

深層学習における適応的共役勾配法

明治大学大学院 *小林悠 KOBAYASHI Yu
01016200 明治大学 飯塚秀明 IIDUKA Hideaki

1. はじめに

経験損失最小化は、深層学習において学習の根幹を担う重要な最適化問題である。その解法として確率的勾配降下法 [1] を筆頭とする確率的勾配を用いた確率的最適化アルゴリズムが知られている。特に、適応的なアルゴリズムと呼ばれる RMSProp [2] や AdaGrad [3], Adam [4] が広く用いられている。Adam は RMSProp や AdaGrad を組み合わせたアルゴリズムであり、局所解への高速な収束性を持つことからその改善手法も多く提案されており、AMSGrad [5] は著名な深層学習ライブラリ¹でも実装されている。

一方、非線形共役勾配法 [6] は、無制約非線形最適化問題のための手法として広く知られている。非線形共役勾配法は、ヘッセ行列の計算を必要としないため、大規模な問題を解く場合でも各反復ごとの計算量も非常に少ないが、最急降下法よりも高速に最適解へ収束するという利点がある。つまり、機械学習のような大規模な訓練データを必要とする問題に対しても有効な手法であることが期待できる。よって、本発表では既存の確率的最適化アルゴリズムの一部に確率的勾配から生成した共役勾配方向を適用したアルゴリズムを提案する。

この発表では、確率的に近似した共役勾配方向を既存の Adam や AMSGrad に取り入れた新しいアルゴリズムを示し、その収束性について議論する。さらに、深層学習を用いた数値実験によって提案手法の性質や性能について議論する。

2. 確率的最適化問題

\mathbb{R}^N を N 次元ユークリッド空間とし、 \mathbb{R}^N の内積 $\langle \cdot, \cdot \rangle$ から誘導されるノルムを $\|\cdot\|$ とする。 $\mathcal{T} := \{1, 2, \dots, T\}$ とする。

仮定 2.1

(A1) $\mathcal{F} \subset \mathbb{R}^N$ は空ではない凸閉集合であり、 \mathcal{F} への射影は容易に計算できる。

(A2) $f: \mathbb{R}^N \rightarrow \mathbb{R}$ は、任意の $\mathbf{x} \in \mathbb{R}^N$ に対して、 $f(\mathbf{x}) := \mathbb{E}[F(\mathbf{x}, \boldsymbol{\xi})]$ と定義され、*well defined* である。ここで、 $F: \mathbb{R}^N \times \Xi \rightarrow \mathbb{R}$ であり、 $\boldsymbol{\xi}$ は $\Xi \subset \mathbb{R}^N$ を台とする確率分布 P に従う乱数ベクトルである。

(A3) 独立同一分布に従う乱数ベクトル $\boldsymbol{\xi}$ の実現値からなる点列 $\boldsymbol{\xi}_0, \boldsymbol{\xi}_1, \dots$ が存在する。

(A4) 任意の $(\mathbf{x}, \boldsymbol{\xi}) \in \mathbb{R}^N \times \Xi$ に対して、確率的勾配 $\mathbf{G}(\mathbf{x}, \boldsymbol{\xi})$ は、 $\mathbf{g}(\mathbf{x}) := \mathbb{E}[\mathbf{G}(\mathbf{x}, \boldsymbol{\xi})]$ 、かつ *well defined* であり、 $\mathbf{g}(\mathbf{x})$ は f の \mathbf{x} における勾配 $\nabla f(\mathbf{x})$ と一致する。

(A5) $M > 0$ が存在し、任意の $(\mathbf{x}, \boldsymbol{\xi}) \in \mathbb{R}^N \times \Xi$ に対して、 $\mathbb{E}[\|\mathbf{G}(\mathbf{x}, \boldsymbol{\xi})\|] \leq M$ を満たす。

問題 2.1 (確率的最適化問題)

$$\text{Minimize } f(\mathbf{x}) := \mathbb{E}[F(\mathbf{x}, \boldsymbol{\xi})] = \frac{1}{T} \sum_{t \in \mathcal{T}} f_t(\mathbf{x})$$

subject to $\mathbf{x} \in \mathcal{F}$.

3. 非線形共役勾配法

反復法は任意の初期点 $\mathbf{x}_0 \in \mathbb{R}^N$ から出発し、次の反復式によって点列を更新する。

$$\mathbf{x}_{n+1} := \mathbf{x}_n + \alpha_n \mathbf{d}_n.$$

非線形共役勾配法は反復法の反復式において探索方向 \mathbf{d}_n を次の式で計算したものである。

$$\mathbf{d}_n := \begin{cases} -\nabla f(\mathbf{x}_n), & n = 0, \\ -\nabla f(\mathbf{x}_n) + \gamma_n \mathbf{d}_{n-1}, & n \geq 1. \end{cases}$$

パラメータ γ_n の選択法としては、以下の公式がよく知られている [6]。ただし、 $\mathbf{y}_{n-1} := \nabla f(\mathbf{x}_n) - \nabla f(\mathbf{x}_{n-1})$ 、 $\lambda > 1/4$ とする。

¹<https://pytorch.org/docs/stable/optim.html>

$$\begin{aligned}\gamma_n^{\text{FR}} &:= \frac{\|\mathbf{g}_n\|^2}{\|\mathbf{g}_{n-1}\|^2} \quad (\text{Fletcher-Reeves}), \\ \gamma_n^{\text{HS}} &:= \frac{\langle \mathbf{g}_n, \mathbf{y}_{n-1} \rangle}{\langle \mathbf{d}_{n-1}, \mathbf{y}_{n-1} \rangle} \quad (\text{Hestenes-Stiefel}), \\ \gamma_n^{\text{PRP}} &:= \frac{\langle \mathbf{g}_n, \mathbf{y}_{n-1} \rangle}{\|\mathbf{g}_{n-1}\|^2} \quad (\text{Polak-Ribière}), \\ \gamma_n^{\text{DY}} &:= \frac{\|\mathbf{g}_n\|^2}{\langle \mathbf{d}_{n-1}, \mathbf{y}_{n-1} \rangle} \quad (\text{Dai-Yuan}), \\ \gamma_n^{\text{HZ}} &:= \frac{\langle \mathbf{g}_n, \mathbf{y}_{n-1} \rangle}{\langle \mathbf{d}_{n-1}, \mathbf{y}_{n-1} \rangle} - \lambda \frac{\|\mathbf{y}_{n-1}\|^2 \langle \mathbf{g}_n, \mathbf{d}_{n-1} \rangle}{\langle \mathbf{d}_{n-1}, \mathbf{y}_{n-1} \rangle^2} \\ &\quad (\text{Hager-Zhang}),\end{aligned}$$

4. Conjugate-gradient-Based Adam

提案アルゴリズムを Algorithm 1 に示す。ここ

Algorithm 1 Conjugate-gradient-Based Adam

Require: $(\alpha_n)_{n \in \mathbb{N}} \subset (0, 1)$, $(\beta_{1n})_{n \in \mathbb{N}} \subset [0, 1)$,

$$\hat{\beta}_1 \in [0, 1)$$

```

1:  $n \leftarrow 0, \mathbf{x}_0 \in \mathbb{R}^N, \mathbf{m}_{-1} := \mathbf{0}, \mathbf{g}_0 := \mathbf{0}$ 
2: loop
3:   for  $t = 1, 2, \dots, T$  do
4:      $\mathbf{G}_t := \mathbf{G}(\mathbf{x}_t, \boldsymbol{\xi}_t)$ 
5:     if  $t = 1$  then
6:        $\mathbf{CG}_t := -\mathbf{G}_t$ 
7:     else
8:        $\mathbf{CG}_t \in \mathbb{R}^N$ 
9:     end if
10:     $\mathbf{m}_t := \beta_{1t} \mathbf{m}_{t-1} + (1 - \beta_{1t}) \mathbf{CG}_t$ 
11:     $\hat{\mathbf{m}}_t := \frac{\mathbf{m}_t}{1 - \hat{\beta}_t}$ 
12:     $\mathbf{H}_t \in \text{diag}((h_{t,i} : h_{t,i} > 0)_{i=1,2,\dots,N})$ 
13:    Find  $\mathbf{d}_t \in \mathbb{R}^N$  that solves  $\mathbf{H}_t \mathbf{d} = -\hat{\mathbf{m}}_t$ 
14:     $\mathbf{x}_t := \Pi_{\mathcal{F}, \mathbf{H}_t}(\mathbf{x}_{t-1} + \alpha_t \mathbf{d}_t)$ 
15:     $t \leftarrow t + 1, n \leftarrow n + 1$ 
16:  end for
17:   $\mathbf{g}_n := \frac{1}{T} \sum_{t \in \mathcal{T}} \nabla f_t(\mathbf{x}_n)$ 
18:   $\bar{\gamma}_n := \gamma(\mathbf{g}_n, \mathbf{g}_{n-1}, \bar{\mathbf{d}}_{n-1})$ 
19:   $\bar{\mathbf{d}}_n := -\mathbf{g}_n + \bar{\gamma}_n \bar{\mathbf{d}}_{n-1}$ 
20: end loop
```

で、確率的共役勾配方向における近似の仕方は、勾配と共役勾配方向とで、それぞれ確率的なものど決定的なものを使い分けることでいくつかの組み合わせを考えることができる。

また、 $\mathbf{H}_t \in \mathbb{R}^{N \times N}$ は学習率を適応的に調整するために計算される行列であり、Adam と AMS-Grad とで計算方法が異なる [7]。この際、確率的勾配 $\mathbf{G}(\mathbf{x}_t, \boldsymbol{\xi}_t)$ を用いている部分で確率的勾配方向 \mathbf{CG}_t ($t \in \mathcal{T}$) を用いることもできる。

さらに、 $\gamma: \mathbb{R}^N \times \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}^N$ は (決定的) 共役勾配方向の更新パラメータを求める関数である。 $\{\gamma_n^{\text{FR}}, \gamma_n^{\text{HS}}, \gamma_n^{\text{PRP}}, \gamma_n^{\text{DY}}, \gamma_n^{\text{HZ}}\}$ から選び、 $\mathbf{d}_{n-1} := \bar{\mathbf{d}}_{n-1}$ とする。

収束解析や数値実験の詳細に関しては、発表内で説明する。

参考文献

- [1] Robbins, H. and Monro, S.: A stochastic approximation method, Herbert Robbins Selected Papers, Springer (1985) 102–109
- [2] Tieleman, T. and Hinton, G.: Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude, COURSE-ERA: Neural networks for machine learning 4 (2) (2012) 26–31
- [3] Duchi, J. and Hazan, E. and Singer, Y.: Adaptive subgradient methods for online learning and stochastic optimization, Journal of Machine Learning Research 2 (2011) 2121–2159
- [4] Kingma, D. P. and Ba, J. L.: Adam, A method for stochastic optimization, arXiv preprint arXiv:1412.6980 (2014).
- [5] Reddi, S. J., Kale, S., and Kumar, S.: On the convergence of adam and beyond, arXiv preprint arXiv:1904.09237 (2019)
- [6] Hager, W. W. and Zhang, H.: A survey of nonlinear conjugate gradient methods, Pacific Journal of Optimization, 2 (1) (2006) 35–58
- [7] Hideaki, I.: Appropriate Learning Rates of Adaptive Learning Rate Optimization Algorithms for Training Deep Neural Networks, arXiv preprint arXiv:2002.09647 (2020)