

Tracking the gradients using the Hessian:  
A new look at variance reducing stochastic  
methods  
Appendix

February 14, 2018

## A Theoretical Analysis: Proofs

Following Bubeck (2015), we consider a single epoch of SVRG and its extensions, that is,  $\bar{\theta} \in \Theta$ , and the iteration, started from  $\theta_0 = \bar{\theta}$ :

$$\theta_t = \Pi_{\Theta} \left( \theta_{t-1} - \gamma \left[ f'_{i_t}(\theta_{t-1}) - z_{i_t}(\theta_{t-1}) + \frac{1}{N} \sum_{j=1}^N z_j(\theta_{t-1}) \right] \right),$$

with  $i_t$  uniformly at random in  $\{1, \dots, N\}$ .

We also recall that, while the results are given using  $R$  the radius of the data, they can be readily transposed to  $L_{\max}$  using  $L_{\max} = R^2$ .

We have, with  $\mathcal{F}_{t-1}$  representing the information up to time  $t$ :

$$\begin{aligned} \mathbf{E} [\|\theta_t - \theta_*\|^2 | \mathcal{F}_{t-1}] &\leq \mathbf{E} \left[ \left\| \theta_{t-1} - \theta_* - \gamma \left[ f'_{i_t}(\theta_{t-1}) - z_{i_t}(\theta_{t-1}) + \frac{1}{N} \sum_{j=1}^N z_j(\theta_{t-1}) \right] \right\|^2 | \mathcal{F}_{t-1} \right] \\ &\quad \text{by contractivity of projections,} \\ &\leq \|\theta_{t-1} - \theta_*\|^2 - 2\gamma F'(\theta_{t-1})^\top (\theta_{t-1} - \theta_*) + \gamma^2 \|F'(\theta_{t-1})\|^2 \\ &\quad + \frac{\gamma^2}{N} \sum_{i=1}^N \left\| f'_i(\theta_{t-1}) - z_i(\theta_{t-1}) - \frac{1}{N} \sum_{j=1}^N f'_j(\theta_{t-1}) + \frac{1}{N} \sum_{j=1}^N z_j(\theta_{t-1}) \right\|^2 \\ &\leq \|\theta_{t-1} - \theta_*\|^2 - 2\gamma F'(\theta_{t-1})^\top (\theta_{t-1} - \theta_*) + \gamma^2 \|F'(\theta_{t-1})\|^2 \\ &\quad + \frac{\gamma^2}{N} \sum_{i=1}^N \|f'_i(\theta_{t-1}) - z_i(\theta_{t-1})\|^2, \text{ by bounding the variance by the second moment.} \end{aligned}$$

In the following sections, we provide proofs for several algorithms we consider in this paper.

## A.1 SVRG

For regular SVRG (we provide the proof for completeness and because we need it later), we have:  $z_i(\theta) = f'_i(\bar{\theta})$  and we consider the bound

$$\begin{aligned} \frac{\gamma^2}{N} \sum_{i=1}^N \|f'_i(\theta_{t-1}) - z_i(\theta_{t-1})\|^2 &\leq \frac{2\gamma^2}{N} \sum_{i=1}^N \|f'_i(\theta_{t-1}) - f'_i(\theta_*)\|^2 + \frac{2\gamma^2}{N} \sum_{i=1}^N \|f'_i(\bar{\theta}) - f'_i(\theta_*)\|^2 \\ &\leq 2\gamma^2 R^2 F'(\theta_{t-1})^\top (\theta_{t-1} - \theta_*) + 2\gamma^2 R^2 [F(\bar{\theta}) - F(\theta_*)] \end{aligned}$$

leading to

$$\begin{aligned} \mathbf{E} [\|\theta_t - \theta_*\|^2 | \mathcal{F}_{t-1}] &\leq \|\theta_{t-1} - \theta_*\|^2 - 2\gamma F'(\theta_{t-1})^\top (\theta_{t-1} - \theta_*) + \gamma^2 \|F'(\theta_{t-1})\|^2 \\ &\quad + 2\gamma^2 R^2 F'(\theta_{t-1})^\top (\theta_{t-1} - \theta_*) + 2\gamma^2 R^2 [F(\bar{\theta}) - F(\theta_*)]. \end{aligned}$$

Thus if  $\gamma \leq 1/(2R^2 + L)$ , we get

$$\mathbf{E} [\|\theta_t - \theta_*\|^2 | \mathcal{F}_{t-1}] \leq \|\theta_{t-1} - \theta_*\|^2 - \gamma [F(\bar{\theta}_{t-1}) - F(\theta_*)] + 2\gamma^2 R^2 [F(\bar{\theta}) - F(\theta_*)].$$

This implies that

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbf{E} [F(\bar{\theta}_{t-1}) - F(\theta_*)] &\leq \frac{1}{\gamma T} \|\bar{\theta} - \theta_*\|^2 + 2\gamma R^2 [F(\bar{\theta}) - F(\theta_*)] \\ \mathbf{E} \left[ F \left( \frac{1}{T} \sum_{t=1}^T \bar{\theta}_{t-1} \right) - F(\theta_*) \right] &\leq \left( \frac{2}{\mu\gamma T} + 2\gamma R^2 \right) [F(\bar{\theta}) - F(\theta_*)]. \end{aligned}$$

This implies that if  $\gamma = \frac{1}{4R^2}$  and  $T \geq 8/(\gamma\mu) = \frac{32R^2}{\mu}$ , then

$$\mathbf{E} \left[ F \left( \frac{1}{T} \sum_{t=1}^T \bar{\theta}_{t-1} \right) - F(\theta_*) \right] \leq \frac{3}{4} [F(\bar{\theta}) - F(\theta_*)].$$

Thus, after  $K = O(\log \frac{1}{\epsilon})$  epochs of SVRG we have attained the required precision, which makes an overall access to gradients of  $KN + KT = (N + \frac{R^2}{\mu}) \log \frac{1}{\epsilon}$ .

## A.2 SVRG-2

We assume that  $\frac{4\beta^2 R^4}{\alpha} D^2 \leq L$  and  $\gamma = 1/(4L)$ .

In this situation, with no approximation, we have  $z_i(\theta) = f'_i(\bar{\theta}) + f''_i(\bar{\theta})(\theta - \bar{\theta})$

and:

$$\begin{aligned}
& \frac{\gamma^2}{N} \sum_{i=1}^N \|f'_i(\theta_{t-1}) - z_i(\theta_{t-1})\|^2 \\
& \leq \frac{\gamma^2}{N} \sum_{i=1}^N R^2 \left[ \frac{\beta}{2} (x_i^\top \theta_{t-1} - x_i^\top \bar{\theta})^2 \right]^2 = \frac{\gamma^2 \beta^2 R^2}{4N} \sum_{i=1}^N (x_i^\top (\theta_{t-1} - \bar{\theta}))^4 \text{ using the bound on } \varphi''', \\
& \leq \frac{\gamma^2 \beta^2 R^2}{N} \sum_{i=1}^N [2(x_i^\top (\theta_{t-1} - \theta_*))^4 + 2(x_i^\top (\theta_* - \bar{\theta}))^4] \\
& \leq \frac{\gamma^2 \beta^2 R^2}{N} \sum_{i=1}^N [2R^2 \|\theta_{t-1} - \theta_*\|^2 (x_i^\top (\theta_{t-1} - \theta_*))^2 + 2R^2 \|\bar{\theta} - \theta_*\|^2 (x_i^\top (\bar{\theta} - \theta_*))^2] \text{ using } \|x_i\| \leq R, \\
& \leq \frac{2\gamma^2 \beta^2 R^4}{N} \|\theta_{t-1} - \theta_*\|^2 \sum_{i=1}^N (x_i^\top (\theta_{t-1} - \theta_*))^2 + \frac{2\gamma^2 \beta^2 R^4}{N} \|\bar{\theta} - \theta_*\|^2 \sum_{i=1}^N (x_i^\top (\bar{\theta} - \theta_*))^2 \\
& \leq \frac{4\gamma^2 \beta^2 R^4}{\alpha} \|\theta_{t-1} - \theta_*\|^2 [F(\theta_{t-1}) - F(\theta_*)] + \frac{4\gamma^2 \beta^2 R^4}{\alpha} \|\bar{\theta} - \theta_*\|^2 [F(\bar{\theta}) - F(\theta_*)] \text{ using } \varphi'' \geq \alpha, \\
& \leq \frac{4\gamma^2 \beta^2 R^4}{\alpha} D^2 [F(\theta_{t-1}) - F(\theta_*)] + \frac{4\gamma^2 \beta^2 R^4}{\alpha} D^2 [F(\bar{\theta}) - F(\theta_*)], \text{ using the compactness of } \Theta.
\end{aligned}$$

With our assumptions, we have  $\gamma \left( L + \frac{4\beta^2 R^4}{\alpha} D^2 \right) \leq 1$ , and we get that

$$\begin{aligned}
\mathbf{E} [\|\theta_t - \theta_*\|^2 | \mathcal{F}_{t-1}] & \leq \|\theta_{t-1} - \theta_*\|^2 - \gamma F'(\theta_{t-1})^\top (\theta_{t-1} - \theta_*) + \frac{4\gamma^2 \beta^2 R^4}{\alpha} D^2 [F(\bar{\theta}) - F(\theta_*)] \\
& \leq \|\theta_{t-1} - \theta_*\|^2 - \gamma [F(\theta_{t-1}) - F(\theta_*)] + \frac{4\gamma^2 \beta^2 R^4}{\alpha} D^2 [F(\bar{\theta}) - F(\theta_*)].
\end{aligned}$$

This leads to, with  $T \geq 4/(\mu\gamma) = \frac{16L}{\mu}$  and using  $\gamma \left( \frac{4\beta^2 R^4}{\alpha} D^2 \right) \leq 1/2$ ,

$$\begin{aligned}
\mathbf{E} \left[ F \left( \frac{1}{T} \sum_{t=1}^T \bar{\theta}_{t-1} \right) - F(\theta_*) \right] & \leq \left( \frac{2}{\mu\gamma T} + \frac{4\gamma\beta^2 R^4}{\alpha} D^2 \right) [F(\bar{\theta}) - F(\theta_*)] \\
& \leq \frac{3}{4} [F(\bar{\theta}) - F(\theta_*)].
\end{aligned}$$

Thus, after  $K = O(\log \frac{1}{\varepsilon})$  epochs of SVRG we have attained the required precision, which makes an overall access to gradients of  $KN + KT = (N + \frac{L}{\mu}) \log \frac{1}{\varepsilon}$ .

### A.3 Stability of SVRG-2

If we make no compactness assumption on  $\Theta$ , then we have:

$$\begin{aligned}
& \frac{\gamma^2}{N} \sum_{i=1}^N \|f'_i(\theta_{t-1}) - z_i(\theta_{t-1})\|^2 \\
& \leq \frac{2\gamma^2}{N} \sum_{i=1}^N \|f'_i(\theta_{t-1}) - f'_i(\bar{\theta})\|^2 + \frac{2\gamma^2}{N} \sum_{i=1}^N \|f''_i(\bar{\theta})(\theta_{t-1} - \bar{\theta})\|^2 \\
& \leq 2\gamma^2 R^2 F'(\theta_{t-1})^\top (\theta_{t-1} - \theta_*) + 2\gamma^2 R^2 [F(\bar{\theta}) - F(\theta_*)] \text{ from the SVRG proof,} \\
& \quad + \frac{2\gamma^2}{N} \sum_{i=1}^N R^2 \|x_i^\top (\theta_{t-1} - \bar{\theta})\|^2 \\
& \leq 2\gamma^2 R^2 F'(\theta_{t-1})^\top (\theta_{t-1} - \theta_*) + 2\gamma^2 R^2 [F(\bar{\theta}) - F(\theta_*)] \\
& \quad \frac{2\gamma^2 R^2}{\alpha} F'(\theta_{t-1})^\top (\theta_{t-1} - \theta_*) + \frac{2\gamma^2 R^2}{\alpha} [F(\bar{\theta}) - F(\theta_*)] \\
& \leq \frac{4\gamma^2 R^2}{\alpha} F'(\theta_{t-1})^\top (\theta_{t-1} - \theta_*) + \frac{4\gamma^2 R^2}{\alpha} [F(\bar{\theta}) - F(\theta_*)]
\end{aligned}$$

Thus, if we take the smaller step-size  $\gamma = \frac{\alpha}{8R^2}$  and  $T = \frac{64R^2}{\alpha\mu}$ , we get the same convergence.

### A.4 Robustness to errors in the Hessian

We assume that  $z_i(\theta) = f'_i(\bar{\theta}) + H_i(\theta - \bar{\theta})$ , with a relative error  $\frac{1}{N} \sum_{i=1}^N (f''_i(\bar{\theta}) - H_i)^2 \preceq R^2 \eta \frac{1}{N} \sum_{i=1}^n f''_i(\bar{\theta})$ . If we take  $H_i = 0$  (plain SVRG), we can take  $\eta = 1$ .

We assume  $\frac{8\beta^2 R^4}{\alpha} D^2 \leq L$  and  $\gamma = 1/(4L)$  with  $8\frac{R^2}{\alpha} \eta \leq L$ . Then

$$\begin{aligned}
& \frac{\gamma^2}{N} \sum_{i=1}^N \|f'_i(\theta_{t-1}) - z_i(\theta_{t-1})\|^2 \\
& \leq 2\frac{\gamma^2}{N} \sum_{i=1}^N \|f'_i(\theta_{t-1}) - f'_i(\bar{\theta}) - f''_i(\bar{\theta})(\theta_{t-1} - \bar{\theta})\|^2 + 2\frac{\gamma^2}{N} \sum_{i=1}^N \|(H_i - f''_i(\bar{\theta}))(\theta_{t-1} - \bar{\theta})\|^2 \\
& \leq \frac{4\gamma^2 \beta^2 R^4}{\alpha} D^2 [F(\theta_{t-1}) - F(\theta_*)] + \frac{4\gamma^2 \beta^2 R^4}{\alpha} D^2 [F(\bar{\theta}) - F(\theta_*)] \\
& \quad + 2\frac{\gamma^2 R^2 \eta}{N} \sum_{i=1}^n (x_i^\top (\bar{\theta} - \theta_{t-1}))^2 \\
& \leq \frac{4\gamma^2 \beta^2 R^4}{\alpha} D^2 [F(\theta_{t-1}) - F(\theta_*)] + \frac{4\gamma^2 \beta^2 R^4}{\alpha} D^2 [F(\bar{\theta}) - F(\theta_*)] \\
& \quad + 4\frac{\gamma^2 R^2 \eta}{\alpha} \left( [F(\bar{\theta}) - F(\theta_*)] + [F(\theta_{t-1}) - F(\theta_*)] \right)
\end{aligned}$$

Thus, with the exact same proof as before (i.e., combining regular SVRG and SVRG2) we reach the desired result.

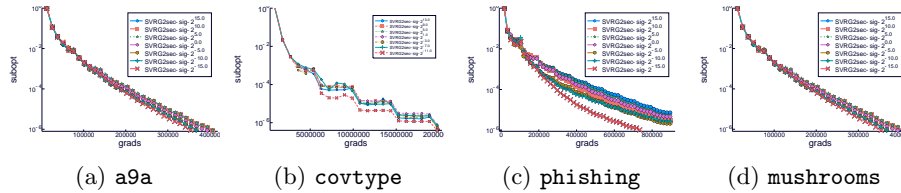


Figure 1: Performance of the SVRG2sec with different choices of  $\sigma$  on: (a) a9a (b) covtype (c) phishing (d) mushrooms.

## B Robustness of the diagonal approximation.

Our robust secant equation has a hyperparameter,  $\sigma^2$ . Since the popularity of an optimization method depends as much of its ease of use as of its convergence rate, we tested the impact of  $\sigma^2$  on the convergence speed. In the supplementary material we show that the impact is generally very limited and that our method is robust to the choice of  $\sigma^2$ . In all other experiments we set  $\sigma^2 = 0.01$ .

## C Additional experiments

We include here the results on two additional LIBSVM datasets, *gisette* and *madelon*.

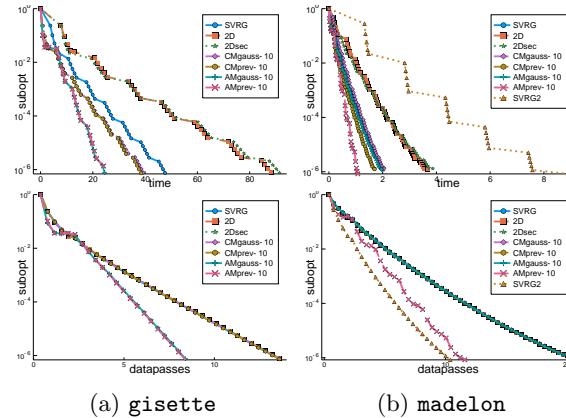


Figure 2: Performance of various SVRG-based methods on LIBSVM test problems: (a) gisette (b) madelon.

## References

- [1] Sébastien Bubeck et al. “Convex optimization: Algorithms and complexity”. In: *Foundations and Trends® in Machine Learning* 8.3-4 (2015), pp. 231–357.