

高斯过程公式推导

假设有一个训练集 \mathcal{D} ，其中 $\mathcal{D} = \{(x_i, y_i) | i = 1, 2, \dots, n\}$ ，我们用 \mathbf{x} 来表示输入向量，用 y 来表示一个标量输出，其中 \mathbf{x} 的尺寸为 $D \times n$ ， y 的尺寸为 $n \times 1$ ，那么相应的 $\mathcal{D} = (\mathbf{X}, \mathbf{y})$ 。

以标准的线性模型为例，对于：

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w}, \quad y = f(\mathbf{x}) + \varepsilon$$

\mathbf{x} 是输入向量， \mathbf{w} 是权重向量， y 是观察到的目标值，一般来说在 $f(\mathbf{x})$ 中会有偏置向量，但是因为偏置部分可通过 $\mathbf{x}_{new}^T = (\mathbf{x}_{old}^T, 1)$ 来完成，因此此处将偏置隐去。其中 ε 为在 $f(\mathbf{x})$ 上的噪声，我们假设它是独立分布的，相应的均值为0，方差为 σ_n^2 ，那么对应的分布符合：

$$\varepsilon \sim \mathcal{N}(0, \sigma_n^2)$$

相应 y 的分布如下：

$$y \sim \mathcal{N}(f(\mathbf{x}), \sigma_n^2) \sim \mathcal{N}(\mathbf{X}^T \mathbf{w}, \sigma_n^2)$$

即：

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \prod_{i=1}^n p(y_i | x_i, \mathbf{w}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{(y_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma_n^2}\right)$$

在贝叶斯分布中，我们需要对于参数设置一个先验分布，对于权重 \mathbf{w} ，我们假设它的分布为：

$$\mathbf{w} \sim \mathcal{N}(0, \Sigma_p)$$

根据贝叶斯定律：

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}$$

那么：

$$p(w|y, X) = \frac{p(y|X, w)p(w)}{p(y|X)}, \quad p(y|X) = \int p(y|X, w)p(w)dw$$

由于边际概率 $p(y|X)$ 与权重 w 无关，因此 $p(y|X)$ 相当于 $p(w|y, X)$ 的一个常系数，于是：

$$p(w|y, X) \propto \exp\left(-\frac{1}{2\sigma_n^2}(y - X^T w)^T(y - X^T w)\right) \exp\left(-\frac{1}{2}w^T \Sigma_p^{-1}w\right)$$

由于：

$$\begin{aligned} & w^T \Sigma_p^{-1}w + \frac{1}{\sigma_n^2}(y - X^T w)^T(y - X^T w) \\ &= w^T \Sigma_p^{-1}w + \frac{1}{\sigma_n^2}w^T X X^T w + \frac{1}{\sigma_n^2}y^T y - \frac{1}{\sigma_n^2}y^T X^T w - \frac{1}{\sigma_n^2}w^T X y \\ &= w^T \left(\frac{1}{\sigma_n^2}X X^T + \Sigma_p^{-1}\right)w + \frac{1}{\sigma_n^2}y^T y - \frac{1}{\sigma_n^2}y^T X^T w - \frac{1}{\sigma_n^2}w^T X y \end{aligned}$$

令 $A = X X^T + \sigma_n^2 \Sigma_p^{-1}$ ，有 $A = A^T$ ，并记上式为：

$$\frac{1}{\sigma_n^2}(w - \bar{w})^T A (w - \bar{w}) = \frac{1}{\sigma_n^2}(w^T A w + \bar{w}^T A \bar{w} - \bar{w}^T A w - w^T A \bar{w})$$

即：

$$\bar{w}^T A \bar{w} - \bar{w}^T A w - w^T A \bar{w} = y^T y - y^T X^T w - w^T X y$$

那么有 $A \bar{w} = X y$ ，即：

$$\begin{aligned} p(w|y, X) &\propto \exp\left(-\frac{1}{2\sigma_n^2}(w - \bar{w})^T A (w - \bar{w})\right) \\ &\sim \mathcal{N}(\bar{w} = A^{-1}Xy, \sigma_n^2 A^{-1}) \end{aligned}$$

值得注意的是，当我们用传统神经网络 **MSE** 的方式来表示误差：

$$\text{loss} = \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T w)^2 = \frac{1}{n} (y - X^T w)^T (y - X^T w)$$

那么对于 **loss** 求导之后：

$$\frac{\partial \text{loss}}{\partial \mathbf{w}} = \frac{1}{n} X(\mathbf{y} - X^T \mathbf{w}) = 0$$

$$\mathbf{w} = (XX^T)^{-1} X\mathbf{y} \approx A^{-1} X\mathbf{y}$$

这样求得的 \mathbf{w} 与高斯分布中 \mathbf{w} 对应的均值相同。

综合分析可以看出：

$$p(\mathbf{x}_{test}^T \mathbf{w} | \mathbf{x}_{test}, \mathbf{y}, X) \sim \mathcal{N}(\mathbf{x}_{test}^T A^{-1} X\mathbf{y}, \sigma_n^2 \mathbf{x}_{test}^T A^{-1} \mathbf{x}_{test})$$

那么相应的 y_{test} 分布对应：

$$p(y_{test} | \mathbf{x}_{test}, \mathbf{y}, X) \sim \mathcal{N}(\mathbf{x}_{test}^T A^{-1} X\mathbf{y}, \sigma_n^2 + \sigma_n^2 \mathbf{x}_{test}^T A^{-1} \mathbf{x}_{test})$$

以上所有分析都基于线性模型，那么当我们先把输入向量 \mathbf{x} 通过函数 $\phi(\mathbf{x})$ 映射到 M 维的特征空间：

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}, \quad \mathbf{y} = f(\mathbf{x}) + \varepsilon$$

相应的 y_{test} 分布为：

$$\begin{cases} \mu(\mathbf{x}_{test}) = \phi(\mathbf{x}_{test})^T A^{-1} \Phi \mathbf{y} \\ \sigma^2(\mathbf{x}_{test}) = \sigma_n^2 + \sigma_n^2 \phi(\mathbf{x}_{test})^T A^{-1} \phi(\mathbf{x}_{test}) \\ \Phi = (\phi(\mathbf{x}_1), \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_n)) \\ A = \Phi \Phi^T + \sigma_n^2 \Sigma_p^{-1} I \end{cases}$$

进一步推导可知，相应的核函数为：

$$E[f(\mathbf{x})] = E[\phi(\mathbf{x})^T \mathbf{w}] = \phi(\mathbf{x})^T E(\mathbf{w}) = 0$$

$$\begin{aligned}
k(x, y) &= E[(f(x) + \varepsilon - E[f(x) + \varepsilon])(f(y) + \varepsilon - E[f(y) + \varepsilon])] \\
&= E[(f(x) + \varepsilon)(f(y) + \varepsilon)] = E[\phi(x)^T w w^T \phi(y)] + \sigma_n^2 \\
&= \phi(x)^T E[w w^T] \phi(y) + \sigma_n^2 = \phi(x)^T \text{cov}(w) \phi(y) + \sigma_n^2 \\
&= \phi(x)^T \Sigma_p \phi(y) + \sigma_n^2
\end{aligned}$$

那么对应的协方差矩阵为：

$$K_\theta = \Phi^T \Sigma_p \Phi + \sigma_n^2 I$$

那么对应的概率分布为：

$$p(y|X, \theta) = \frac{1}{(2\pi)^{n/2} |K_\theta|^{1/2}} \exp\left(-\frac{1}{2} y^T K_\theta^{-1} y\right)$$

根据《Gaussian Processes for Machine Learning》书中

Appendix.3 可知：

$$K_\theta^{-1} = (\sigma_n^2 I + \Phi^T \Sigma_p \Phi)^{-1} = \frac{1}{\sigma_n^2} I - \frac{1}{\sigma_n^2} \Phi^T \left(\Sigma_p^{-1} + \Phi \frac{1}{\sigma_n^2} \Phi^T \right)^{-1} \Phi \frac{1}{\sigma_n^2}$$

$$K_\theta^{-1} = \frac{1}{\sigma_n^2} I - \frac{1}{\sigma_n^2} \Phi^T (\sigma_n^2 \Sigma_p^{-1} + \Phi \Phi^T)^{-1} \Phi$$

$$K_\theta^{-1} = \frac{1}{\sigma_n^2} I - \frac{1}{\sigma_n^2} \Phi^T A^{-1} \Phi$$

$$|K_\theta| = |\sigma_n^2 I| |\Sigma_p| \left| \Sigma_p^{-1} + \Phi \frac{1}{\sigma_n^2} \Phi^T \right| = |\sigma_n^2 I| |\Sigma_p| \left| \frac{A}{\sigma_n^2} \right|$$

$$\log |K_\theta| = \log |\sigma_n^2 I| + \log \left| \frac{\sigma_p^2}{M} I \right| + \log \left| \frac{A}{\sigma_n^2} \right|$$

$$\log |K_\theta| = N * \log(\sigma_n^2) + \log |A| - M * \log\left(\frac{M \sigma_p^2}{\sigma_p^2}\right)$$

令 $\Sigma_p = \frac{\sigma_p^2}{M} I$ ，有：

$$\log p(y|X, \theta) = -\frac{1}{2} y^T K_\theta^{-1} y - \frac{N}{2} \log(2\pi) - \frac{1}{2} \log |K_\theta|$$

$$\log p(y|X, \theta) = -\frac{1}{2} y^T \left(\frac{1}{\sigma_n^2} I - \frac{1}{\sigma_n^2} \Phi^T A^{-1} \Phi \right) y - \frac{N}{2} \log(2\pi) - \frac{1}{2} \log |K_\theta|$$

$$\log p(y|X, \theta) = -\frac{1}{2\sigma_n^2} (y^T y - y^T \Phi^T A^{-1} \Phi y) - \frac{N}{2} \log(2\pi) - \frac{1}{2} \log |K_\theta|$$

$$\begin{aligned} \log p(y|X, \theta) = & -\frac{1}{2\sigma_n^2} (y^T y - y^T \Phi^T A^{-1} \Phi y) - \frac{N}{2} \log(2\pi\sigma_n^2) - \frac{1}{2} \log |A| \\ & + \frac{M}{2} \log \left(\frac{M\sigma_n^2}{\sigma_p^2} \right) \end{aligned}$$

以上值得注意的是，在建模的过程中，需要模型尽可能地拟合训练集，目标函数为：

$$\text{maximize. } \log p(y|X, \theta)$$

也就是要求：

$$\text{minimize. } -\log p(y|X, \theta)$$

$$\begin{aligned} \text{minimize. } & \frac{1}{2\sigma_n^2} (y^T y - y^T \Phi^T A^{-1} \Phi y) + \frac{N}{2} \log(2\pi\sigma_n^2) + \frac{1}{2} \log |A| \\ & - \frac{M}{2} \log \left(\frac{M\sigma_n^2}{\sigma_p^2} \right) \end{aligned}$$

完成建模后，目标更换为含条件约束的目标优化问题：

$$\begin{aligned} & \text{minimize. } f(x) \\ & \text{s.t. } \begin{cases} c_1(x) < 0 \\ \dots \\ c_{N_c}(x) < 0 \end{cases} \end{aligned}$$

暂时先不考虑约束条件，假设目前已经求得的目标函数 $f(x)$ 最小值为 τ ，那么求解 $f(x)$ 更小值的概率可以求得，由于之前我们已经对于 $y = f(x)$ 进行了高斯过程建模，已知其对应的 $\mu(x)$ 和 $\sigma(x)$ ，improvement function 可以表示为：

$$I(y, \tau) = \begin{cases} \tau - y, & y < \tau \\ 0, & \text{otherwise} \end{cases}$$

那么相应的 Expected Improvement(EI)为：

$$\begin{aligned} E[I(y, \tau)] &= \int_{-\infty}^{+\infty} I(y, \tau) p(y|X, \theta) dy \\ E[I(y, \tau)] &= \int_{-\infty}^{\tau} (\tau - y) p(y|X, \theta) dy \end{aligned}$$

对于 x ，GP 会预测 $y = f(x) \sim \mathcal{N}(\mu, \sigma^2)$ ，令 $y_* = \frac{y - \mu}{\sigma}$ ，则 $y_* \sim \mathcal{N}(0, 1)$ 为标准正态分布，令 $\tau_* = \frac{\tau - \mu}{\sigma}$ ，进一步令 $\Phi(y) = \int_{-\infty}^y \frac{1}{\sqrt{2\pi}} \exp(-0.5 * y^2) dy$ 为 CDF(cumulative distribution function)函数，令 $\phi(y) = \frac{1}{\sqrt{2\pi}} \exp(-0.5 * y^2)$ 为 PDF(probability distribution function)函数，可以进一步推导：

$$\begin{aligned} I(y, \tau) &= \sigma I(y_*, \tau_*) \\ E[I(y, \tau)] &= \sigma E[I(y_*, \tau_*)] \\ E[I(y_*, \tau_*)] &= \int_{-\infty}^{\tau_*} (\tau_* - y_*) \phi(y_*) dy_* \end{aligned}$$

$$E[I(y_*, \tau_*)] = \int_{-\infty}^{\tau_*} \tau_* \phi(y_*) dy_* - \int_{-\infty}^{\tau_*} y_* \phi(y_*) dy_*$$

$$E[I(y_*, \tau_*)] = \tau_* \Phi(\tau_*) - \int_{-\infty}^{\tau_*} y_* \phi(y_*) dy_*$$

$$\text{其中 } y_* \phi(y_*) = \frac{y_*}{\sqrt{2\pi}} \exp(-0.5 * y_*^2) = \left(\frac{-1}{\sqrt{2\pi}} \exp(-0.5 * y_*^2) \right)'$$

$$\int_{-\infty}^{\tau_*} y_* \phi(y_*) dy_* = (-\phi(\tau_*)) - (-\phi(-\infty)) = -\phi(\tau_*)$$

$$E[I(y_*, \tau_*)] = \tau_* \Phi(\tau_*) + \phi(\tau_*)$$

$$E[I(y, \tau)] = \sigma E[I(y_*, \tau_*)] = (\tau - y) \Phi\left(\frac{\tau - \mu}{\sigma}\right) + \sigma \phi\left(\frac{\tau - \mu}{\sigma}\right)$$

以上 **Expected Improvement** 对应没有考虑约束条件的情况，当将约束条件纳入考虑时，可以采用 **weighted Expected**

Improvement(wEI) 计算方式来构造目标函数：

$$I_c(y, c, \tau) = \begin{cases} \tau - y, & y < \tau \text{ and } c < 0 \\ 0, & \text{otherwise} \end{cases}$$

对于每一个约束条件 $c_i(x)$ 都构造一个相应的高斯过程模型，那么每一个 $c_i(x)$ 都有 $c_i(x) \sim \mathcal{N}(\mu_i, \sigma_i)$ ，那么 $c_i(x) < 0$ 约束满足的概率为：

$$\Pr(c_i(x) < 0) = \Phi\left(\frac{0 - \mu_i}{\sigma_i}\right) = \Phi\left(\frac{-\mu_i}{\sigma_i}\right)$$

相应的期望为：

$$E[I_c(y, c, \tau)] = \int_{-\infty}^{\tau} (\tau - y) p(y|X, \theta) \prod_{i=1}^{N_c} \Pr(c_i(x) < 0) dy$$

$$E[I_c(y, c, \tau)] = \left\{ (\tau - y) \Phi\left(\frac{\tau - \mu}{\sigma}\right) + \sigma \phi\left(\frac{\tau - \mu}{\sigma}\right) \right\} \prod_{i=1}^{N_c} \Phi\left(\frac{-\mu_i}{\sigma_i}\right)$$

为了求得最佳 x 需要 **maximize**. $E[I_c(y, c, \tau)]$ 。

值得注意的是 $E[I_c(y, c, \tau)]$ 公式中，我们假设了约束之间相互独立，

在实际问题中往往约束之间存在某种关系，使得其并不相互独立。