

贝叶斯优化

周志健

1 贝叶斯优化(bayesian optimization)

1.1 简介

贝叶斯优化是一种全局黑盒优化方法。该方法利用先验数据建立目标函数的代理模型，不需要依赖目标函数的梯度，这是此方法可以优化黑盒的根本原因。贝叶斯优化方法主要包含代理模型（高斯过程回归模型）和采集函数这两个核心部分。

1.2 高斯过程(gaussian process)

当我们有 n 个取样点 x_1, \dots, x_n ，以及相应的目标函数值， $y(x_1), \dots, y(x_n)$ ，那么我们可以得到下面这个高斯分布，即：

$$\mathbf{y} \sim N(\mathbf{0}, \mathbf{K}) \quad (1)$$

其中 $\mathbf{y} = (y(x_1) \cdots y(x_n))^T$ ， $\mathbf{K}_{ij} = k(x_i, x_j)$ ， $k(x_i, x_j)$ 为协方差核函数。本文中，协方差核函数选择常用的 RBF 核：

$$k_{\text{RBF}}(x_i, x_j) = \sigma^2 \exp\left(-\frac{r^2}{2l^2}\right) \quad (2)$$

其中 l 和 σ 表示超参数，它们决定了高斯过程曲线的形状， $r = |x_i - x_j|$ 。

因此，训练高斯过程回归模型时，实际是通过梯度法优化其边际似然函数的对数化形式，由此确定参数 l 和 σ 的值，即：

$$\log p(\mathbf{y} | \mathbf{x}_{1:n}) = -\frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}| - \frac{n}{2} \log 2\pi \quad (3)$$

其中 $p(\mathbf{y} | \mathbf{x}_{1:n})$ 是基于训练集的高斯过程的边际似然函数。

指定一个采样点 x_{n+1} ，训练好的高斯过程回归模型可以给出其目标函数值的预测分布，即：

$$y(x_{n+1}) \sim N(\mu(x_{n+1}), \sigma(x_{n+1})) \quad (4)$$

其中均值和方差分别由下面两个公式给出：

$$\mu(x_{n+1}) = \mathbf{K}_{n+1}^T \mathbf{K}^{-1} \mathbf{y} \quad (5)$$

$$\sigma(x_{n+1}) = k(x_{n+1}, x_{n+1}) - \mathbf{K}_{n+1}^T \mathbf{K}^{-1} \mathbf{K}_{n+1} \quad (6)$$

其中 $\mathbf{K}_{n+1} = (k(x_{n+1}, x_1), \dots, k(x_{n+1}, x_n))^T$ 。

1.3 采集函数(acquisition function)

采集函数通过采样点的均值和方差指导下一次采样，通过最大化采集函数以尽量少的迭代次数找到目标函数的最优值，即：

$$x_{\text{next}} = \arg \max_{\tau \in \Omega} \text{AC}(x) \quad (7)$$

其中 $\text{AC}(\cdot)$ 表示采集函数。

1.3.1 单目标 AC

$$UCB = \mu(x) + k\sigma(x) \quad (8)$$

其中 $\mu(x)$ 指的是均值， $\sigma(x)$ 指的是方差。

k 值人为设定，用以权衡 exploration 和 exploitation。

1.3.2 多目标 AC

Pareto 理论：

目标函数空间中存在成对的被支配点和支配点，支配点在任一方向上都优于或等于被支配点。

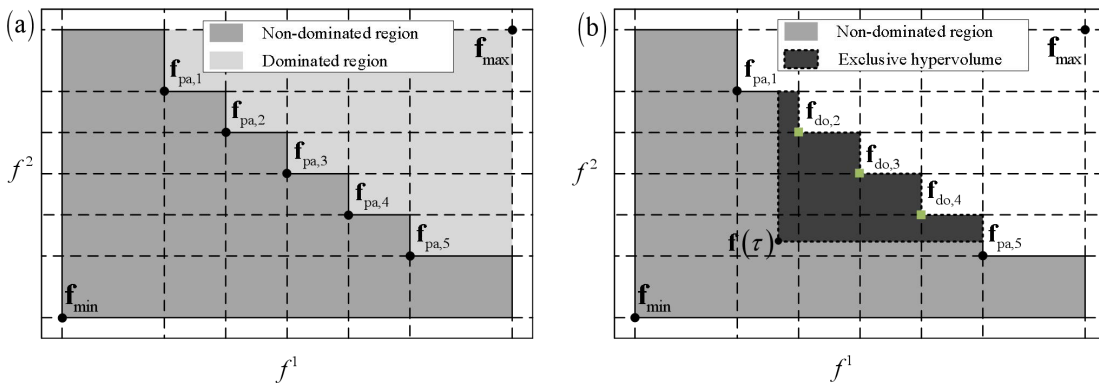


图 1 两个目标函数的帕累托前沿示意图

Emmi(the expected maximin improvement)AC 函数:

$$E[I(\mathbf{f}(x))] \quad (9)$$

其中 $I(\mathbf{y}(x))$ 如下:

$$I(\mathbf{f}(x)) \equiv - \max_{\mathbf{x}_i \in \mathcal{P}_x} \min_{j=1, \dots, m} (f_j(x) - f_j(\mathbf{x}_i)) \times 1 \left[- \max_{\mathbf{x}_i \in \mathcal{P}_x} \min_{j=1, \dots, m} (f_j(x) - f_j(\mathbf{x}_i)) > 0 \right] \quad (10)$$

注意此处 $\mathbf{f} = (f_1(x) \cdots f_n(x))^T$, \mathcal{P}_x 指训练数据中的 pareto 集合。

1.4 伪代码和流程图

Algorithm 1 单目标贝叶斯优化

- 1: 随机选取 n 个样本构成训练集 $D_{0:0}$
 - 2: **for** $i = 1, 2 \dots$ **do**
 - 3: 获取训练集 $X_{1:N}$ 和训练集对应的目标函数值 y
 - 4: 优化似然函数: $\log p(\mathbf{y}|X_{1:n}) = -\frac{1}{2}\mathbf{y}^T K^{-1}\mathbf{y} - \frac{1}{2}\log|K| - \frac{n}{2}\log 2\pi$
 - 5: 求解 $x_i = \argmax(UCB)$
 - 6: 将 x_i 代入黑匣子中获得返回值 $y(x_i)$
 - 7: 将 $y(x_i)$ 和 x_i 加入训练集 $D_{0:i}$ 中构成新的训练集 $D_{0:i+1}$
 - 8: **end for**
-

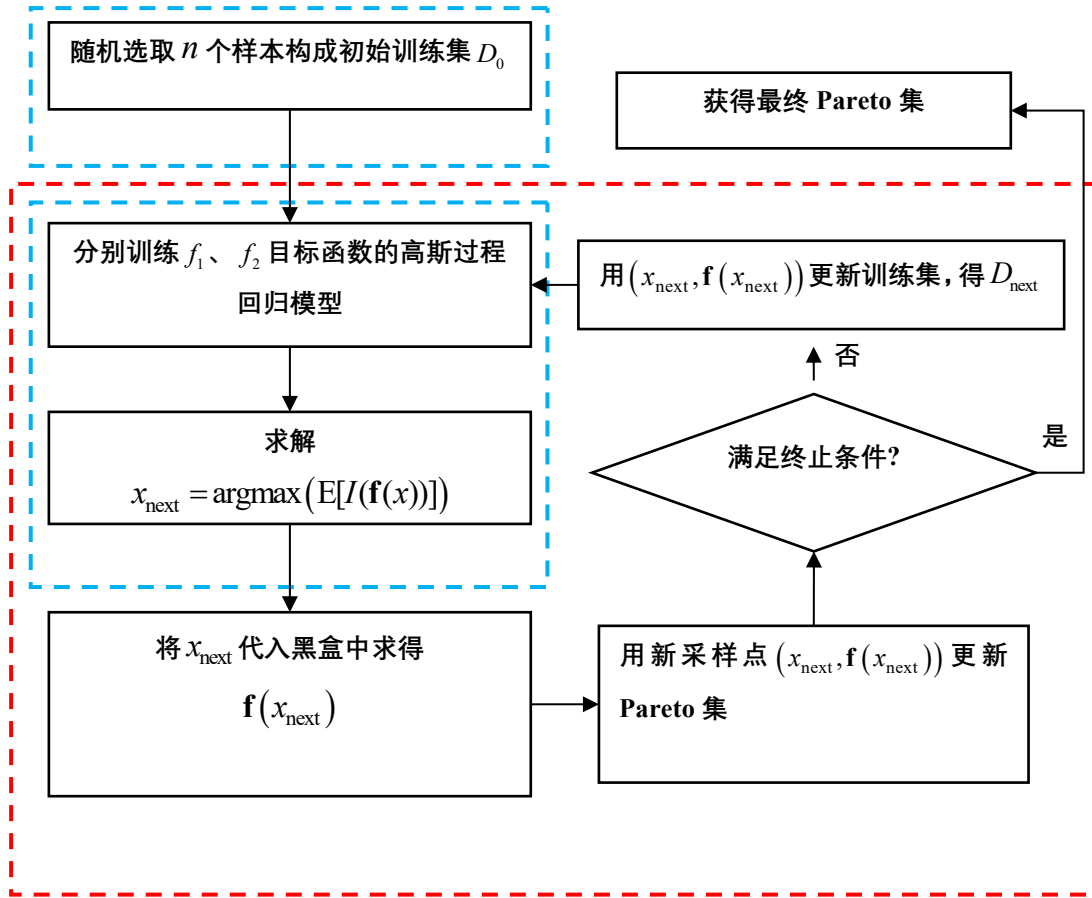


图 2 两个目标函数的贝叶斯优化

2 均方误差

$$\text{MSE} = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)^2 \quad (10)$$

其中 $f(x_i)$ 为模型预测值， y_i 为真实值。

3 总结

MSE 表

手动实现 GP	Sklearn 库 GP
24963.5	41512.1

数据量少时，手动实现的算法效果较好，但算法优化差，时间和空间复杂度都很高。特别是未经优化的矩阵运算导致算法无法应对大数据集。