

Adaptive step size momentum method for deconvolution

Trung Vu and Raviv Raich

School of EECS, Oregon State University, Corvallis, OR, 97331-5501, USA vutru@oregonstate.edu, raich@eecs.oregonstate.edu



1 Introduction

- Deconvolution is the process of reversing the effects of convolution.
- Common iterative algorithms: gradient descent, momentum method, and Newton-Raphson method [1].
- Momentum method: fast convergence, low computational cost, but requires the prior knowledge of the function curvature in choosing the step sizes.
- We propose an adaptive schedule that uses the gradient information to compute the step size for momentum method at each iteration accordingly.

5 Algorithm

Algorithm 1 Adaptive step size schedule for momentum.

1: Given initial guess $\boldsymbol{w}^{(0)}$ and $\boldsymbol{w}^{(1)}$. 2: repeat for k = 1, 2, ...3: $\Delta \boldsymbol{w} = \boldsymbol{w}^{(k)} - \boldsymbol{w}^{(k-1)}$ 4: $\nabla f = \sum_{m,t} \frac{\partial c}{\partial (\boldsymbol{w}^T \boldsymbol{x}_{mt})} \boldsymbol{x}_{mt} + \lambda \boldsymbol{w}$ 5: for m = 1, ..., M, t = 1, ..., n do 6: $p_{mt} = \boldsymbol{x}_{mt}^T \nabla f, q_{mt} = \boldsymbol{x}_{mt}^T \Delta \boldsymbol{w}$ 7: $c_{mt} = \partial^2 c / \partial (\boldsymbol{w}^T \boldsymbol{x}_{mt})^2$

$$\triangleright O(h)$$

$$\triangleright O(Mnh)$$

$$\triangleright O(Mnh)$$

- In a deconvolution setting, the special structure of the objective function allows us to implement the algorithm efficiently without heavy computations of the Hessian.
- Our method asymptotically recovers the optimal convergence rate while only requires twice the number of operations per iteration in gradient descent.

2 Problem Formulation (deconvolution)

- ★ Settings: a training set { x_m, y_m }^M_{m=1} and a convolution kernel w of size h $x_m = [x_m(1), x_m(2), \dots, x_m(n)]^T \qquad x_{mt} = [x_m(t), x_m(t-1), \dots, x_m(t-h+1)]^T$ $y_m = [y_m(1), y_m(2), \dots, y_m(n)]^T$
- Goal: minimize the following objective function

$$f(\boldsymbol{w}) = \sum_{m=1}^{M} \sum_{t=1}^{n} c(\boldsymbol{w}^{T} \boldsymbol{x}_{mt}, \boldsymbol{y}_{m}(t)) + \Omega(\boldsymbol{w})$$

Assumption:

$$0 \preceq \frac{\partial^2 c}{\partial^2 a} \preceq \mu I, \qquad \lambda I \preceq \frac{d^2 \Omega}{d \boldsymbol{w} d \boldsymbol{w}^T} \preceq \gamma I, \qquad \forall \mu, \lambda, \gamma > 0$$

Bounded Hessian:

$$\lambda I \preceq \nabla^2 f(\boldsymbol{w}) \preceq \mu \sum_{m=1}^M R_{\hat{\boldsymbol{x}}_m} + \gamma I, \qquad \forall \boldsymbol{u}$$

8:
$$u = \nabla f^T \nabla f, v = \Delta \boldsymbol{w}^T \nabla f, t = \Delta \boldsymbol{w}^T \Delta \boldsymbol{w}$$
 $\triangleright O(h)$
9: $a = \sum_{m=1}^M \sum_{t=1}^n c_{mt} p_{mt}^2 + \lambda u$ $\triangleright O(Mn)$
10: $b = \sum_{m=1}^M \sum_{t=1}^n c_{mt} q_{mt} + \lambda v$ $\triangleright O(Mn)$
11: $d = \sum_{m=1}^M \sum_{t=1}^n c_{mt} q_{mt}^2 + \lambda t$ $\triangleright O(Mn)$
12: $\alpha^{(k)} = \frac{du - bv}{ad - b^2}, \beta^{(k)} = \frac{bu - av}{ad - b^2}$ $\triangleright O(1)$
13: Update $\boldsymbol{w}^{(k+1)}$ using (1). $\triangleright O(h)$
14: until convergence

6 Computational Complexity

Table 1: Computational complexity of fixed step size gradient (GD), adaptive step size gradient (AGD), fixed step size momentum (MO), adaptive step size momentum (AMO), and Newton's method. ϵ is the relative accuracy.

Method	# Ops. / Iter.	Cvg. rate	# Iters. needed
GD	O(Mnh)	$\frac{r-1}{r+1}$	$\frac{r+1}{2}log(1/\epsilon)$
AGD	O(Mnh)	$\frac{\kappa-1}{\kappa+1}$	$\frac{\kappa+1}{2}log(1/\epsilon)$
MO	O(Mnh)	$\frac{\sqrt{r}-1}{\sqrt{r}+1}$	$\frac{\sqrt{r+1}}{2}log(1/\epsilon)$
AMO	O(Mnh)	$\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}$	$\frac{\sqrt{\kappa}+1}{2}log(1/\epsilon)$
Newton	$O(Mnh^3)$	quadratic	O(1)

3 Preliminary (momentum/heavy-ball method)

Consider the following minimization

$$f(\boldsymbol{w}): \mathbb{R}^d \to \mathbb{R}, \qquad lI \preceq \nabla^2 f(\boldsymbol{w}) \preceq LI, \qquad \forall \boldsymbol{w}$$

Momentum updates

 $\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} - \alpha^{(k)} \nabla f(\boldsymbol{w}^{(k)}) + \beta^{(k)} (\boldsymbol{w}^{(k)} - \boldsymbol{w}^{(k-1)})$ (1)

• Polyak [2] showed that an optimal convergence rate of $\frac{\sqrt{L}-\sqrt{l}}{\sqrt{L}+\sqrt{l}}$ can be obtained on a quadratic by using **constant** step sizes

$$\alpha^{(k)} = \left(\frac{2}{\sqrt{L} + \sqrt{l}}\right)^2 \qquad \beta^{(k)} = \left(\frac{\sqrt{L} - \sqrt{l}}{\sqrt{L} + \sqrt{l}}\right)^2 \qquad (2$$

However, this is not optimal for non-quadratic objectives, where the asymptotic convergence is determined by the local Hessian at the solution [3].

4 Adaptive Step Size Schedule

- ✤ <u>Idea</u>: perform line search on the second-order Taylor expansion of the objective function
- Gradient descent

 $f(\boldsymbol{w} - \alpha \nabla f(\boldsymbol{w})) \approx f(\boldsymbol{w}) - \alpha \nabla f(\boldsymbol{w})^T \nabla f(\boldsymbol{w}) + \frac{1}{2} \alpha^2 \nabla f(\boldsymbol{w})^T \nabla^2 f(\boldsymbol{w}) \nabla f(\boldsymbol{w})$

7 Numerical Example



Fig. 1: Top - an image generated by randomly inserting a sequence of 0, 1, 0, 1. Bottom - the corresponding label series.





Fig. 2: The log-scale decrease in the distance to the solution on domain value side through iterations.



- 1. M. R. Banham and A. K. Katsaggelos, "Digital image restoration," IEEE Signal Processing Magazine, vol. 14, no. 2, pp. 24-41, 1997.
- 2. B. Polyak, "Some methods of speeding up the convergence of iteration methods," in Ussr Computational Mathematics and Mathematical Physics, vol. 4, no. 1, pp. 1–17, 1964.
- 3. B. O'Donoghue and E. Candès, "Adaptive restart for accelerated gradient schemes," Foundation of Computational Mathematics, vol. 15, no. 3, pp. 715-732, 2015.