

# Speed Maintained SVRG

November 12, 2017

## Speed Maintained SVRG

In this section we describe a novel algorithm: Speed Maintained SVRG (SMSVRG), which can set the appropriate iteration number in each epoch automatically and has superior convergence properties in our experiments.

Algorithm 1 just requires two parameters: learning rate  $\eta$ , mini-epoch size  $m_0$ . There are two loops in SMSVRG. In the outer loop (we call each outer iteration as a mini-epoch),  $m_0$  SGD iterations are computed. And then we compute the inequality

$$\frac{\|\tilde{\omega}_{s+1} - \tilde{\omega}_s\|}{\|\tilde{\omega}_s - \tilde{\omega}_{s-1}\|} < \frac{\|\tilde{\omega}_s - \tilde{\omega}_{s-1}\|}{\|\tilde{\omega}_{s-1} - \tilde{\omega}_{s-2}\|} \quad (1)$$

If the inequality holds, we compute a snapshot of the full gradient and step into a new mini-epoch. Otherwise, a new mini-epoch begins directly. It is apparent that we cannot compute the inequality until we have finished two mini-epochs. Note mini-epoch is different from epoch in SVRG, thus we define one epoch as the group of all mini-epochs between two full gradient computations.

In practical experiments, we modify the inequality as follows:

$$\frac{\|\tilde{\omega}_{s+1} - \tilde{\omega}_s\|}{\|\tilde{\omega}_s - \tilde{\omega}_{s-1}\|} < \left( \frac{\|\tilde{\omega}_s - \tilde{\omega}_{s-1}\|}{\|\tilde{\omega}_{s-k} - \tilde{\omega}_{s-k-1}\|} \right)^{\frac{1}{k}} \quad (2)$$

which can reduce the errors incurred by variance and have a better performance than the former.

## Numerical Experiments

### Experimental settings

In this section, we conduct some experiments to illustrate the efficiency of our proposed algorithm, i.e. SMSVRG. We evaluate our algorithm on eight training datasets, which are public on the LIBSVM website<sup>1</sup>. In our experiments, SBSBRG is applied for two standard machine learning tasks:  $l_2$ -regularized logistic regression and  $l_2$ -regularized ridge regression.

The  $l_2$ -regularized logistic regression task is conducted on the four datasets: ijcnn1, a9a, mushrooms, w8a. Since the label of each instance in these datasets is set to be 1 or -1, the loss function of  $l_2$ -regularized logistic regression task is:

$$\min_{\omega} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \omega^T x_i}) + \lambda \|\omega\|^2. \quad (3)$$

The  $l_2$ -regularized ridge regression task is conducted on the four datasets: abalone, cadata, cpusmall, space\_ga. The loss function of  $l_2$ -regularized ridge regression task is:

$$\min_{\omega} \frac{1}{n} \sum_{i=1}^n (\omega^T x_i - y_i)^2 + \lambda \|\omega\|^2. \quad (4)$$

We scale the value of all features to  $[-1, 1]$  and set the weighting parameter  $\lambda$  to  $10^{-4}$  for all evaluations.

<sup>1</sup> <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

---

**Algorithm 1** Speed Maintained SVRG

---

**Require:** learning rate  $\eta$ , minimal epoch size  $m_0$

**Initialize:**  $\tilde{\omega}_0 = \mathbf{0}$ ,  $k=0$

```
1:  $\tilde{\mu} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{\omega}_0)$ 
2: for  $s = 0, 1, \dots$  do
3:    $\omega_0 = \tilde{\omega}_s$ 
4:   for  $t = 0, 1, \dots, m_0$  do
5:     Randomly pick  $i_t \in \{1, 2, \dots, n\}$ 
6:      $\omega_t = \omega_{t-1} - \eta(\nabla f_{i_t} - \nabla f_{i_t}(\tilde{\omega}_s) + \tilde{\mu})$ 
7:   end for
8:    $\tilde{\omega}_{s+1} = \omega_{\frac{n}{2}}$ 
9:   if  $s > 1$  and  $\frac{\|\tilde{\omega}_{s+1} - \tilde{\omega}_s\|}{\|\tilde{\omega}_s - \tilde{\omega}_{s-1}\|} < \frac{\|\tilde{\omega}_s - \tilde{\omega}_{s-1}\|}{\|\tilde{\omega}_{s-1} - \tilde{\omega}_{s-2}\|}$  then
10:     $\tilde{\mu} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{\omega}_{s+1})$ 
11:   end if
12: end for
13: return  $\tilde{\omega}_{s+1}$ 
```

---

## Numerical Results of SMSVRG

We compare SMSVRG with SVRG for optimize (3) and (4). For SMSVRG, we set  $m_0$  to be  $n/2$ , while  $n$  represents the size of datasets. For SVRG, we set epoch size  $m$  as four different values:  $n, 2n, 4n, 6n$ . Our experiments on SVRG show that epoch sizes under  $n$  perform always worse than that equals  $n$ , so we set the lower bound to be  $n$ . Besides, according to the experiments we found epoch sizes bigger than  $6n$  perform almost the same, so we set the upper bound of  $m$  to be  $6n$ . And then it is natural to choose  $m = 2n$  and  $m = 4n$  in this range. For each dataset, we experiment for different learning rates, i.e.  $\eta$  to confirm the theoretical results and insights. In all figures, the  $x$ -axis denotes the computational cost, which is measured by the number of gradient computation divided by the size of training data, i.e.  $n$ . The  $y$ -axis denotes training loss residual, i.e.  $F(\tilde{\omega}_s) - F(\omega^*)$ . Note that the optimum  $\omega^*$  is estimated by running the gradient descent for a long time. In all figures, the dashed lines correspond to SVRG with fixed epoch size given in the legends of the figures, while the green solid lines correspond to SMSVRG

It can be seen from Figures 1(a) to 1(p) that SMSVRG can always have the similar performance as SVRG with most suitable epoch size. We observe that when  $\eta$  is big, setting  $m$  to be a small value, i.e.  $n$ , can achieve better performance. The main reason is that when  $\eta$  is big, the variance becomes big simultaneously, so  $m$  must be set small to constrain the variance. As  $\eta$  diminishes, the optimal value of  $m$  increases, which means that the algorithm can tolerate more variance induced by extra iterations. As illustrated in Figures, our method is comparable to and sometimes even better than SVRG with best-tuned epoch sizes when learning rate is large or medium. However, if  $\eta$  is set to be too small, SMSVRG performs slightly inferior to SVRG with large epoch sizes, but outperforms SVRG with recommended epoch sizes, i.e.  $n$  and  $2n$ . It is noting that setting  $\eta$  to be too small is not a practical approach when using SVRG or its variants, because the convergence rate will be extremely low. Therefore, the sub-optimal performance of SMSVRG with very small  $\eta$  is acceptable.

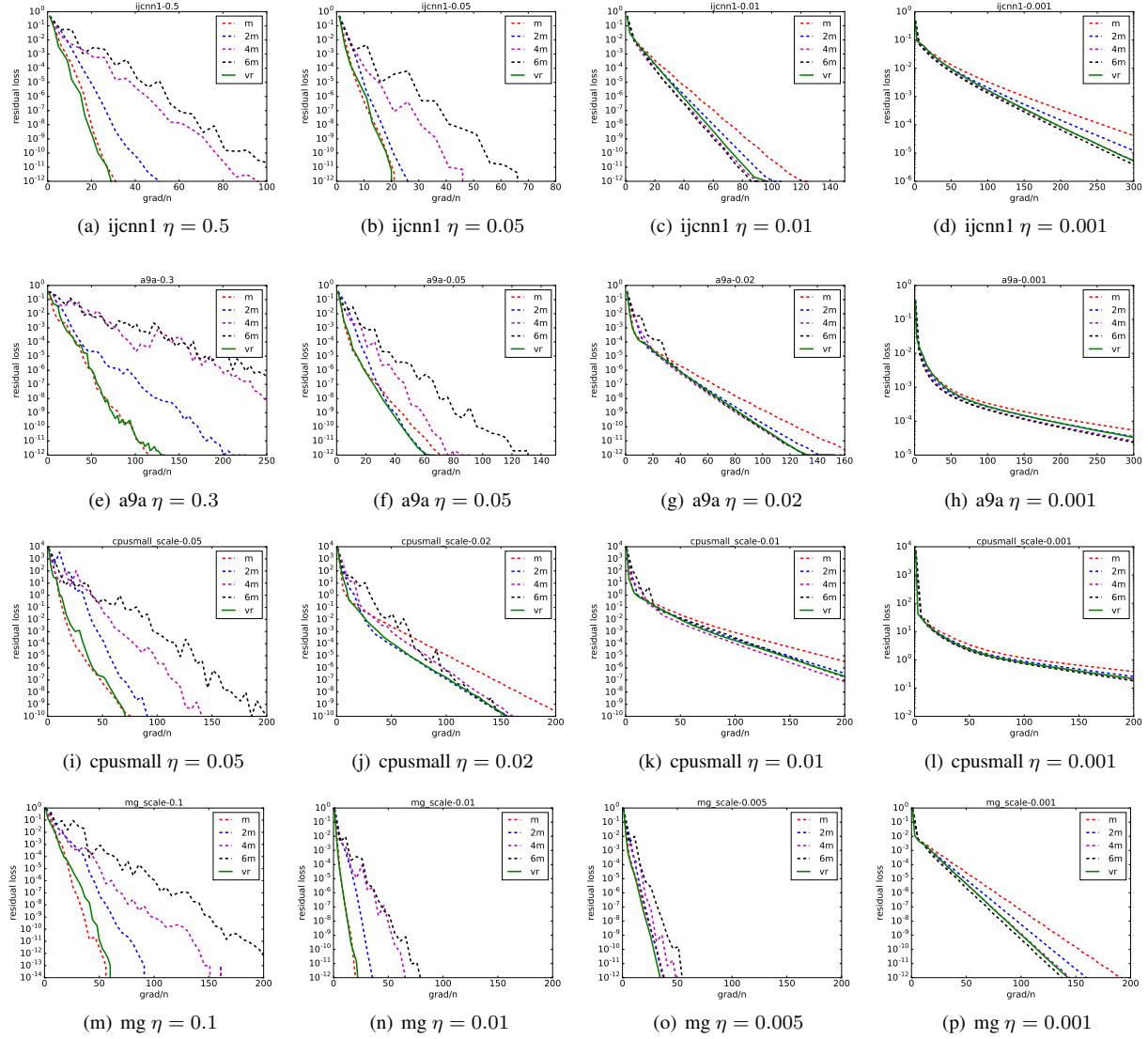


Figure 1: sCSVRG can automatically set a appropriate  $m$  for different learning rates