变量 $y(\mathbf{x}')$ 满足如下线性变换：

$$y(\mathbf{x}')|y(\mathbf{X}) = f(\mathbf{x}')|y(\mathbf{X}) + \epsilon' \quad \text{其中：} \epsilon' \sim \mathcal{N}(0, \eta^2) \quad (2.11)$$

因此有：

$$p(y(\mathbf{x}')|y(\mathbf{X})) = \mathcal{N}(\mu', \sigma^{2'} + \eta^2) \quad (2.12)$$

第3章 Variational GPs

前面在介绍高斯过程时，假定似然函数服从高斯分布。本章我们来看，当似然函数不服从高斯分布时，如何求解后验分布。

3.1 问题引出

假设我们要解决的是二分类问题，贝叶斯推断中总是包含先验和似然这两项，因此我们分别来看。

3.1.1 先验

对于二分类问题来说，先验项依然为随机变量 $f(\mathbf{X})$ 的联合分布，假设其服从如下高斯分布：

$$f(\mathbf{X}) \sim \mathcal{N}(m(\mathbf{X}), k(\mathbf{X}, \mathbf{X})) \quad (3.1)$$

其中 $m(\mathbf{X})$ 为均值，一般令均值为 0， $k(\mathbf{X}, \mathbf{X})$ 为协方差，可以采用 RBF 核函数，因此其概率密度函数如下所示：

$$p(\mathbf{f}; \ell, \sigma^2) = \frac{1}{\sqrt{(2\pi)^n \det(\mathbf{K})}} \exp\left(-\frac{1}{2} \mathbf{f}^\top \mathbf{K}^{-1} \mathbf{f}\right) \quad (3.2)$$

上述建模过程中未考虑测试噪声，和上一章类似，引入高斯噪声后：

$$f(\mathbf{X}) \sim \mathcal{N}(m(\mathbf{X}), k(\mathbf{X}, \mathbf{X}) + \eta^2 \mathbf{I}_n) \quad (3.3)$$

因此引入高斯噪声后，可以用 $k(\mathbf{X}, \mathbf{X}) + \eta^2 \mathbf{I}_n$ 替代 $k(\mathbf{X}, \mathbf{X})$ 用于有噪声建模。

3.1.2 似然

定义随机变量 y 对应二分类问题中观察到的 label， f 为随机变量 y 对应的潜在随机变量。由于二分类问题服从伯努利分布，其对应似然函数形式如下：

$$p(y_i | f_i) = g(f_i)^{y_i} (1 - g(f_i))^{(1-y_i)} \quad (3.4)$$

其中：

$$g(f_i) = \frac{e^{f_i}}{1 + e^{f_i}} \quad (3.5)$$

对于整个训练集来说有：

$$p(\mathbf{y} | \mathbf{f}) = \prod_{i=1}^n p(y_i | f_i) \quad (3.6)$$

3.1.3 后验

给定先验和似然后，由贝叶斯法则，可以知：

$$p(\mathbf{f}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f})}{\int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f}} \quad (3.7)$$

求解后验分布需要对分母的积分进行简化，求得其相应的解析表达式。由于先验为高斯分布，而似然为伯努利分布，伯努利分布非高斯分布的共轭分布，因此无法求得上述积分形式的解析表达式。

3.2 变分近似

虽然无法直接求得后验分布 $p(\mathbf{f}|\mathbf{y})$ 的解析表达式，但可以通过另外一个分布 $q(\mathbf{f})$ 来近似后验分布，希望该分布尽量和后验分布一致，称该近似分布为变分分布。假设变分分布服从均值为 \mathbf{v} ，协方差矩阵为 \mathbf{S} 的高斯分布，则有：

$$q(\mathbf{f}) = \mathcal{N}(\mathbf{v}, \mathbf{S}) \quad (3.8)$$

给定变分分布后，可以通过 KL 距离衡量变分分布和后验分布近似程度：

$$\begin{aligned} KL(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) &= \mathbb{E}_{\mathbf{f} \sim q} \left[\log \left(\frac{q(\mathbf{f})}{p(\mathbf{f}|\mathbf{y})} \right) \right] \\ &= \mathbb{E}_{\mathbf{f} \sim q} [\log(q(\mathbf{f}))] - \mathbb{E}_{\mathbf{f} \sim q} [\log(p(\mathbf{f}|\mathbf{y}))] \\ &= \mathbb{E}_{\mathbf{f} \sim q} [\log(q(\mathbf{f}))] - \mathbb{E}_{\mathbf{f} \sim q} \left[\log \left(\frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f})}{p(\mathbf{y})} \right) \right] \\ &= \mathbb{E}_{\mathbf{f} \sim q} [\log(q(\mathbf{f}))] - \mathbb{E}_{\mathbf{f} \sim q} [\log(p(\mathbf{y}|\mathbf{f})p(\mathbf{f}))] + \mathbb{E}_{\mathbf{f} \sim q} [\log p(\mathbf{y})] \\ &= \mathbb{E}_{\mathbf{f} \sim q} [\log(q(\mathbf{f}))] - \mathbb{E}_{\mathbf{f} \sim q} [\log(p(\mathbf{y}|\mathbf{f})p(\mathbf{f}))] + \log p(\mathbf{y}) \end{aligned} \quad (3.9)$$

令

$$\text{ELBO} = -(\mathbb{E}_{\mathbf{f} \sim q} [\log(q(\mathbf{f}))] - \mathbb{E}_{\mathbf{f} \sim q} [\log(p(\mathbf{y}|\mathbf{f})p(\mathbf{f}))]) \quad (3.10)$$

因此有：

$$KL(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y})) = \log p(\mathbf{y}) - \text{ELBO} \quad (3.11)$$

有上述形式可知，由于 $KL(q(\mathbf{f})||p(\mathbf{f}|\mathbf{y}))$ 大于等于 0， $\log p(\mathbf{y})$ 和参数无关，因此要最小化变分分布和后验分布的 KL 距离，可以最大化 ELBO。

3.3 ELBO 解析式

由上一小节可知，ELBO 公式形式如下：

$$\begin{aligned}
 \text{ELBO} &= -(\mathbb{E}_{\mathbf{f} \sim q} [\log(q(\mathbf{f}))] - \mathbb{E}_{\mathbf{f} \sim q} [\log(p(\mathbf{y}|\mathbf{f})p(\mathbf{f}))]) \\
 &= -\mathbb{E}_{\mathbf{f} \sim q} [\log(q(\mathbf{f}))] + \mathbb{E}_{\mathbf{f} \sim q} [\log(p(\mathbf{y}|\mathbf{f}))] + \mathbb{E}_{\mathbf{f} \sim q} [\log(p(\mathbf{f}))] \\
 &= \mathbb{E}_{\mathbf{f} \sim q} [\log(p(\mathbf{y}|\mathbf{f}))] - \mathbb{E}_{\mathbf{f} \sim q} \left[\frac{q(\mathbf{f})}{p(\mathbf{f})} \right] \\
 &= \mathbb{E}_{\mathbf{f} \sim q} [\log(p(\mathbf{y}|\mathbf{f}))] - KL(q(\mathbf{f})||p(\mathbf{f})) \\
 &= \int \log(p(\mathbf{y}|\mathbf{f}))q(\mathbf{f})d\mathbf{f} - KL(q(\mathbf{f})||p(\mathbf{f})) \tag{3.12}
 \end{aligned}$$

可见 ELBO 形式中的似然项和 KL 距离项均为积分形式，求得上述两项的解析式后，即可通过梯度上升方法求解使得 ELBO 取得最大值的参数。

3.3.1 似然解析式

给定 ELBO 形式后，我们首先看如何求得似然项的解析式

$$\begin{aligned}
 \int \log(p(\mathbf{y}|\mathbf{f}))q(\mathbf{f})d\mathbf{f} &= \int \log\left(\prod_{i=1}^n p(y_i|f_i)\right)q(\mathbf{f})d\mathbf{f} \\
 &= \int \left(\sum_{i=1}^n \log(p(y_i|f_i))\right)q(\mathbf{f})d\mathbf{f} \\
 &= \sum_{i=1}^n \left(\int \log(p(y_i|f_i))q(\mathbf{f})d\mathbf{f}\right) \tag{3.13}
 \end{aligned}$$

对上述求和形式中每一项的积分形式进行化简：

$$\begin{aligned}
 \int \log(p(y_i|f_i))q(\mathbf{f})d\mathbf{f} &= \int \cdots \int \log(p(y_i|f_i))q([f_1, \dots, f_n]^\top)df_1 \dots df_n \\
 &= \iint \log(p(y_i|f_i))q([f_{-i}, f_i]^\top)df_{-i}df_i \\
 &= \int \log(p(y_i|f_i)) \left(\int q([f_{-i}, f_i]^\top)df_{-i} \right) df_i \\
 &= \int \log(p(y_i|f_i))q(f_i)df_i \tag{3.14}
 \end{aligned}$$

上述公式中 f_{-i} 表示向量 \mathbf{f} 中除 f_i 的所有元素。由多元高斯分布的性质可知，变分分布 $q(\mathbf{f})$ 的边缘分布为：

$$\begin{aligned}
 q(f_i) &= \mathcal{N}(f_i; \mu_i, \Sigma_{i,i}^2) \\
 &= \frac{1}{\sqrt{2\pi\Sigma_{i,i}^2}} \exp\left(-\frac{(f_i - \mu_i)^2}{2\Sigma_{i,i}^2}\right) \tag{3.15}
 \end{aligned}$$

因此有：

$$\int \log(p(y_i|f_i))q(f_i)df_i = \int \log(p(y_i|f_i))\frac{1}{\sqrt{2\pi\Sigma_{i,i}^2}}\exp\left(-\frac{(f_i-\mu_i)^2}{2\Sigma_{i,i}^2}\right)df_i \quad (3.16)$$

可以通过 Gauss-Hermit quadrature 求得上述积分形式的数值近似解析式，其定义如下：

$$\int_{-\infty}^{+\infty} e^{-x^2} h(x) dx \approx \sum_{i=1}^n w_i h(x_i) \quad (3.17)$$

为了应用上述形式求解似然项的积分近似解析式，定义随机变量 f_i 满足如下形式：

$$f_i = \Sigma_{i,i}u + \mu_i \quad (3.18)$$

由高斯线性变换可知，随机变量 u 服从均值为 0，方差为 1 的高斯分布，具体推理如下：

$$\begin{aligned} f_i &\sim \mathcal{N}(\mu_i + \Sigma_{i,i} \cdot 0, \Sigma_{i,i}1\Sigma_{i,i}) \\ &\sim \mathcal{N}(\mu_i, \Sigma_{i,i}^2) \end{aligned} \quad (3.19)$$

给定 f_i 的变换形式后，似然项的积分形式可以简化为如下形式：

$$\begin{aligned} &\int \log(p(y_i|f_i))q(f_i)df_i \\ &= \int \log(p(y_i|f_i))\frac{1}{\sqrt{2\pi\Sigma_{i,i}^2}}\exp\left(-\frac{(f_i-\mu_i)^2}{2\Sigma_{i,i}^2}\right)df_i \\ &= \int \log(p(y_i|\Sigma_{i,i}u + \mu_i))\frac{1}{\sqrt{2\pi\Sigma_{i,i}^2}}\exp\left(-\frac{(\Sigma_{i,i}u + \mu_i - \mu_i)^2}{2\Sigma_{i,i}^2}\right)d(\Sigma_{i,i}u + \mu_i) \\ &= \int \log(p(y_i|\Sigma_{i,i}u + \mu_i))\frac{1}{\sqrt{2\pi\Sigma_{i,i}^2}}\exp\left(-\frac{1}{2}u^2\right) \cdot \Sigma_{i,i}du \\ &= \int \log(p(y_i|\Sigma_{i,i}u + \mu_i))\frac{1}{\sqrt{2\pi}}\exp\left(-\frac{1}{2}u^2\right)du \end{aligned} \quad (3.20)$$

令 $t = 1/\sqrt{2}u$ ，则 $u = \sqrt{2}t$ ，因此有：

$$\begin{aligned} &\int \log(p(y_i|\Sigma_{i,i}u + \mu_i))\frac{1}{\sqrt{2\pi}}\exp\left(-\frac{1}{2}u^2\right) \cdot du \\ &= \int \log(p(y_i|\sqrt{2}\Sigma_{i,i}t + \mu_i))\frac{1}{\sqrt{2\pi}}\exp(-t^2) \cdot \sqrt{2}dt \\ &= \int \log(p(y_i|\sqrt{2}\Sigma_{i,i}t + \mu_i))\frac{1}{\sqrt{\pi}}\exp(-t^2)dt \end{aligned} \quad (3.21)$$

因此令

$$h(t) = \frac{1}{\sqrt{\pi}}\log(p(y_i|\sqrt{2}\Sigma_{i,i}t + \mu_i)) \quad (3.22)$$

其中 $p(y_i|\sqrt{2}\Sigma_{i,i}t + \mu_i)$ 服从伯努利分布，因此有：

$$\begin{aligned} & \int \log(p(y_i|\sqrt{2}\Sigma_{i,i}t + \mu_i)) \frac{1}{\sqrt{\pi}} \exp(-t^2) dt \\ & \approx \sum_{j=1}^m w_j \cdot \frac{1}{\sqrt{\pi}} \log \left(\left(\frac{e^{(\sqrt{2}\Sigma_{i,i}t_j + \mu_i)}}{1 + e^{(\sqrt{2}\Sigma_{i,i}t_j + \mu_i)}} \right)^{y_i} \left(1 - \frac{e^{(\sqrt{2}\Sigma_{i,i}t_j + \mu_i)}}{1 + e^{(\sqrt{2}\Sigma_{i,i}t_j + \mu_i)}} \right)^{(1-y_i)} \right) \end{aligned} \quad (3.23)$$

上述形式即为似然项积分形式对应的近似解析式

3.3.2 KL 距离解析式

由上一小节可知，KL 距离项 $KL(q(f)||p(f))$ ，其中 $q(f)$ 和 $p(f)$ 分别服从如下分布：

$$\begin{aligned} q(f) &= N(\mathbf{v}, \mathbf{S}) = \frac{1}{\sqrt{(2\pi)^d |\mathbf{S}|}} \exp \left(-\frac{1}{2} (\mathbf{f} - \mathbf{v})^\top \mathbf{S}^{-1} (\mathbf{f} - \mathbf{v}) \right) \\ p(f) &= N(\mathbf{m}, \mathbf{K}) = \frac{1}{\sqrt{(2\pi)^d |\mathbf{K}|}} \exp \left(-\frac{1}{2} (\mathbf{f} - \mathbf{m})^\top \mathbf{K}^{-1} (\mathbf{f} - \mathbf{m}) \right) \end{aligned}$$

因此有：

$$\begin{aligned} KL(q(f)||p(f)) &= \int q(f) \log \left(\frac{q(f)}{p(f)} \right) d\mathbf{f} \\ &= \int q(f) (\log(q(f)) - \log(p(f))) d\mathbf{f} \end{aligned}$$

其中：

$$\begin{aligned} \log(q(f)) &= \log \left(\frac{1}{\sqrt{(2\pi)^d |\mathbf{S}|}} \exp \left(-\frac{1}{2} (\mathbf{f} - \mathbf{v})^\top \mathbf{S}^{-1} (\mathbf{f} - \mathbf{v}) \right) \right) \\ &= \log \left(\frac{1}{\sqrt{(2\pi)^d |\mathbf{S}|}} \right) - \frac{1}{2} (\mathbf{f} - \mathbf{v})^\top \mathbf{S}^{-1} (\mathbf{f} - \mathbf{v}) \\ &= -\frac{1}{2} \log((2\pi)^d) - \frac{1}{2} \log(|\mathbf{S}|) - \frac{1}{2} (\mathbf{f} - \mathbf{v})^\top \mathbf{S}^{-1} (\mathbf{f} - \mathbf{v}) \end{aligned}$$

$$\begin{aligned} \log(p(f)) &= \log \left(\frac{1}{\sqrt{(2\pi)^d |\mathbf{K}|}} \exp \left(-\frac{1}{2} (\mathbf{f} - \mathbf{m})^\top \mathbf{K}^{-1} (\mathbf{f} - \mathbf{m}) \right) \right) \\ &= \log \left(\frac{1}{\sqrt{(2\pi)^d |\mathbf{K}|}} \right) - \frac{1}{2} (\mathbf{f} - \mathbf{m})^\top \mathbf{K}^{-1} (\mathbf{f} - \mathbf{m}) \\ &= -\frac{1}{2} \log((2\pi)^d) - \frac{1}{2} \log(|\mathbf{K}|) - \frac{1}{2} (\mathbf{f} - \mathbf{m})^\top \mathbf{K}^{-1} (\mathbf{f} - \mathbf{m}) \end{aligned}$$

结合上式可知：

$$\begin{aligned} & \log(q(\mathbf{f})) - \log(p(\mathbf{f})) \\ &= \frac{1}{2} \log\left(\frac{|\mathbf{K}|}{|\mathbf{S}|}\right) - \frac{1}{2}(\mathbf{f} - \mathbf{v})^\top \mathbf{S}^{-1}(\mathbf{f} - \mathbf{v}) + \frac{1}{2}(\mathbf{f} - \mathbf{m})^\top \mathbf{K}^{-1}(\mathbf{f} - \mathbf{m}) \end{aligned}$$

结合 Matrix Cookbook 第 8.2 节的性质：假设 $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \Sigma)$ ，则：

$$\mathbb{E}[(\mathbf{x} - \mathbf{m}')^\top \mathbf{A}(\mathbf{x} - \mathbf{m}')] = (\mathbf{m} - \mathbf{m}')^\top \mathbf{A}(\mathbf{m} - \mathbf{m}') + \text{Tr}(\mathbf{A}\Sigma) \quad (3.24)$$

因此有：

$$\begin{aligned} & KL(q(\mathbf{f})||p(\mathbf{f})) \\ &= \int q(\mathbf{f}) \left(\frac{1}{2} \log\left(\frac{|\mathbf{K}|}{|\mathbf{S}|}\right) - \frac{1}{2}(\mathbf{f} - \mathbf{v})^\top \mathbf{S}^{-1}(\mathbf{f} - \mathbf{v}) + \frac{1}{2}(\mathbf{f} - \mathbf{m})^\top \mathbf{K}^{-1}(\mathbf{f} - \mathbf{m}) \right) d\mathbf{f} \\ &= \frac{1}{2} \log\left(\frac{|\mathbf{K}|}{|\mathbf{S}|}\right) + \mathbb{E}_{\mathbf{f} \sim q} \left[-\frac{1}{2}(\mathbf{f} - \mathbf{v})^\top \mathbf{S}^{-1}(\mathbf{f} - \mathbf{v}) + \frac{1}{2}(\mathbf{f} - \mathbf{m})^\top \mathbf{K}^{-1}(\mathbf{f} - \mathbf{m}) \right] \\ &= \frac{1}{2} \log\left(\frac{|\mathbf{K}|}{|\mathbf{S}|}\right) + \frac{1}{2} \left((\mathbf{v} - \mathbf{m})^\top \mathbf{K}^{-1}(\mathbf{v} - \mathbf{m}) + \text{Tr}(\mathbf{K}^{-1}\mathbf{S}) \right) - \frac{1}{2} \text{Tr}(\mathbf{S}^{-1}\mathbf{S}) \\ &= \frac{1}{2} \left(\log\left(\frac{|\mathbf{K}|}{|\mathbf{S}|}\right) + (\mathbf{v} - \mathbf{m})^\top \mathbf{K}^{-1}(\mathbf{v} - \mathbf{m}) + \text{Tr}(\mathbf{K}^{-1}\mathbf{S}) - n \right) \end{aligned} \quad (3.25)$$

上式即为 KL 距离的解析式。

3.4 预测应用

求得后验分布 $p(\mathbf{f}(\mathbf{X})|\mathbf{y}(\mathbf{X}))$ 的变分分布 $q(\mathbf{f}(\mathbf{X}))$ 近似后，我们即可应用该分布求解函数在 \mathbf{x}' 处的函数分布，令随机变量 f' 表示 $f(\mathbf{x}')$ ，随机向量 \mathbf{f} 表示 $\mathbf{f}(\mathbf{X})$ ，随机向量 \mathbf{y} 表示 $\mathbf{y}(\mathbf{X})$ ，则有：

$$\begin{aligned} p(f(\mathbf{x}')|\mathbf{y}(\mathbf{X})) &= \int p(f(\mathbf{x}'), \mathbf{f}(\mathbf{X})|\mathbf{y}(\mathbf{X})) d\mathbf{f}(\mathbf{X}) = \int p(f', \mathbf{f}|\mathbf{y}) d\mathbf{f} \\ &= \int p(f'|\mathbf{f}, \mathbf{y}) p(\mathbf{f}|\mathbf{y}) d\mathbf{f} = \int p(f'|\mathbf{f}) p(\mathbf{f}|\mathbf{y}) d\mathbf{f} \\ &= \int p(f'|\mathbf{f}) q(\mathbf{f}) d\mathbf{f} \end{aligned}$$

由多元高斯分布条件概率公式可知：

$$p(f'|\mathbf{f}) = \mathcal{N}(\mu', \sigma'^2)$$

其中：

$$\begin{aligned}\mu' &= \mathbf{m} + k(\mathbf{x}', \mathbf{X})k(\mathbf{X}, \mathbf{X})^{-1}(\mathbf{f} - \mathbf{m}) \\ \sigma^{2'} &= k(\mathbf{x}', \mathbf{x}') - k(\mathbf{x}', \mathbf{X})k(\mathbf{X}, \mathbf{X})^{-1}k(\mathbf{x}', \mathbf{X})^\top\end{aligned}$$

由于现实应用中一般 \mathbf{m} 为 $\mathbf{0}$ ，因此有：

$$\begin{aligned}\mu' &= \underbrace{k(\mathbf{x}', \mathbf{X})k(\mathbf{X}, \mathbf{X})^{-1}}_{\mathbf{a}^\top} \mathbf{f} \\ \sigma^{2'} &= \underbrace{k(\mathbf{x}', \mathbf{x}') - k(\mathbf{x}', \mathbf{X})k(\mathbf{X}, \mathbf{X})^{-1}k(\mathbf{x}', \mathbf{X})^\top}_{b} \\ f'|f &\sim \mathcal{N}(\mathbf{a}^\top \mathbf{f}, b)\end{aligned}$$

由于随机向量 \mathbf{f} 服从变分分布 $\mathbf{f} \sim \mathcal{N}(\mathbf{v}, \mathbf{S})$ ，根据高斯线性变换可知：

$$f'|f \sim \mathcal{N}(\mathbf{a}^\top \mathbf{v}, \mathbf{a}^\top \mathbf{S} \mathbf{a} + b) \quad (3.26)$$

因此有：

$$\begin{aligned}p(f'|f) &= \mathcal{N}(\mathbf{a}^\top \mathbf{v}, \mathbf{a}^\top \mathbf{S} \mathbf{a} + b) \\ p(f(\mathbf{x}')|y(\mathbf{X})) &= \int p(f'|f)q(\mathbf{f})d\mathbf{f} \\ &= \mathcal{N}(\mathbf{a}^\top \mathbf{v}, \mathbf{a}^\top \mathbf{S} \mathbf{a} + b) \int q(\mathbf{f})d\mathbf{f} \\ &= \mathcal{N}(\mathbf{a}^\top \mathbf{v}, \mathbf{a}^\top \mathbf{S} \mathbf{a} + b)\end{aligned} \quad (3.27)$$

上式即为给定 \mathbf{x}' 时其取值函数对应的分布，这里未考虑高斯噪声。

第4章 Sparse VGPs

从上一章节可知通过变分高斯过程，我们可以在预测阶段给定任意样本点时，给出其对应函数值随机变量服从的高斯分布密度函数。通过分析可见，在求解过程中需要构建训练数据集上所有样本点对应的协方差矩阵 $k(\mathbf{X}, \mathbf{X})$ 的逆，该计算的时间复杂度为 $O(n^3)$ ，可见当 n 较大时，时间复杂度较高。针对上述问题稀疏变分高斯过程在前述方案基础上通过引入 $m(m \ll n)$ 个稀疏归纳点，有效解决了变分高斯过程存在的时间复杂度和空间复杂度较高的问题。

4.1 问题引出

由于协方差矩阵计算复杂度较高，针对该问题，我们定义 m 个稀疏点，所有稀疏点构成矩阵 $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_M]^\top \in \mathbb{R}^{m \times d}$ ，令随机向量 $f(\mathbf{Z})$ 服从如下高斯分布：

$$f(\mathbf{Z}) \sim \mathcal{N}(m(\mathbf{Z}), k(\mathbf{Z}, \mathbf{Z})) \quad (4.1)$$

本节介绍引入 m 个稀疏点后，如何定义问题求解所需的高斯过程先验，似然及后验分布。上述建模过程中未考虑测试噪声，和前述章节类似，引入高斯噪声后：

$$f(\mathbf{Z}) \sim \mathcal{N}(m(\mathbf{Z}), k(\mathbf{Z}, \mathbf{Z}) + \eta^2 \mathbf{I}_n) \quad (4.2)$$

因此引入高斯噪声后，可以用 $k(\mathbf{Z}, \mathbf{Z}) + \eta^2 \mathbf{I}_n$ 替代 $k(\mathbf{Z}, \mathbf{Z})$ 用于有噪声建模。

4.1.1 稀疏先验

假设训练数据和稀疏归纳点服从如下联合高斯分布：

$$\begin{bmatrix} f(\mathbf{X}) \\ f(\mathbf{Z}) \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} m(\mathbf{X}) \\ m(\mathbf{Z}) \end{bmatrix}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) & k(\mathbf{X}, \mathbf{Z}) \\ k(\mathbf{X}, \mathbf{Z})^\top & k(\mathbf{Z}, \mathbf{Z}) \end{bmatrix}\right) \quad (4.3)$$

4.1.2 似然函数

和前序章节类似，通过似然建立先验和训练数据之间的关系，因此有：

$$p(y(\mathbf{X})|f(\mathbf{X}), f(\mathbf{Z})) = p(y(\mathbf{X})|f(\mathbf{X})) \quad (4.4)$$

该分布对应的具体形式和应用相关，对于二分类问题，其概率密度函数为伯努利分布，因此有：

$$p(y_i|f_i) = g(f_i)^{y_i} (1 - g(f_i))^{(1-y_i)} \quad (4.5)$$

其中：

$$g(f_i) = \frac{e^{f_i}}{1 + e^{f_i}} \quad (4.6)$$

对于整个训练集来说有：

$$p(y(\mathbf{X})|f(\mathbf{X})) = \prod_{i=1}^n p(y_i|f_i) \quad (4.7)$$

对于回归问题，一般采用高斯分布，同时假设观察值来自于引入高斯噪声后的真实值，因此有：

$$\begin{aligned} p(y(\mathbf{X})|f(\mathbf{X})) &= \mathcal{N}(f(\mathbf{X}), \eta^2 \mathbf{I}_n) \\ &= \mathcal{N}(m(\mathbf{X}), k(\mathbf{X}, \mathbf{X}) + \eta^2 \mathbf{I}_n) \end{aligned} \quad (4.8)$$

其中 η^2 为各向同性的噪声方差。

4.1.3 后验分布

后验分布为给定训练集后，求如下概率密度函数：

$$p(f(\mathbf{X}), f(\mathbf{Z})|y(\mathbf{X})) = \frac{p(y(\mathbf{X})|f(\mathbf{X}), f(\mathbf{Z}))p(f(\mathbf{X}), f(\mathbf{Z}))}{\iint p(y(\mathbf{X})|f(\mathbf{X}), f(\mathbf{Z}))p(f(\mathbf{X}), f(\mathbf{Z}))df(\mathbf{X})df(\mathbf{Z})} \quad (4.9)$$

令符号 f 表示随机变量 $f(\mathbf{X})$ ， f_Z 表示随机变量 $f(\mathbf{Z})$ ， y 表示随机变量 $y(\mathbf{X})$ ，因此有：

$$p(f, f_Z|y) = \frac{p(y|f, f_Z)p(f, f_Z)}{\iint p(y|f, f_Z)p(f, f_Z)dfdf_Z} \quad (4.10)$$

由于稀疏先验为高斯分布，若问题为分类问题，则似然函数非高斯先验。因此无法求得上述后验分布的解析式。

4.2 变分近似

由于无法直接求得后验分布，因此考虑采用变分分布 $q(f, f_Z)$ 近似后验分布 $p(f, f_Z|y)$ ，并定义：

$$q(f, f_Z) = p(f|f_Z)q(f_Z) \quad (4.11)$$

其中 $q(f_Z)$ 服从如下高斯分布：

$$q(f_Z) = \mathcal{N}(v, S) \quad (4.12)$$

由高斯过程先验结合多元高斯分布条件概率可知：

$$p(f|f_Z) = \mathcal{N}(\mu, \Sigma) \quad (4.13)$$

其中：

$$\mu = m(\mathbf{X}) + k(\mathbf{X}, \mathbf{Z})k(\mathbf{Z}, \mathbf{Z})^{-1}(\mathbf{f}_Z - m(\mathbf{Z})) \quad (4.14)$$

$$\Sigma = k(\mathbf{X}, \mathbf{X}) - k(\mathbf{X}, \mathbf{Z})k(\mathbf{Z}, \mathbf{Z})^{-1}k(\mathbf{X}, \mathbf{Z})^\top \quad (4.15)$$

给定变分分布后，可通过 KL 距离度量变分分布和后验分布的近似程度：

$$\begin{aligned} & KL(q(\mathbf{f}, \mathbf{f}_Z) \| p(\mathbf{f}, \mathbf{f}_Z | \mathbf{y})) \\ &= \iint q(\mathbf{f}, \mathbf{f}_Z) \log \frac{q(\mathbf{f}, \mathbf{f}_Z)}{p(\mathbf{f}, \mathbf{f}_Z | \mathbf{y})} d\mathbf{f} d\mathbf{f}_Z \\ &= \iint q(\mathbf{f}, \mathbf{f}_Z) \log \frac{q(\mathbf{f}, \mathbf{f}_Z)p(\mathbf{y})}{p(\mathbf{y} | \mathbf{f}, \mathbf{f}_Z)p(\mathbf{f}, \mathbf{f}_Z)} d\mathbf{f} d\mathbf{f}_Z \\ &= \iint q(\mathbf{f}, \mathbf{f}_Z) \log(p(\mathbf{y})) d\mathbf{f} d\mathbf{f}_Z + \iint q(\mathbf{f}, \mathbf{f}_Z) \log \frac{q(\mathbf{f}, \mathbf{f}_Z)}{p(\mathbf{y} | \mathbf{f}, \mathbf{f}_Z)p(\mathbf{f}, \mathbf{f}_Z)} d\mathbf{f} d\mathbf{f}_Z \\ &= \log(p(\mathbf{y})) - \iint q(\mathbf{f}, \mathbf{f}_Z) \log \frac{p(\mathbf{y} | \mathbf{f}, \mathbf{f}_Z)p(\mathbf{f}, \mathbf{f}_Z)}{q(\mathbf{f}, \mathbf{f}_Z)} d\mathbf{f} d\mathbf{f}_Z \\ &= \log(p(\mathbf{y})) - \iint q(\mathbf{f}, \mathbf{f}_Z) \log(p(\mathbf{y} | \mathbf{f}, \mathbf{f}_Z)) d\mathbf{f} d\mathbf{f}_Z - \iint q(\mathbf{f}, \mathbf{f}_Z) \log \frac{p(\mathbf{f}, \mathbf{f}_Z)}{q(\mathbf{f}, \mathbf{f}_Z)} d\mathbf{f} d\mathbf{f}_Z \\ &= \log(p(\mathbf{y})) - \iint q(\mathbf{f}, \mathbf{f}_Z) \log(p(\mathbf{y} | \mathbf{f})) d\mathbf{f} d\mathbf{f}_Z + KL(q(\mathbf{f}, \mathbf{f}_Z) \| p(\mathbf{f}, \mathbf{f}_Z)) \end{aligned}$$

令

$$\text{ELBO} = \iint q(\mathbf{f}, \mathbf{f}_Z) \log(p(\mathbf{y} | \mathbf{f})) d\mathbf{f} d\mathbf{f}_Z - KL(q(\mathbf{f}, \mathbf{f}_Z) \| p(\mathbf{f}, \mathbf{f}_Z)) \quad (4.16)$$

因此有：

$$KL(q(\mathbf{f}, \mathbf{f}_Z) \| p(\mathbf{f}, \mathbf{f}_Z | \mathbf{y})) = \log(p(\mathbf{y})) - \text{ELBO} \quad (4.17)$$

因此要最小化变分分布和后验分布的差异，需要最大化 ELBO。

4.3 ELBO 解析式

由 ELBO 解析式 4.16 可知，其包含两项分别为似然项和 KL 距离项，下面分别求取两项的解析式

4.3.1 似然解析式

似然项为二重积分形式，下面给出其解析式求解过程：

$$\begin{aligned} & \iint q(\mathbf{f}, \mathbf{f}_Z) \log(p(\mathbf{y} | \mathbf{f})) d\mathbf{f} d\mathbf{f}_Z \\ &= \int \log(p(\mathbf{y} | \mathbf{f})) \left(\int q(\mathbf{f}, \mathbf{f}_Z) d\mathbf{f}_Z \right) d\mathbf{f} \\ &= \int \log(p(\mathbf{y} | \mathbf{f})) q(\mathbf{f}) d\mathbf{f} \end{aligned} \quad (4.18)$$

为了求解上述积分形式的解析式，需要求得 $q(f)$ 和 $p(y|f)$ 的解析式：

$$\begin{aligned} q(f) &= \int q(f, f_Z) df_Z = \int p(f|f_Z) q(f_Z) df_Z \\ &= \int \mathcal{N}(f; \mu, \Sigma) \cdot \mathcal{N}(f_Z; \nu, S) df_Z \end{aligned} \quad (4.19)$$

其中：

$$\begin{aligned} \mu &= m(X) + k(X, Z)k(Z, Z)^{-1}(f_Z - m(Z)) \\ &= m(X) + k(X, Z)k(Z, Z)^{-1}f_Z - k(X, Z)k(Z, Z)^{-1}m(Z) \end{aligned} \quad (4.20)$$

由于：

$$f_Z \sim \mathcal{N}(\nu, S) \quad (4.21)$$

由高斯线性变换可知：

$$\begin{aligned} \mathcal{N}(f; \mu, \Sigma) &= \mathcal{N}(\mu_{qf}, \Sigma_{qf}) \\ \mu_{qf} &= m(X) + k(X, Z)k(Z, Z)^{-1}\nu - k(X, Z)k(Z, Z)^{-1}m(Z) \\ &= m(X) + k(X, Z)k(Z, Z)^{-1}(\nu - m(Z)) \\ \Sigma_{qf} &= k(X, Z)k(Z, Z)^{-1}S(k(X, Z)k(Z, Z)^{-1})^\top + \Sigma \end{aligned} \quad (4.22)$$

其中 Σ 为公式4.15，上述形式不包含随机变量 f_Z ，因此有：

$$\begin{aligned} q(f) &= \int \mathcal{N}(f; \mu, \Sigma) \cdot \mathcal{N}(f_Z; \nu, S) df_Z \\ &= \mathcal{N}(f; \mu_{qf}, \Sigma_{qf}) \int \mathcal{N}(f_Z; \nu, S) df_Z \\ &= \mathcal{N}(f; \mu_{qf}, \Sigma_{qf}) \end{aligned} \quad (4.23)$$

给定上述形式后，似然项为：

$$\begin{aligned} &\int \log(p(y|f)) q(f) df \\ &= \sum_{i=1}^n \left(\int \log(p(y_i|f_i)) q(f_i) df_i \right) \end{aligned} \quad (4.24)$$

由于：

$$q(f_i) = \mathcal{N}(f_i; \mu_i, \Sigma_{ii}^2) \quad (4.25)$$

因此若求得上述一元高斯分布中的 μ_i 和 Σ_{ii}^2 则后面的近似公式和变分高斯过程中求对数似然的形式保持一致，因此这里重点给出 μ_i 和 Σ_{ii}^2 的公式，由 $q(f)$ 分布公式可以较容易

推断出 μ_i 和 Σ_{ii}^2 形式，并最终得出似然项的数值积分近似解析式。

4.3.2 KL 距离解析式

下面我们来看如何求得 KL 距离项解析式

$$\begin{aligned}
 KL(q(f, f_Z) || p(f, f_Z)) &= \iint q(f, f_Z) \log \left(\frac{q(f, f_Z)}{p(f, f_Z)} \right) df df_Z \\
 &= \iint q(f, f_Z) \log \left(\frac{p(f|f_Z)q(f_Z)}{p(f|f_Z)p(f_Z)} \right) df df_Z \\
 &= \iint \log \left(\frac{q(f_Z)}{p(f_Z)} \right) q(f, f_Z) df df_Z \\
 &= \int \log \left(\frac{q(f_Z)}{p(f_Z)} \right) \left(\int q(f, f_Z) df \right) df_Z \\
 &= \int q(f_Z) \log \left(\frac{q(f_Z)}{p(f_Z)} \right) df_Z
 \end{aligned} \tag{4.26}$$

上式和上一章求取 KL 距离项的形式完全一致，因此其最终的解析式如下：

$$KL(q(f, f_Z) || p(f, f_Z)) = \frac{1}{2} \left(\log \left(\frac{|\mathbf{K}|}{|\mathbf{S}|} \right) + (\mathbf{v} - \mathbf{m})^\top \mathbf{K}^{-1} (\mathbf{v} - \mathbf{m}) + \text{Tr}(\mathbf{K}^{-1} \mathbf{S}) - m \right) \tag{4.27}$$

4.4 预测应用

给定稀疏点对应的高斯过程先验及后验分布的变分近似后，我们即可以进行测试样本点对应的概率密度函数预测。

$$\begin{aligned}
 p(f'|y) &= \iint p(f', f, f_Z | y) df df_Z \\
 &= \iint p(f' | f, f_Z, y) p(f, f_Z | y) df df_Z \\
 &= \iint p(f' | f, f_Z) p(f, f_Z | y) df df_Z \\
 &= \iint p(f' | f, f_Z) q(f, f_Z) df df_Z \\
 &= \iint p(f' | f, f_Z) p(f | f_Z) q(f_Z) df df_Z \\
 &= \int \left(\int p(f' | f, f_Z) p(f | f_Z) df \right) q(f_Z) df_Z \\
 &= \int \left(\int p(f', f | f_Z) df \right) q(f_Z) df_Z \\
 &= \int p(f' | f_Z) q(f_Z) df_Z
 \end{aligned} \tag{4.28}$$

从如上形式可见，最终进行预测时只和稀疏归纳点及变分分布有关，且有：

$$\begin{bmatrix} f' \\ f_Z \end{bmatrix} = \mathcal{N} \left(\begin{bmatrix} m(f') \\ m(f_Z) \end{bmatrix}, \begin{bmatrix} k(X', X') & k(X', Z) \\ k(X', Z)^\top & k(Z, Z) \end{bmatrix} \right)$$

因此有：

$$p(f'|f_Z) = \mathcal{N}(\mu_{f'}, \Sigma_{f'}) \quad (4.29)$$

其中：

$$\mu_{f'} = m(X') + k(X', Z)k(Z, Z)^{-1}(f_Z - m(Z)) \quad (4.30)$$

$$\Sigma_{f'} = k(X', X') - k(X', Z)k(Z, Z)^{-1}k(X', Z)^\top \quad (4.31)$$

由于 f_Z 服从高斯分布 $\mathcal{N}(v, S)$ ，由高斯线性变换可知：

$$p(f'|f_Z) = \mathcal{N}(\mu_p, \Sigma_p) \quad (4.32)$$

其中：

$$\mu_p = m(X') + k(X', Z)k(Z, Z)^{-1}(v - m(Z)) \quad (4.33)$$

$$\begin{aligned} \Sigma_p &= k(X', Z)k(Z, Z)^{-1}S(k(X', Z)k(Z, Z)^{-1})^\top + \Sigma_{f'} \\ &= k(X', Z)k(Z, Z)^{-1}Sk(Z, Z)^{-1}k(X', Z)^\top + \Sigma_{f'} \\ &= k(X', X') - k(X', Z)k(Z, Z)^{-1}(k(Z, Z) - S)k(Z, Z)^{-1}k(X', Z)^\top \end{aligned} \quad (4.34)$$

从上述可见 $p(f'|f_Z)$ 和随机向量 f_Z 无关，因此有：

$$\begin{aligned} p(f'|y) &= \int p(f'|f_Z)q(f_Z)df_Z \\ &= p(f'|f_Z) \int q(f_Z)df_Z = p(f'|f_Z) \end{aligned} \quad (4.35)$$

上式即为给定任意测试点时，预测的函数分布。