

Bayesian Optimization and TPE

参考文献 [1] に基づいて parameter 最適化のアルゴリズムについてまとめる。

1 Sequential Model-based Global Optimization (SMBO)

評価に非常に時間がかかる予測モデル $f: \mathbf{x} \rightarrow \mathbb{R}$ の hyper parameter \mathbf{x} の最適化を考えた場合、直感的な解決方法は予測モデル f の \mathbf{x} に対する performance を低コストに予測できる適当な回帰モデル (response surface model と呼ばれることが多い) を作成し、それを用いて \mathbf{x} の評価を行うことで最適化を行うことである。このような手法は model-based optimization と呼ばれるが、このモデル作成と予測値の収集の過程を繰り返していく手法が Sequential model-based optimization (SMBO) と呼ばれるものである。

以上の話からわかるように、model-based な最適化では以下の 2 点を適当に定めなければならない。

- Hyper parameter に対する予測にどのようなモデルを仮定するか
- そのモデルの評価にどのような指標 (関数) を用いるか

hyperopt や optuna では、予測確率モデルとして Tree-structured Parzen Estimator (TPE) を、評価指標として Expected Improvement [2] (EI)

$$\text{EI}_{y^*}(x) \equiv \int_{-\infty}^{\infty} \max(y^* - y, 0) p_M(y|x) dy \quad (1)$$

を用いている。これは、 f の予測モデル M の下での $f(x)$ の threshold y^* からの期待スコア改善量を表している^{*1}。

なお、その他の評価指標として minimizing the Conditional Entropy of the Minimizer [3] や the bandit-based criterion [4] などがあるが、ここでは EI のみを考える。

^{*1} 一般に評価指標 y は小さいほどよいものが想定されている。

Input: Target algorithm $f: \mathbf{x} \rightarrow \mathbb{R}$, initial SMBO model M_0 , number of trials T , surrogate function S

Output: Observation history \mathcal{H} of parameters and performance ($[(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots]$)

$\mathcal{H} \leftarrow \emptyset$;

for t *in* T **do**

$\mathbf{x}^* \leftarrow \text{argmin}_{\mathbf{x}} S(\boldsymbol{\theta}, M_{t-1})$;

$y \leftarrow f(\mathbf{x}^*)$ (Expensive step);

$\mathcal{H} \leftarrow \mathcal{H} \cup (\mathbf{x}^*, y)$;

$M_t \leftarrow \text{FitModel}(\mathcal{H})$

end

return \mathcal{H}

Algorithm 1: Algorithm framework of generic Sequential Model-Based Optimization (SMBO)

2 The Gaussian Process Approaches (GP)

Gaussian process approach では p にガウス分布を、 y^* に observation の内の最小値

$$y^* = \min (f(x_i), 1 \leq i \leq n), \quad (2)$$

を用いる。

3 Tree-structured Parzen Estimator Approach (TPE)

TPE approach は Bayes の定理

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}, \quad (3)$$

を用いて $p(y|x)$ を計算をする手法である。簡単のため、

$$p(x|y) = \begin{cases} l(x) & (y < y^*) \\ g(x) & (y \geq y^*) \end{cases} \quad (4)$$

とし、 y^* を y の γ 分位点 $p(y < y^*) \equiv \gamma$ とおく。 $l(x), g(x)$ はそれぞれ予測値が y^* 以下、以上になったものの parameter の分布を表す。

このとき EI は

$$\text{EI}_{y^*}(x) = \int_{-\infty}^{y^*} (y^* - y) p(y|x) dy, \quad (5)$$

$$= \int_{-\infty}^{y^*} (y^* - y) \frac{p(x|y)p(y)}{p(x)} dy, \quad (6)$$

$$= \frac{l(x) \left[\gamma y^* - \int_{-\infty}^{y^*} y p(y) dy \right]}{\gamma l(x) + (1 - \gamma) g(x)}, \quad (7)$$

$$\propto \left(\gamma + (1 - \gamma) \frac{g(x)}{l(x)} \right)^{-1}. \quad (8)$$

と変形できることがわかる。従って、EI を最大化するには良い予測値の parameter による分布の値 $l(x)$ を大きく、逆に悪いものの $g(x)$ の値を小さくなるような x を取れば良い。

$l(x), g(x)$ を予測する手法の一つとして Parzen Estimator (カーネル密度推定法とも呼ぶ) がある。カーネル密度推定法は未知の母集団からの sampling とみなせるデータ点の集合から母集団の確率密度分布を推定するノンパラメトリック手法の一つで、 $x_i (i = 1, \dots, n)$ を確率密度 $p(x)$ からの独立な標本としたとき、 $p(x)$ は以下のように表される。

$$p(x) = \frac{1}{nh} \sum_{i=1}^n K \left(\frac{x - x_i}{h} \right) \quad (9)$$

つまり、標本まわりにある band width h でカーネル関数 K を仮定し、それを足し合わせて平均して全体の確率密度関数とする方法である。カーネルとしては標準ガウス関数

$$K(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}, \quad (10)$$

が使われることが多い。

参考文献

- [1] J. Bergstra, R. Bardenet, Y. Bengio and B. Kégl, *Algorithms for hyper-parameter optimization*, in *Advances in Neural Information Processing Systems* (J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira and K. Q. Weinberger, eds.), vol. 24, Curran Associates, Inc., 2011.
- [2] D. R. Jones, *A taxonomy of global optimization methods based on response surfaces*, *Journal of Global Optimization* **21** (2001) 345–383.
- [3] J. Villemonteix, E. Vázquez and E. Walter, *An informational approach to the global optimization of expensive-to-evaluate functions*, *CoRR* **abs/cs/0611143** (2006) [[arXiv:cs/0611143](#)].
- [4] N. Srinivas, A. Krause, S. M. Kakade and M. W. Seeger, *Gaussian process bandits without regret: An experimental design approach*, *CoRR* **abs/0912.3995** (2009) [[arXiv:0912.3995](#)].