

確率雑音反応法による連続系での最適化

01013100 日本IBM東京基礎研究所* 岡野 裕之†
01105930 筑波大学社会学系‡ 香田 正人§

1 はじめに

多変数連続関数の最小化アルゴリズムを提案する。このアルゴリズムは、各変数に印加したガウス雑音を用いて目的関数の微係数を確率的に求め、各変数を勾配方向に更新する勾配法の1つである。本報告ではまずアルゴリズムを述べ、排他的論理和 (XOR) の判別関数を例題として用い、焼鈍し法 (SA)、降下法 (DA) と比較しながらアルゴリズムの性質を示す。さらに1次元部分ブラウン運動 (fBm) を使って分析する。

2 確率雑音反応法

任意の目的関数 $f(x) \in R, x \in R^n$ を、

$$\frac{dx_i}{dt} = -\frac{\partial f(x)}{\partial x_i} \quad (1)$$

のような力学系に基づく勾配法で最小化する。一般的には $f(x)$ の偏導関数を解析的に求めるのは困難なので、何らかの近似が必要となる。そこで、各変数 x_i に $N(0, 1)$ のガウス雑音 ξ_i を印加して $x_i^* = x_i + \xi_i$ とし、 x_i をすべて x_i^* で置き換えた目的関数を $f^*(x^*)$ とする。すると $\frac{\partial f^*(x^*)}{\partial x_i} = \frac{\partial f^*(x^*)}{\partial \xi_i}$ となるが、ノビコフの定理によれば、任意のガウス確率過程 ξ の関数 $H(\xi)$ について $E[\frac{\delta H(\xi)}{\delta \xi_i}] = E[H(\xi)\xi_i]$ が成り立つ。($E[\cdot]$ は平均、 $\frac{\delta H(\xi)}{\delta \xi_i}$ は汎関数微分。) これを用いれば、 $f^*(x^*)$ の各変数に関する微係数を、平均として考えれば雑音成分との積で代用してよいことになる:

$$E\left[\frac{\partial f^*(x^*)}{\partial x_i}\right] = E\left[\frac{\partial f^*(x^*)}{\partial \xi_i}\right] = E\left[f^*(x^*)\xi_i\right] \quad (2)$$

したがって、次のようなアルゴリズムによって任意の目的関数の勾配法を実現できる:

1. 変数ベクトル x に初期解を設定する。
2. **For** N 回ループ **do** // 反復回数 N のループ
3. **begin**
4. 変分ベクトル δx を 0 に初期化する。
5. **For** R 回ループ **do** // 変分 $f^*(x^*)\xi_i$ の蓄積
6. **begin**
7. ガウス雑音 $\xi_i, i = 1, 2, \dots, n$ を生成する。
8. $\delta x_i := \delta x_i - \mu f^*(x^*)\xi_i$ // 変分ベクトル更新
9. **end;**

10. $x_i := x_i + \delta x_i$ // 変数ベクトルに変分を加える
 11. **end;** // それまでに得られた最良解を出力
- この方法を確率雑音反応法 (Stochastic Noise Reaction, SNR) と呼ぶ [1].

2.1 比較のために用いる方法

SA および DA では上記ステップ 4, 7, 10 がなく、ステップ 8 で変数ベクトルを更新する。SA, DA 共に、 $f^*(x^*) < f(x)$ であれば $x := x^*$ とする。SA ではさらに確率 $\exp(-(f^*(x^*) - f(x))/T)$ で $x := x^*$ とする。

3 SNR の性質

SNR は勾配のある領域 ($|\frac{\partial f(x)}{\partial x_i}| > 0$) では目的関数を確率的に減少する一方、平坦な領域 ($\frac{\partial f(x)}{\partial x_i} \simeq 0$) では高確率でその場にとどまる。つまりよい結果を得るには、勾配のある (最適値近くの) 領域を選ぶ必要がある。この条件は SA, DA でも同様である。例えば $\tanh^2(x)$ を最小化する場合、初期値 $[-2, +2]$ の場合にはほぼ 0 に収束する (図 1)。($\mu = 0.1, N = 100, R = 100$)

次に 2 入力 2 隠れノード 1 出力の 3 層ニューラルネットワークで XOR を判別する関数を考える。重みとしきい値を変数 $x_i, 2$ つの入力を $a, b \in \{0, 1\}$ とすると、ネットの出力を与える関数は $h(a, b) = g(v_1 x_7 + v_2 x_8 + x_9)$ と定義できる。ここで $v_1 = \tanh(l(a)x_1 + l(b)x_2 + x_3), v_2 = \tanh(l(a)x_4 + l(b)x_5 + x_6), g(y) = (1 + \tanh(y))/2, l(y) = 2y - 1$ とする。このような関数 $h(a, b)$ が XOR を判別するための x を得る (判別関数の学習) には、次の目的関数を最小化すればよい: $f(x) = (h(0, 0) - 1)^2 + (h(1, 0) - 1)^2 + (h(0, 1) - 1)^2 + h(1, 1)^2/4$ 。各変数 x_i を雑音 $N(0, 1)$ で初期化し、値域を $[-10, +10]$ として実験したところ、SNR は 97% の試行で学習に成功した。この実験では、どれかの変数が値域を越えた時点で変数を再初期化して実行を続けた。

3.1 その他の方法との比較

反復 $N/10$ ごとに $u = 1.0, 0.9, \dots, 0.1$ として温度 $T = -u/\log(p)$ ($p = 0.5$) を設定し、 T を減少するごとにそれまでの最良値で x を再初期化した。すると SA は判別関数の学習にすべての試行で成功した。また単純な降下法 DA でも 92% の割合で成功した。変数の値域を与えない場合の成功率は、SNR が 66%, SA が 78%, DA が 77% となった。いずれの方法でも変数の値域の設定が重要であることが分かる。

* 〒 242-8502 神奈川県大和市下鶴間 1623-14

† E-mail: okanoh@jp.ibm.com

‡ 〒 305-8573 茨城県つくば市天王台 1-1-1

§ E-mail: koda@shako.sk.tsukuba.ac.jp

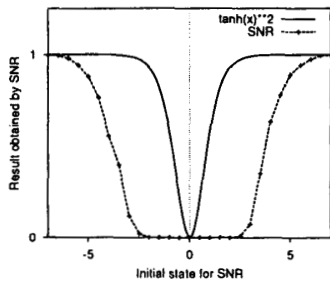


図 1. SNR による $\tanh^2(x)$ の最小化

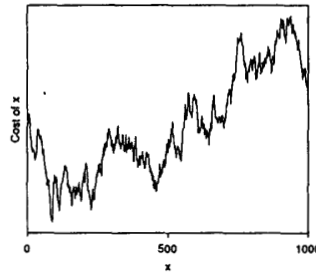


図 2. フラクタル関数の例 ($H = 0.5$ の fBm)

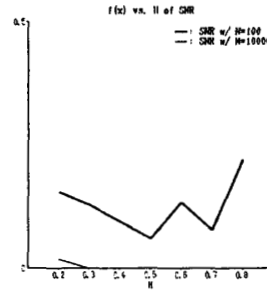


図 3. SNR による fBm の最小値探索

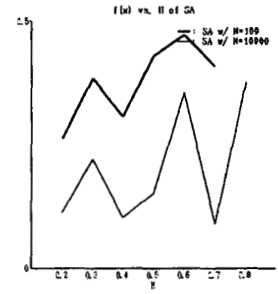


図 4. SA による fBm の最小値探索

4 フラクタル関数の最小化

前述のように、SNR や SA の性能は対象とする目的関数や、適切な値域が設定されるかどうかで異なる。特に勾配が 0 になる領域が広い場合、結果は初期解に大きく依存する。逆に目的関数が「遠くに行く程変化が大きい」(フラクタルな)性質を持つ場合、SA が速く収束することが知られており [2]、SNR も同様の性質に従うと予想される。そこで、フラクタル関数として文献 [2] でも使われている 1 次元の fBm を使って SNR の性質を評価する。fBm は次式で表される：

$$E[(f(x') - f(x))^2] \propto d(x', x)^{2H} \quad (3)$$

H はフラクタル性の強弱を表すパラメータで、値が小さいと関数の地形は荒らく、値が大きいと滑かになる。 $H = 0.5$ の fBm は特にブラウン運動と呼ばれ、 $N(0, 1)$ の雑音を次々に足し合わせることで容易に生成できる(図 2)。任意の目的関数 $f(x)$ について x, x' をランダムに多数選び、 $E[(f(x') - f(x))^2]$ と距離 $d(x', x)$ を log-log でプロットすると、フラクタル性があればそれらが直線に沿い、その傾きが $2H$ となる。アルゴリズムの性質がこの定量的なパラメータ H に関連して理解できれば、問題ごとのアルゴリズムの選択に役立てることができる。

4.1 実験の設定

1 次元の fBm($x \in [0, 1], x \in [0, 10000)$) を用いた。 x の初期値は試行ごとにランダムに生成し、実行時に値域を越えたと再初期化して実行を続けた。fBm($[x]$) の値を配列で保持し、要素間を線形補間することで実現したので、解空間は実質 10000 個の実数ということになる。SNR のパラメータ μ は、焼鈍し温度と同様に反復 $N/10$ ごとに $\mu = 100, 90, \dots, 10$ と設定し、 u を減少するごとにそれまでの最良値で x を再初期化した。SA の温度 T は前述のとおりとした。また $R = 100$ とし、反復回数 $N = 10000, 100$ の 2 通りを行った。 $N = 100$ の場合は解空間と同規模の探索をしていることになる。以上の設定で各 H の fBm を 1 系列ずつ生成し、それぞれについて 100 回の試行の平均値を結果とした。

4.2 SNR による最適化とフラクタル性

H を 0.2 から 0.8 まで変化させて、SNR によって得られた結果をプロットした(図 3)。 $N = 100$ の結果は平均して 0.1 付近だが、 $N = 10000$ ではほぼ最小値が得られている。大きい H でよりよい解が得られたが、 H と結果の間には強い相関は見られない。SNR には山登りの機能が明示的に含まれていないにも関わらず、起伏のある地形での最小化が行えることに注目する。

同様に SA についても実験を行った(図 4)。 $N = 100$ と比較して $N = 10000$ の結果がよくない。これは雑音の分散が小さい(近傍が狭い)ために、探索範囲が初期解付近に限定されたためと考えられる。 H が小さい程よい解が得られており、SNR とは傾向が逆である。SNR と同様に、 H と結果の間には強い相関は見られない。

5 おわりに

ガウス雑音を用いて確率的に勾配を求めることで、任意の目的関数を最小化する確率雑音反応法(SNR)を提案した。9 変数の XOR 判別関数の学習や、複雑な地形のフラクタル関数を用いてその性質を調べたところ、SNR を使ってそれらの関数のよい近似解を得られることが分かった。その際、変数の値域を正しく設定する必要があること、フラクタル性の強弱(H の大小)には解の質が大きくは左右されないことが分かった。また、SNR には山登りの機能が明示的には組み込まれていないが、アルゴリズムのパラメータ μ を適切に設定することで、複雑な地形の関数が最小化できることが示された。

今後の課題として、フラクタル性が弱い関数の地形を特徴づける指標を導入し、アルゴリズムの性能をより詳細に分析することが挙げられる。

参考文献

- [1] M. Koda and H. Okano, "A New Stochastic Learning Algorithm for Neural Networks," J. of Operations Research Society of Japan, to appear.
- [2] G. Sorkin "Efficient Simulated Annealing on Fractal Energy Landscapes," Algorithmica, 6, 367-418 (1991).