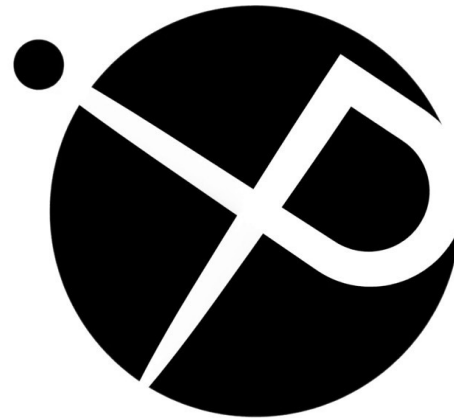


Optimization for Data Science

Stochastic Variance Reduced Gradient Methods

Lecturer: Robert M. Gower & Alexandre Gramfort

Tutorials: Quentin Bertrand, Nidham Gazagnadou



References for this class



O. Sebbouh, N. Gazagnadou, S. Jelassi, F. Bach, R. M. G. **Towards closing the gap between the theory and practice of SVRG**, Neurips 2019.



M. Schmidt, N. Le Roux, F. Bach (2016), Mathematical Programming **Minimizing Finite Sums with the Stochastic Average Gradient**.



RMG, P. Richtárik and Francis Bach (2018) **Stochastic quasi-gradient methods: variance reduction via Jacobian sketching**

EXE: `variance_reduced_exe` + `convergence_prob_exe`

Solving the Finite Sum Training Problem

Optimization Sum of Terms

A Datum Function

$$f_i(w) := \ell(h_w(x^i), y^i) + \lambda R(w)$$

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \ell(h_w(x^i), y^i) + \lambda R(w) &= \frac{1}{n} \sum_{i=1}^n (\ell(h_w(x^i), y^i) + \lambda R(w)) \\ &= \frac{1}{n} \sum_{i=1}^n f_i(w) \end{aligned}$$

Finite Sum Training Problem

$$\min_{w \in \mathbf{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(w) =: f(w)$$

SGD recap

SGD 0.0 Constant stepsize

Set $w^0 = 0$, choose $\alpha > 0$

for $t = 0, 1, 2, \dots, T - 1$

sample $j \in \{1, \dots, n\}$

$$w^{t+1} = w^t - \alpha \nabla f_j(w^t)$$

Output w^T

$$\begin{aligned} \|w^{t+1} - w^*\|_2^2 &= \|w^t - w^* - \alpha \nabla f_j(w^t)\|_2^2 \\ &= \|w^t - w^*\|_2^2 - 2\alpha \langle \nabla f_j(w^t), w^t - w^* \rangle + \alpha^2 \|\nabla f_j(w^t)\|_2^2. \end{aligned}$$

Taking expectation with respect to j

$$\mathbb{E}_j [\|w^{t+1} - w^*\|_2^2] = \|w^t - w^*\|_2^2 - 2\alpha \langle \nabla f(w^t), w^t - w^* \rangle + \alpha^2 \mathbb{E}_j [\|\nabla f_j(w^t)\|_2^2]$$

SGD recap

SGD 0.0 Constant stepsize

Set $w^0 = 0$, choose $\alpha > 0$

for $t = 0, 1, 2, \dots, T - 1$

sample $j \in \{1, \dots, n\}$

$$w^{t+1} = w^t - \alpha \nabla f_j(w^t)$$

Output w^T

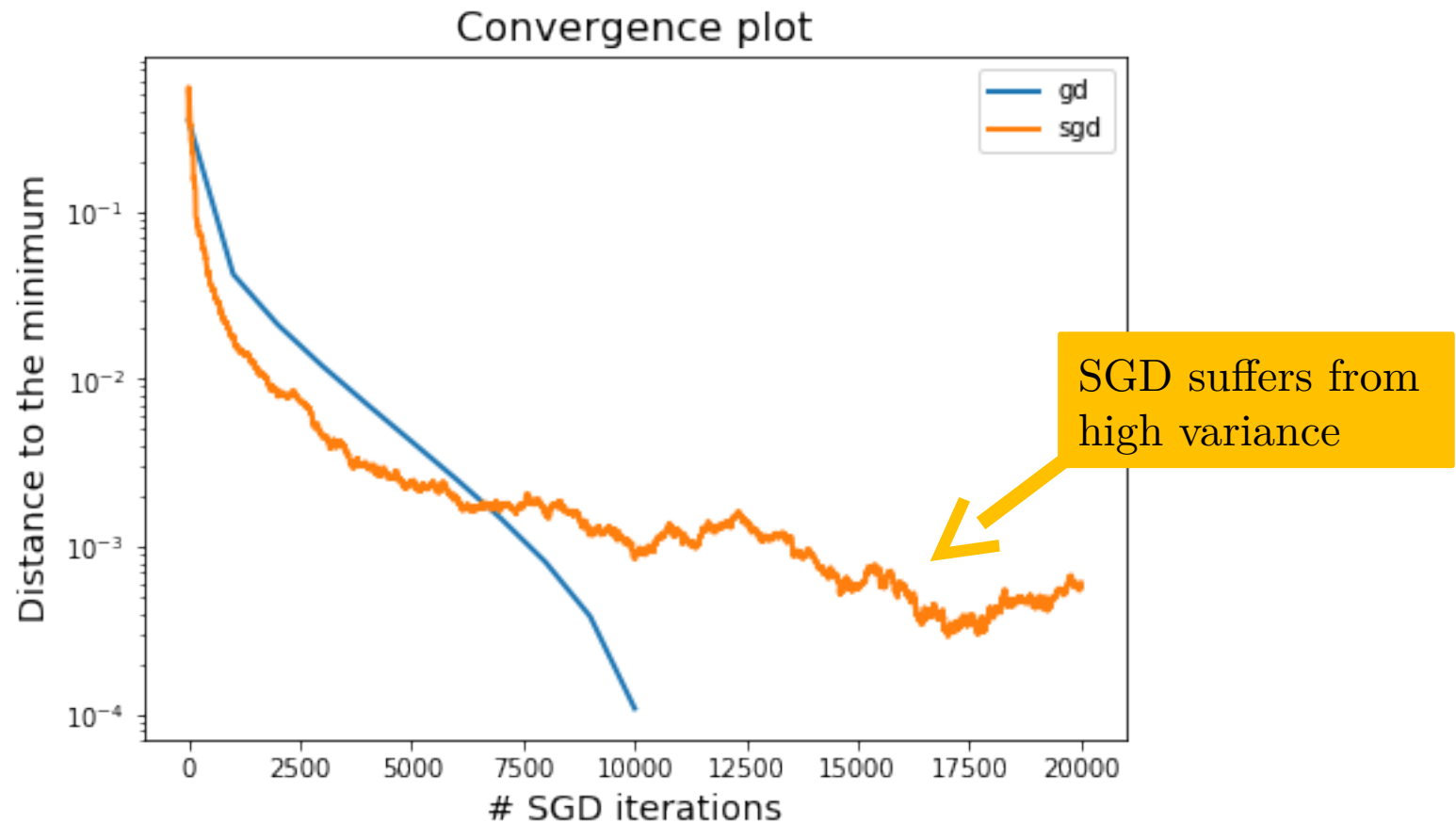
$$\begin{aligned} \|w^{t+1} - w^*\|_2^2 &= \|w^t - w^* - \alpha \nabla f_j(w^t)\|_2^2 \\ &= \|w^t - w^*\|_2^2 - 2\alpha \langle \nabla f_j(w^t), w^t - w^* \rangle + \alpha^2 \|\nabla f_j(w^t)\|_2^2. \end{aligned}$$

Taking expectation with respect to j

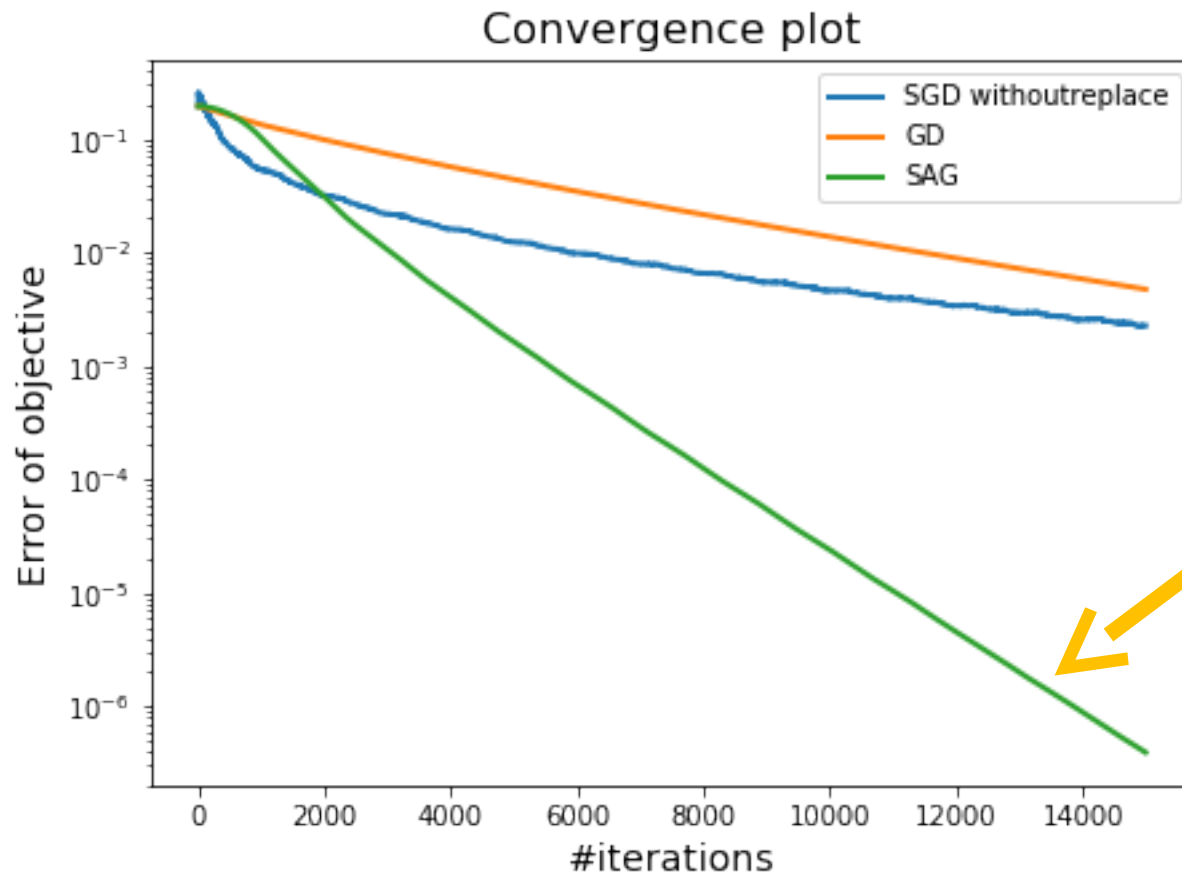
$$\mathbb{E}_j [\|w^{t+1} - w^*\|_2^2] = \|w^t - w^*\|_2^2 - 2\alpha \langle \nabla f(w^t), w^t - w^* \rangle + \alpha^2 \mathbb{E}_j [\|\nabla f_j(w^t)\|_2^2]$$

The Problem: This variance
does not converge

SGD initially fast, slow later



Can we get best of both?



Today we learn
about methods
like this one

Stochastic variance reduced methods

Build an Estimate of the Gradient



Instead of using directly $\nabla f_j(w^t) \approx \nabla f(w^t)$
Use $\nabla f_j(w^t)$ to update estimate $g_t \approx \nabla f(w^t)$



Build an Estimate of the Gradient



Instead of using directly $\nabla f_j(w^t) \approx \nabla f(w^t)$
Use $\nabla f_j(w^t)$ to update estimate $g_t \approx \nabla f(w^t)$



$$w^{t+1} = w^t - \gamma g^t$$

Build an Estimate of the Gradient



Instead of using directly $\nabla f_j(w^t) \approx \nabla f(w^t)$
 Use $\nabla f_j(w^t)$ to update estimate $g_t \approx \nabla f(w^t)$



$$w^{t+1} = w^t - \gamma g^t$$

We would like gradient estimate such that:

Similar

$$g^t \approx \nabla f(w^t)$$

Converges
in L^2

$$\mathbb{E} ||g^t||_2^2 \xrightarrow{w^t \rightarrow w^*} 0$$

Build an Estimate of the Gradient



Instead of using directly $\nabla f_j(w^t) \approx \nabla f(w^t)$
 Use $\nabla f_j(w^t)$ to update estimate $g_t \approx \nabla f(w^t)$



$$w^{t+1} = w^t - \gamma g^t$$

We would like gradient estimate such that:

Typically unbiased
 $\mathbf{E}[g^t] = \nabla f(w^t)$

Similar

$$g^t \approx \nabla f(w^t)$$

Converges
 in L^2

$$\mathbb{E}||g^t||_2^2 \xrightarrow{w^t \rightarrow w^*} 0$$

Build an Estimate of the Gradient



Instead of using directly $\nabla f_j(w^t) \approx \nabla f(w^t)$
 Use $\nabla f_j(w^t)$ to update estimate $g_t \approx \nabla f(w^t)$



$$w^{t+1} = w^t - \gamma g^t$$

We would like gradient estimate such that:

Typically unbiased
 $\mathbf{E}[g^t] = \nabla f(w^t)$

Similar

$$g^t \approx \nabla f(w^t)$$

Converges
 in L^2

$$\mathbb{E}||g^t||_2^2 \xrightarrow{w^t \rightarrow w^*} 0$$

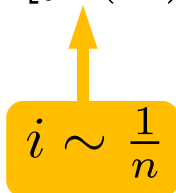
Solves problem of
 $\alpha_t \xrightarrow{t \rightarrow \infty} 0$

Controlled Stochastic Reformulation

Covariate functions:

$$z_i : w \mapsto z_i(w) \in \mathbb{R}, \quad \text{for } i = 1, \dots, n$$

$$\frac{1}{n} \sum_{i=1}^n f_i(w) = \mathbb{E}[f_{\textcolor{red}{i}}(w)] = \mathbb{E}[f_{\textcolor{red}{i}}(w)] - \mathbb{E}[z_{\textcolor{red}{i}}(w)] + \mathbb{E}[z_{\textcolor{red}{i}}(w)]$$


$$i \sim \frac{1}{n}$$

Controlled Stochastic Reformulation

Covariate functions:

$$z_i : w \mapsto z_i(w) \in \mathbb{R}, \quad \text{for } i = 1, \dots, n$$

$$\frac{1}{n} \sum_{i=1}^n f_i(w) = \mathbb{E}[f_{\textcolor{red}{i}}(w)] = \mathbb{E}[f_{\textcolor{red}{i}}(w)] - \mathbb{E}[z_{\textcolor{red}{i}}(w)] + \mathbb{E}[z_{\textcolor{red}{i}}(w)]$$

Controlled Stochastic Reformulation

Covariate functions:

$$z_i : w \mapsto z_i(w) \in \mathbb{R}, \quad \text{for } i = 1, \dots, n$$

$$\frac{1}{n} \sum_{i=1}^n f_i(w) = \mathbb{E}[f_{\textcolor{red}{i}}(w)] = \mathbb{E}[f_{\textcolor{red}{i}}(w)] - \mathbb{E}[z_{\textcolor{red}{i}}(w)] + \mathbb{E}[z_{\textcolor{red}{i}}(w)]$$
$$\textcolor{red}{i} \sim \frac{1}{n} \quad = \mathbb{E}[f_{\textcolor{red}{i}}(w) - z_{\textcolor{red}{i}}(w) + \mathbb{E}[z_{\textcolor{red}{i}}(w)]]$$

Controlled Stochastic Reformulation

Covariate functions:

$$z_i : w \mapsto z_i(w) \in \mathbb{R}, \quad \text{for } i = 1, \dots, n$$

$$\frac{1}{n} \sum_{i=1}^n f_i(w) = \mathbb{E}[f_{\color{red}{i}}(w)] = \mathbb{E}[f_{\color{red}{i}}(w)] - \mathbb{E}[z_{\color{red}{i}}(w)] + \mathbb{E}[z_{\color{red}{i}}(w)]$$

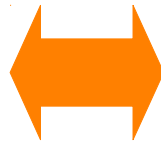
\uparrow
 $i \sim \frac{1}{n}$

$$= \mathbb{E}[f_{\color{red}{i}}(w) - z_{\color{red}{i}}(w) + \mathbb{E}[z_{\color{red}{i}}(w)]]$$

Cancel out
 $\swarrow \searrow$

Original finite sum problem

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(w)$$



Controlled Stochastic Reformulation

$$\min_{w \in \mathbb{R}^d} \mathbb{E}[f_{\color{red}{i}}(w) - z_{\color{red}{i}}(w) + \mathbb{E}[z_{\color{red}{i}}(w)]]$$

Use covariates to **control the variance**

Variance reduction as SGD

$$\min_{w \in \mathbb{R}^d} \mathbb{E} [f_{\textcolor{red}{i}}(w) - z_{\textcolor{red}{i}}(w) + \mathbb{E}[z_{\textcolor{red}{i}}(w)]]$$

Variance reduction as SGD

$$\min_{w \in \mathbb{R}^d} \mathbb{E} [f_{\textcolor{red}{i}}(w) - z_{\textcolor{red}{i}}(w) + \mathbb{E}[z_{\textcolor{red}{i}}(w)]]$$



Sample $\textcolor{red}{i} \sim \frac{1}{n}$

$$w^{t+1} = w^t - \gamma g_{\textcolor{red}{i}}(w^t)$$

Variance reduction as SGD

$$\min_{w \in \mathbb{R}^d} \mathbb{E} [f_{\mathbf{i}}(w) - z_{\mathbf{i}}(w) + \mathbb{E}[z_{\mathbf{i}}(w)]]$$



Sample $\mathbf{i} \sim \frac{1}{n}$

$$w^{t+1} = w^t - \gamma g_{\mathbf{i}}(w^t)$$



$$g_{\mathbf{i}}(w) := \nabla f_{\mathbf{i}}(w) - \nabla z_{\mathbf{i}}(w) + \mathbb{E}[\nabla z_{\mathbf{i}}(w)]$$

Variance reduction as SGD

$$\min_{w \in \mathbb{R}^d} \mathbb{E} [f_{\mathbf{i}}(w) - z_{\mathbf{i}}(w) + \mathbb{E}[z_{\mathbf{i}}(w)]]$$



By design we have that
 $\mathbb{E}[g_{\mathbf{i}}(w^t)] = \nabla f(w^t)$

Sample $\mathbf{i} \sim \frac{1}{n}$

$$w^{t+1} = w^t - \gamma g_{\mathbf{i}}(w^t)$$

$$g_{\mathbf{i}}(w) := \nabla f_{\mathbf{i}}(w) - \nabla z_{\mathbf{i}}(w) + \mathbb{E}[\nabla z_{\mathbf{i}}(w)]$$

Variance reduction as SGD

$$\min_{w \in \mathbb{R}^d} \mathbb{E} [f_{\mathbf{i}}(w) - z_{\mathbf{i}}(w) + \mathbb{E}[z_{\mathbf{i}}(w)]]$$



By design we have that
 $\mathbb{E}[g_{\mathbf{i}}(w^t)] = \nabla f(w^t)$

Sample $\mathbf{i} \sim \frac{1}{n}$

$$w^{t+1} = w^t - \gamma g_{\mathbf{i}}(w^t)$$

$$g_{\mathbf{i}}(w) := \nabla f_{\mathbf{i}}(w) - \nabla z_{\mathbf{i}}(w) + \mathbb{E}[\nabla z_{\mathbf{i}}(w)]$$

How to choose $z_{\mathbf{i}}(w)$?

Choosing the covariate as a linear approximation

Sample $i \sim \frac{1}{n}$

$$w^{t+1} = w^t - \gamma g_i(w^t) := \nabla f_i(w) - \nabla z_i(w) + \mathbb{E}[\nabla z_i(w)]$$

Choosing the covariate as a linear approximation

Sample $i \sim \frac{1}{n}$

$$w^{t+1} = w^t - \gamma g_i(w^t) \quad := \nabla f_i(w) - \nabla z_i(w) + \mathbb{E}[\nabla z_i(w)]$$

We would like:

$$g_i(w) \approx \nabla f(w)$$

Choosing the covariate as a linear approximation

Sample $i \sim \frac{1}{n}$

$$w^{t+1} = w^t - \gamma g_i(w^t) \quad := \nabla f_i(w) - \nabla z_i(w) + \mathbb{E}[\nabla z_i(w)]$$

We would like:

$$g_i(w) \approx \nabla f(w) \quad \Rightarrow \quad \nabla z_i(w) \approx \nabla f_i(w)$$

Choosing the covariate as a linear approximation

Sample $i \sim \frac{1}{n}$

$$w^{t+1} = w^t - \gamma g_i(w^t) \quad := \nabla f_i(w) - \nabla z_i(w) + \mathbb{E}[\nabla z_i(w)]$$

We would like:

$$g_i(w) \approx \nabla f(w) \quad \Rightarrow \quad \nabla z_i(w) \approx \nabla f_i(w)$$

Linear approximation around w

$$z_i(w) = f_i(\tilde{w}) + \langle \nabla f_i(\tilde{w}), w - \tilde{w} \rangle$$

A reference point/ snap shot

SVRG: Stochastic Variance reduced method gradient



Johnson & Zhang, 2013 NIPS

$$w^{t+1} = w^t - \gamma g_{\textcolor{red}{i}}(w^t)$$

Reference point

$$\tilde{w} \in \mathbb{R}^d$$

Sample

$$\nabla f_{\textcolor{red}{i}}(w^t), \quad \text{i.i.d sample with prob } \frac{1}{n}$$

Grad. estimate

$$g_{\textcolor{red}{i}}(w^t) = \nabla f_{\textcolor{red}{i}}(w^t) - \nabla f_{\textcolor{red}{i}}(\tilde{w}) + \nabla f(\tilde{w})$$

It's unbiased
because:

SVRG: Stochastic Variance reduced method gradient



Johnson & Zhang, 2013 NIPS

$$w^{t+1} = w^t - \gamma g_{\mathbf{i}}(w^t)$$

Reference point

$$\tilde{w} \in \mathbb{R}^d$$

Sample

$$\nabla f_{\mathbf{i}}(w^t), \quad \text{i.i.d sample with prob } \frac{1}{n}$$

Grad. estimate

$$g_{\mathbf{i}}(w^t) = \nabla f_{\mathbf{i}}(w^t) - \nabla f_{\mathbf{i}}(\tilde{w}) + \nabla f(\tilde{w})$$

It's unbiased
because:

$$\mathbb{E}[g_{\mathbf{i}}(w)] = \mathbb{E}[\nabla f_{\mathbf{i}}(w)] - \mathbb{E}[\nabla f_{\mathbf{i}}(\tilde{w})] + \nabla f(\tilde{w})$$

SVRG: Stochastic Variance reduced method gradient



Johnson & Zhang, 2013 NIPS

$$w^{t+1} = w^t - \gamma g_{\mathbf{i}}(w^t)$$

Reference point

$$\tilde{w} \in \mathbb{R}^d$$

Sample

$$\nabla f_{\mathbf{i}}(w^t), \quad \text{i.i.d sample with prob } \frac{1}{n}$$

Grad. estimate

$$g_{\mathbf{i}}(w^t) = \nabla f_{\mathbf{i}}(w^t) - \nabla f_{\mathbf{i}}(\tilde{w}) + \nabla f(\tilde{w})$$

It's unbiased
because:

$$\begin{aligned} \mathbb{E}[g_{\mathbf{i}}(w)] &= \mathbb{E}[\nabla f_{\mathbf{i}}(w)] - \mathbb{E}[\nabla f_{\mathbf{i}}(\tilde{w})] + \nabla f(\tilde{w}) \\ &= \nabla f(w) - \nabla f(\tilde{w}) + \nabla f(\tilde{w}) \end{aligned}$$

SVRG: Stochastic Variance reduced method gradient



Johnson & Zhang, 2013 NIPS

$$w^{t+1} = w^t - \gamma g_{\mathbf{i}}(w^t)$$

Reference point

$$\tilde{w} \in \mathbb{R}^d$$

Sample

$$\nabla f_{\mathbf{i}}(w^t), \quad \text{i.i.d sample with prob } \frac{1}{n}$$

Grad. estimate

$$g_{\mathbf{i}}(w^t) = \nabla f_{\mathbf{i}}(w^t) - \nabla f_{\mathbf{i}}(\tilde{w}) + \nabla f(\tilde{w})$$

It's unbiased
because:

$$\begin{aligned} \mathbb{E}[g_{\mathbf{i}}(w)] &= \mathbb{E}[\nabla f_{\mathbf{i}}(w)] - \mathbb{E}[\nabla f_{\mathbf{i}}(\tilde{w})] + \nabla f(\tilde{w}) \\ &= \nabla f(w) - \cancel{\nabla f(\tilde{w})} + \cancel{\nabla f(\tilde{w})} \end{aligned}$$

free-SVRG: Stochastic Variance Reduced Gradients



Jonhson & Zhang
NIPS 2013



Sebbouh, et. al 2019
Neurips 2019

Set $\tilde{w}^0 = 0 = x_0^m$, choose $\gamma > 0, m \in \mathbb{N}$,

$\alpha_t > 0$ with $\sum_{t=0}^{m-1} \alpha_t = 1$

for $s = 1, 2, \dots, T$

$x_s^0 = x_{s-1}^m$

for $t = 0, 1, 2, \dots, m - 1$

i.i.d sample $i \sim \frac{1}{n}$

$g^t = \nabla f_i(x_s^t) - \nabla f_i(\tilde{w}^{s-1}) + \nabla f(\tilde{w}^{s-1})$

$x_s^{t+1} = x_s^t - \gamma g^t$

$\tilde{w}^{s+1} = \sum_{t=0}^{m-1} \alpha_t x_s^t$

Output \tilde{w}^{T+1}

Adding
indices in
 k and t



Most iterates cost $O(1)$



Tune inner loop size m

free-SVRG: Stochastic Variance Reduced Gradients



Jonhson & Zhang
NIPS 2013



Sebbouh, et. al 2019
Neurips 2019

Set $\tilde{w}^0 = 0 = x_0^m$, choose $\gamma > 0, m \in \mathbb{N}$,

$\alpha_t > 0$ with $\sum_{t=0}^{m-1} \alpha_t = 1$

for $s = 1, 2, \dots, T$

$x_s^0 = x_{s-1}^m$

for $t = 0, 1, 2, \dots, m - 1$

i.i.d sample $i \sim \frac{1}{n}$

$g^t = \nabla f_i(x_s^t) - \nabla f_i(\tilde{w}^{s-1}) + \nabla f(\tilde{w}^{s-1})$

$x_s^{t+1} = x_s^t - \gamma g^t$

$\tilde{w}^{s+1} = \sum_{t=0}^{m-1} \alpha_t x_s^t$

Output \tilde{w}^{T+1}

Adding
indices in
 t and s

Adding
indices in
 k and t



Most iterates cost $O(1)$



Tune inner loop size m

free-SVRG: Stochastic Variance Reduced Gradients



Jonhson & Zhang
NIPS 2013



Sebbouh, et. al 2019
Neurips 2019

Set $\tilde{w}^0 = 0 = x_0^m$, choose $\gamma > 0, m \in \mathbb{N}$,

$\alpha_t > 0$ with $\sum_{t=0}^{m-1} \alpha_t = 1$

for $s = 1, 2, \dots, T$

$x_s^0 = x_{s-1}^m$

for $t = 0, 1, 2, \dots, m - 1$

i.i.d sample $i \sim \frac{1}{n}$

$g^t = \nabla f_i(x_s^t) - \nabla f_i(\tilde{w}^{s-1}) + \nabla f(\tilde{w}^{s-1})$

$x_s^{t+1} = x_s^t - \gamma g^t$

$\tilde{w}^{s+1} = \sum_{t=0}^{m-1} \alpha_t x_s^t$

Output \tilde{w}^{T+1}

Adding
indices in
 t and s

Reference point is an
average of inner iterates



Most iterates cost $O(1)$



Tune inner loop size m

SAGA: Stochastic Average Gradient



Defazio, Bach, & Lacoste-Julien, 2014 NIPs

$$w^{t+1} = w^t - \gamma g_{\mathbf{i}}(w^t)$$

Sample

$\nabla f_{\mathbf{i}}(w^t)$, i.i.d sample with prob $\frac{1}{n}$

Grad. estimate

$$g_{\mathbf{i}}(w^t) = \nabla f_{\mathbf{i}}(w^t) - \nabla f_{\mathbf{i}}(w^{t_{\mathbf{i}}}) + \frac{1}{n} \sum_{j=1}^n \nabla f_j(w^{t_j})$$

Store grad.

$$\nabla f_{\mathbf{i}}(w^{t_{\mathbf{i}}}) = \nabla f_{\mathbf{i}}(w^t)$$

SAGA: Stochastic Average Gradient



Defazio, Bach, & Lacoste-Julien, 2014 NIPs

$$w^{t+1} = w^t - \gamma g_{\mathbf{i}}(w^t)$$

Sample

$\nabla f_{\mathbf{i}}(w^t)$, i.i.d sample with prob $\frac{1}{n}$

Grad. estimate

$$g_{\mathbf{i}}(w^t) = \nabla f_{\mathbf{i}}(w^t) - \nabla f_{\mathbf{i}}(w^{t_{\mathbf{i}}}) + \frac{1}{n} \sum_{j=1}^n \nabla f_j(w^{t_j})$$

$$\nabla z_{\mathbf{i}}(w^t) = \nabla f_{\mathbf{i}}(w^{t_{\mathbf{i}}})$$

Store grad.

$$\nabla f_{\mathbf{i}}(w^{t_{\mathbf{i}}}) = \nabla f_{\mathbf{i}}(w^t)$$

SAGA: Stochastic Average Gradient



Defazio, Bach, & Lacoste-Julien, 2014 NIPs

$$w^{t+1} = w^t - \gamma g_i(w^t)$$

Sample

$\nabla f_i(w^t)$, i.i.d sample with prob $\frac{1}{n}$

Grad. estimate

$$g_i(w^t) = \nabla f_i(w^t) - \nabla f_i(w^{t_i}) + \frac{1}{n} \sum_{j=1}^n \nabla f_j(w^{t_j})$$

$$z_i(w) = f_i(w^{t_i}) + \langle \nabla f_i(w^{t_i}), w - w^{t_i} \rangle$$

$$\nabla z_i(w^t) = \nabla f_i(w^{t_i})$$

Store grad.

$$\nabla f_i(w^{t_i}) = \nabla f_i(w^t)$$

SAGA: Stochastic Average Gradient



Defazio, Bach, & Lacoste-Julien, 2014 NIPs

$$w^{t+1} = w^t - \gamma g_i(w^t)$$

Sample

$\nabla f_i(w^t)$, i.i.d sample with prob $\frac{1}{n}$

Grad. estimate

$$g_i(w^t) = \nabla f_i(w^t) - \nabla f_i(w^{t_i}) + \frac{1}{n} \sum_{j=1}^n \nabla f_j(w^{t_j})$$

$$z_i(w) = f_i(w^{t_i}) + \langle \nabla f_i(w^{t_i}), w - w^{t_i} \rangle$$



$$\nabla z_i(w^t) = \nabla f_i(w^{t_i})$$

$$\mathbb{E}[\nabla z_i(w^t)]$$

Store grad.

$$\nabla f_i(w^{t_i}) = \nabla f_i(w^t)$$

SAGA: Stochastic Average Gradient

Set $w^0 = 0, g_i = \nabla f_i(w^0)$, for $i = 1 \dots, n$

Choose $\gamma > 0$

for $t = 0, 1, 2, \dots, T - 1$

sample $i \in \{1, \dots, n\}$

$$g^t = \nabla f_i(w^t) - g_i + \frac{1}{n} \sum_{j=1}^n g_j$$

$$w^{t+1} = w^t - \gamma g^t$$

$$g_i = \nabla f_i(w^t)$$

Output w^T



No inner loop, rolling update



Stores a $d \times n$ matrix

Covariates

Let x and z be random variables. We say that x and z are covariates if:

Variance Reduced Estimate:

$$\text{cov}(x, z) \geq 0$$

$$x_z = x - z + \mathbb{E}[z]$$

Covariates

$$\text{cov}(x, z) := \mathbb{E}[(x - \mathbb{E}[x])(z - \mathbb{E}[z])]$$

Let x and z be random variables. We say that x and z are covariates if:

$$\text{cov}(x, z) \geq 0$$

Variance Reduced Estimate:

$$x_z = x - z + \mathbb{E}[z]$$

EXE:

1. Show that $\mathbb{E}[x_z] = \mathbb{E}[x]$
2. $\text{VAR}[x_z] = \mathbb{E}[(x_z - \mathbb{E}[x_z])^2] = ?$
3. When is $\text{VAR}[x_z] \leq \text{VAR}[x]$

Covariates

$$\text{cov}(x, z) := \mathbb{E}[(x - \mathbb{E}[x])(z - \mathbb{E}[z])]$$

Let x and z be random variables. We say that x and z are covariates if:

$$\text{cov}(x, z) \geq 0$$

Variance Reduced Estimate:

$$x_z = x - z + \mathbb{E}[z]$$

EXE:

1. Show that $\mathbb{E}[x_z] = \mathbb{E}[x]$
2. $\text{VAR}[x_z] = \mathbb{E}[(x_z - \mathbb{E}[x_z])^2] = ?$
3. When is $\text{VAR}[x_z] \leq \text{VAR}[x]$

$$\begin{aligned} \mathbb{E}[(x_z - \mathbb{E}[x_z])^2] &= \mathbb{E}[(x - \mathbb{E}[x] - (z - \mathbb{E}[z]))^2] \\ &= \mathbb{E}[(x - \mathbb{E}[x])^2] - 2\mathbb{E}[(x - \mathbb{E}[x])(z - \mathbb{E}[z])] \\ &\quad + \mathbb{E}[(z - \mathbb{E}[z])^2] \\ &= \text{VAR}[x] - 2\text{cov}(x, z) + \text{VAR}[z] \end{aligned}$$

SAG: Stochastic Average Gradient (Biased version)



M. Schmidt, N. Le Roux, F. Bach (2016), Math prog

$$w^{t+1} = w^t - \gamma g_{\mathbf{i}}(w^t)$$

Sample

$\nabla f_{\mathbf{i}}(w^t)$, i.i.d sample with prob $\frac{1}{n}$

Grad. estimate

$$g_{\mathbf{i}}(w^t) = \frac{1}{n} \sum_{j=1}^n \nabla f_j(w^{t_j})$$

$$\mathbb{E}[g^t] \neq \nabla f(w^t)$$

~~$$x_z = x - z + \mathbb{E}[z]$$~~

Store grad.

$$\nabla f_{\mathbf{i}}(w^{t_{\mathbf{i}}}) = \nabla f_{\mathbf{i}}(w^t)$$

SAG: Stochastic Average Gradient

Set $w^0 = 0, g_i = \nabla f_i(w^0)$, for $i = 1, \dots, n$

Choose $\gamma > 0$

for $t = 0, 1, 2, \dots, T - 1$

sample $i \in \{1, \dots, n\}$

$g_i = \nabla f_i(w^t)$ (update grad)

$g^t = \frac{1}{n} \sum_{j=1}^n g_j$

$w^{t+1} = w^t - \gamma g^t$

Output w^T



Very easy to implement



Stores a $d \times n$ matrix

SAG: Stochastic Average Gradient

Set $w^0 = 0, g_i = \nabla f_i(w^0)$, for $i = 1, \dots, n$

Choose $\gamma > 0$

for $t = 0, 1, 2, \dots, T - 1$

sample $i \in \{1, \dots, n\}$

$g_i = \nabla f_i(w^t)$ (update grad)

$g^t = \frac{1}{n} \sum_{j=1}^n g_j$

$w^{t+1} = w^t - \gamma g^t$

Output w^T



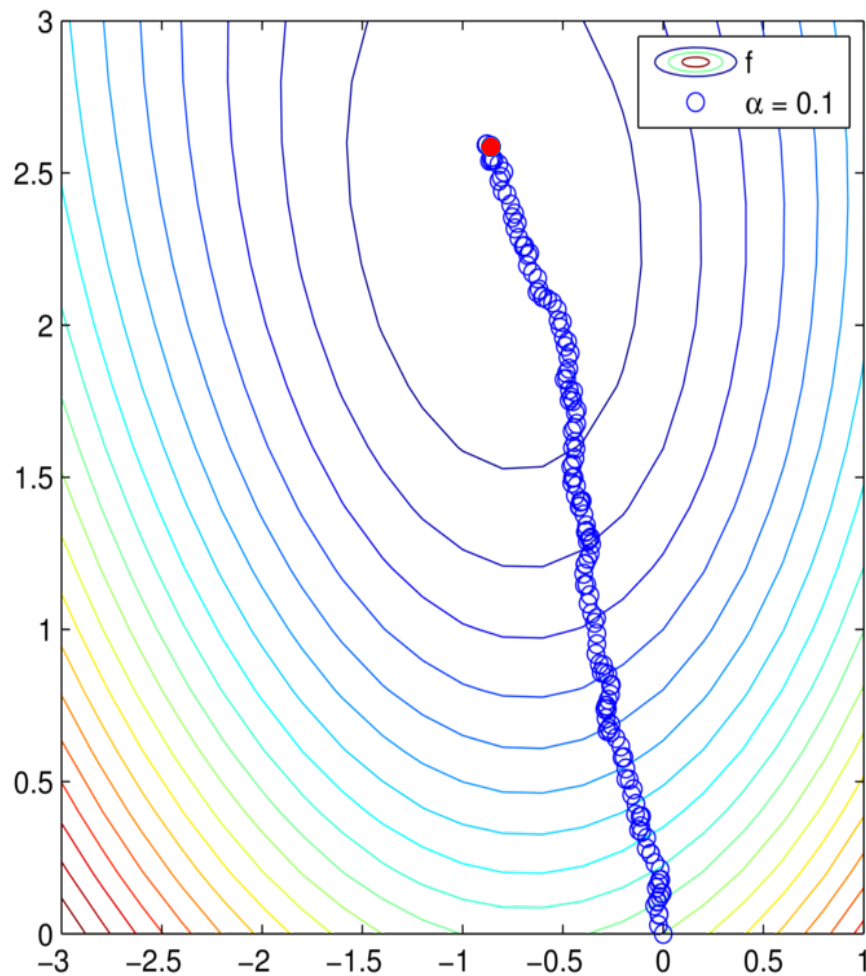
Very easy to implement



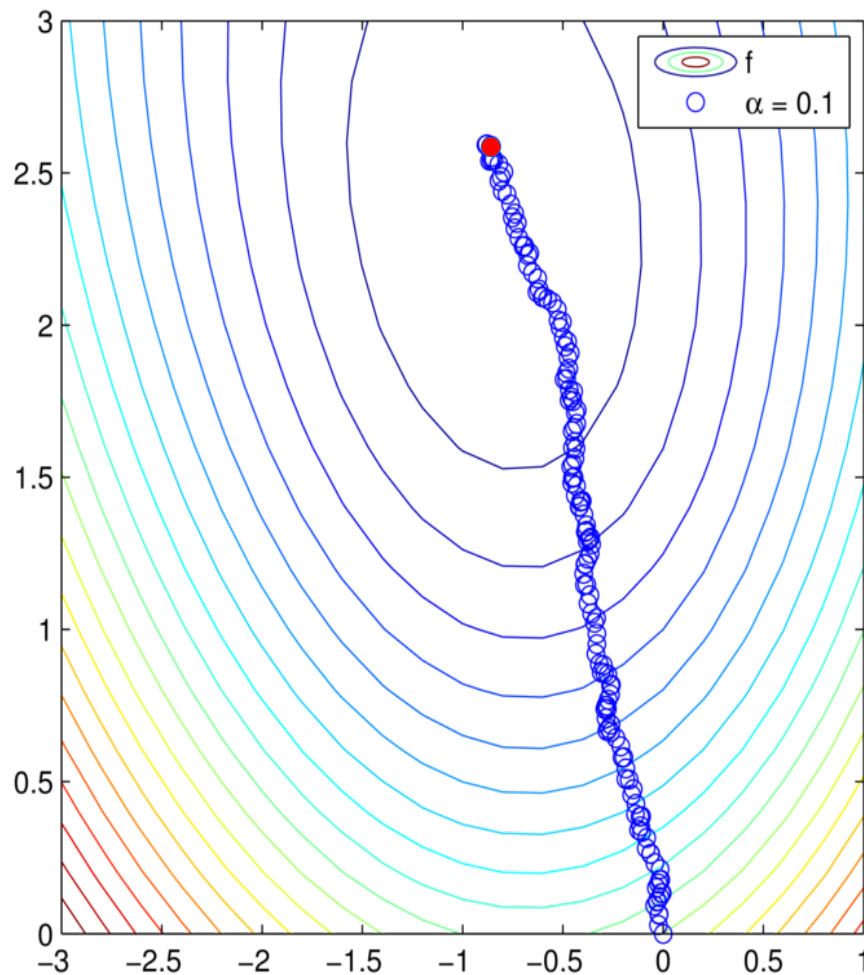
Stores a $d \times n$ matrix

EXE: Introduce a variable $G = (1/n) \sum_{j=1}^n g_j$. Re-write the SAG algorithm so G is updated efficiently at each iteration.

The Stochastic Average Gradient



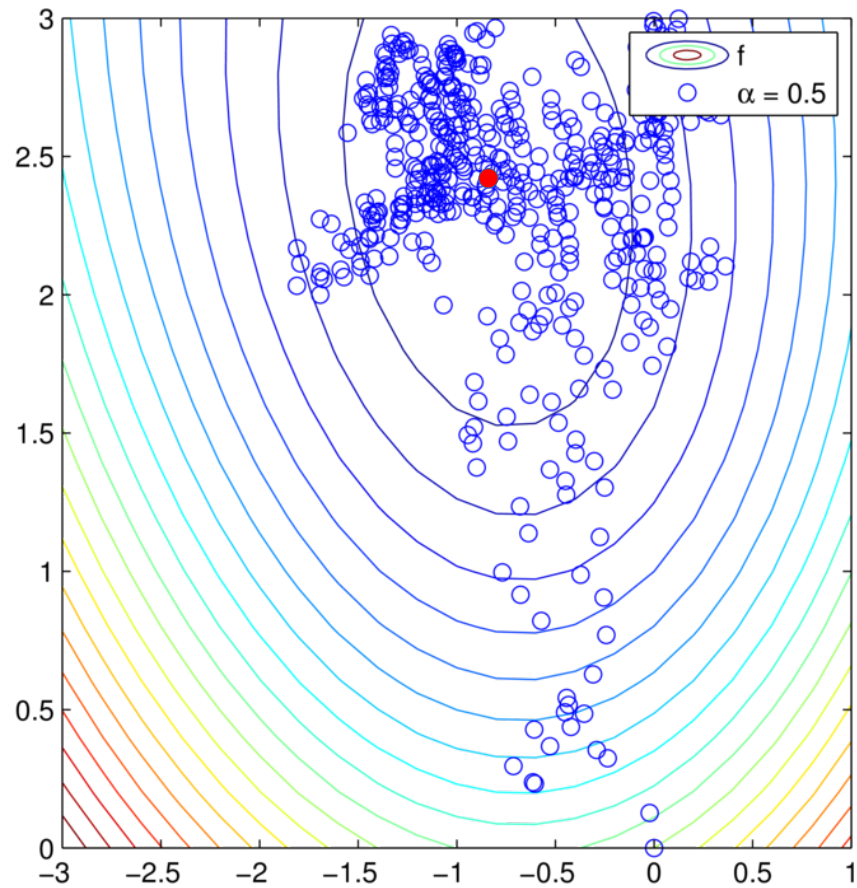
The Stochastic Average Gradient



How to prove this converges? Is this the only option?

Stochastic Gradient Descent

$\alpha = 0.5$



Convergence Theorems

Assumptions for Convergence

Strong Convexity

$$f(w) \geq f(y) + \langle \nabla f(y), w - y \rangle + \frac{\mu}{2} \|w - y\|_2^2$$

Smoothness + convexity

$$f_i(w) \leq f_i(y) + \langle \nabla f_i(y), w - y \rangle + \frac{L_i}{2} \|w - y\|_2^2$$

$$f_i(w) \geq f_i(y) + \langle \nabla f_i(y), w - y \rangle \quad \text{for } i = 1, \dots, n$$

$$L_{\max} := \max_{i=1, \dots, n} L_i$$

free-SVRG: Stochastic Variance Reduced Gradients



Jonhson & Zhang
NIPS 2013



Sebbouh, et. al 2019
Neurips 2019

Set $\tilde{w}^0 = 0 = x_0^m$, choose $\gamma > 0, m \in \mathbb{N}$,

$\alpha_t > 0$ with $\sum_{t=0}^{m-1} \alpha_t = 1$

for $s = 1, 2, \dots, T$

$x_s^0 = x_{s-1}^m$

for $t = 0, 1, 2, \dots, m - 1$

i.i.d sample $i \sim \frac{1}{n}$

$g^t = \nabla f_i(x_s^t) - \nabla f_i(\tilde{w}_{s-1}) + \nabla f(\tilde{w}_{s-1})$

$x_s^{t+1} = x_s^t - \gamma g^t$

$\tilde{w}^{s+1} = \sum_{t=0}^{m-1} \alpha_t x_s^t$

Output \tilde{w}^{T+1}

Adding
indices in
 k and t



Most iterates cost $O(1)$



Tune inner loop size m

free-SVRG: Stochastic Variance Reduced Gradients



Jonhson & Zhang
NIPS 2013



Sebbouh, et. al 2019
Neurips 2019

Set $\tilde{w}^0 = 0 = x_0^m$, choose $\gamma > 0, m \in \mathbb{N}$,

$\alpha_t > 0$ with $\sum_{t=0}^{m-1} \alpha_t = 1$

for $s = 1, 2, \dots, T$

$x_s^0 = x_{s-1}^m$

for $t = 0, 1, 2, \dots, m-1$

i.i.d sample $i \sim \frac{1}{n}$

$g^t = \nabla f_i(x_s^t) - \nabla f_i(\tilde{w}_{s-1}) + \nabla f(\tilde{w}_{s-1})$

$x_s^{t+1} = x_s^t - \gamma g^t$

$\tilde{w}^{s+1} = \sum_{t=0}^{m-1} \alpha_t x_s^t$

Output \tilde{w}^{T+1}

Adding
indices in
 k and t

$$\alpha_k = \frac{(1 - \gamma\mu)^{m-1-t}}{\sum_{i=0}^{m-1} (1 - \gamma\mu)^{m-1-i}}$$



Most iterates cost $O(1)$



Tune inner loop size m

Convergence Theorem for SVRG

Theorem

If $f(w)$ is μ -strongly convex, $f_i(w)$ is L_{\max} -smooth

$$\Psi(x, \tilde{w}) := \|x - w^*\|^2 + \text{cnst} \times (f(\tilde{w}) - f(w^*))$$

$$\text{where } \text{cnst} := 8L_{\max}\gamma^2 \sum_{i=1}^{m-1} (1 - \gamma\mu)^i$$

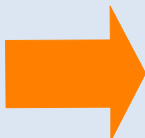
Convergence Theorem for SVRG

Theorem

If $f(w)$ is μ -strongly convex, $f_i(w)$ is L_{\max} -smooth

$$\Psi(x, \tilde{w}) := \|x - w^*\|^2 + \text{cnst} \times (f(\tilde{w}) - f(w^*))$$

If $\gamma \leq \frac{1}{6L_{\max}}$ then

 $\mathbb{E}[\Psi(x_s^m, \tilde{w}_s)] \leq \max \left\{ (1 - \gamma\mu)^m, \frac{1}{2} \right\}^t \Psi(x_0^0, \tilde{w}_0)$

where $\text{cnst} := 8L_{\max}\gamma^2 \sum_{i=1}^{m-1} (1 - \gamma\mu)^i$

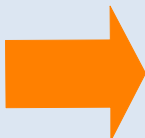
Convergence Theorem for SVRG

Theorem

If $f(w)$ is μ -strongly convex, $f_i(w)$ is L_{\max} -smooth

$$\Psi(x, \tilde{w}) := \|x - w^*\|^2 + \text{cnst} \times (f(\tilde{w}) - f(w^*))$$

If $\gamma \leq \frac{1}{6L_{\max}}$ then

 $\mathbb{E}[\Psi(x_s^m, \tilde{w}_s)] \leq \max \left\{ (1 - \gamma\mu)^m, \frac{1}{2} \right\}^t \Psi(x_0^0, \tilde{w}_0)$

where $\text{cnst} := 8L_{\max}\gamma^2 \sum_{i=1}^{m-1} (1 - \gamma\mu)^i$

Free to choose the number of inner iterates m

Convergence Theorem for SVRG

Theorem

If $f(w)$ is μ -strongly convex, $f_i(w)$ is L_{\max} -smooth

$$\Psi(x, \tilde{w}) := \|x - w^*\|^2 + \text{cnst} \times (f(\tilde{w}) - f(w^*))$$

If $\gamma \leq \frac{1}{6L_{\max}}$ then

$$\Rightarrow \mathbb{E}[\Psi(x_s^m, \tilde{w}_s)] \leq \max \left\{ (1 - \gamma\mu)^m, \frac{1}{2} \right\}^t \Psi(x_0^0, \tilde{w}_0)$$

$$\text{where } \text{cnst} := 8L_{\max}\gamma^2 \sum_{i=1}^{m-1} (1 - \gamma\mu)^i$$

Free to choose the number of inner iterates m

Corollary If $\gamma = 1/6L_{\max}$ and $m = n$

$$t = O\left(\frac{6}{m} \frac{L_{\max}}{\mu}\right) \log\left(\frac{1}{\epsilon}\right) \Rightarrow \frac{\mathbb{E}[\|x_t^m - w^*\|^2]}{\Psi(x_0^0, \tilde{w}_0)} \leq \epsilon$$

Convergence SAGA

Theorem SAGA

If $f(w)$ is μ -strongly convex, $f_i(w)$ is L_{\max} -smooth and $\alpha = 1/(3L_{\max})$ then

$$\mathbb{E} [\|w^t - w^*\|_2^2] \leq \left(1 - \min \left\{ \frac{1}{4n}, \frac{\mu}{3L_{\max}} \right\}\right)^t C_0$$

where $C_0 = \frac{2n}{3L_{\max}}(f(w^0) - f(w^*)) + \|w^0 - w^*\|_2^2 \geq 0$

An even more practical convergence result!

Much easier proof due to unbiased gradients



A. Defazio, F. Bach and J. Lacoste-Julien (2014)
NIPS, **SAGA: A Fast Incremental Gradient Method With Support for Non-Strongly Convex Composite Objectives.**

Convergence SAG

Theorem SAG

If $f(w)$ is μ -strongly convex, $f_i(w)$ is L_{\max} -smooth and $\alpha = 1/(16L_{\max})$ then

$$\mathbb{E} [\|w^t - w^*\|_2^2] \leq \left(1 - \min \left\{ \frac{1}{8n}, \frac{\mu}{16L_{\max}} \right\}\right)^t C_0$$

where $C_0 = \frac{3}{2}(f(w^0) - f(w^*)) + \frac{4L_{\max}}{n} \|w^0 - w^*\|_2^2 \geq 0$

A practical convergence result!

Because of biased gradients, difficult proof that relies on computer assisted steps



M. Schmidt, N. Le Roux, F. Bach (2016)
Mathematical Programming
Minimizing Finite Sums with the Stochastic Average Gradient.

Comparisons in total complexity for strongly convex

Approximate solution

$$\mathbb{E}[f(w^T)] - f(w^*) \leq \epsilon \quad \text{or} \quad \mathbb{E}\|w^t - w^*\|^2 \leq \epsilon$$

SGD

$$O\left(\frac{1}{\epsilon}\right)$$

Gradient descent

$$O\left(\frac{nL}{\mu} \log\left(\frac{1}{\epsilon}\right)\right)$$

SVRG/SAGA/SAG

$$O\left(\left(n + \frac{L_{\max}}{\mu}\right) \log\left(\frac{1}{\epsilon}\right)\right)$$

Variance reduction faster than GD when

$$L \geq \mu + L_{\max}/n$$

How did I get these complexity results from the convergence results?



Section 1.3.5, R.M. Gower, Ph.d thesis: Sketch and Project: Randomized Iterative Methods for Linear Systems and Inverting Matrices University of Edinburgh, 2016

Practicals implementation of SAG for Linear Classifiers

Finite Sum Training Problem

$$\min_{w \in \mathbf{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(\langle w, x^i \rangle, y^i) + \frac{\lambda}{2} \|w\|_2^2$$

L2 regularizer +
linear hypothesis

Practicals implementation of SAG for Linear Classifiers

Finite Sum Training Problem

L2 regularizer +
linear hypothesis

$$\min_{w \in \mathbf{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(\langle w, x^i \rangle, y^i) + \frac{\lambda}{2} \|w\|_2^2$$

$$\nabla f_i(w) = \ell'(\langle w, x^i \rangle, y^i) x^i + \lambda w$$

Practicals implementation of SAG for Linear Classifiers

Finite Sum Training Problem

L2 regularizer +
linear hypothesis

$$\min_{w \in \mathbf{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(\langle w, x^i \rangle, y^i) + \frac{\lambda}{2} \|w\|_2^2$$

$$\nabla f_i(w) = \underbrace{\ell'(\langle w, x^i \rangle, y^i) x^i}_{\text{Nonlinear in } w} + \underbrace{\lambda w}_{\text{Linear in } w}$$

Practicals implementation of SAG for Linear Classifiers

Finite Sum Training Problem

L2 regularizer + linear hypothesis

$$\min_{w \in \mathbf{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(\langle w, x^i \rangle, y^i) + \frac{\lambda}{2} \|w\|_2^2$$

$$\nabla f_i(w) = \underbrace{\ell'(\langle w, x^i \rangle, y^i) x^i}_{\text{Nonlinear in } w} + \underbrace{\lambda w}_{\text{Linear in } w}$$

Reduce
Storage
to $O(n)$

Only store real number

Stoch. gradient estimate

Full gradient estimate

$$\beta_i = \ell'(\langle w^{t_i}, x^i \rangle, y^i)$$

$$\nabla f_i(w^{t_i}) = \beta_i x^i + \lambda w^t$$

$$g^t = \frac{1}{n} \sum_{j=1}^n \beta_j x_j + \lambda w^t$$

Proving Convergence of SVRG

Proof:

$$\begin{aligned}
\|x_s^{t+1} - w^*\|_2^2 &= \|x_s^t - w^* - \gamma g^t\|_2^2 \\
&= \|x_s^t - w^*\|_2^2 - 2\gamma \langle g^t, x_s^t - w^* \rangle + \gamma^2 \|g^t\|_2^2.
\end{aligned}$$

Taking expectation with respect to j

Unbiased estimator

$$\begin{aligned}
\mathbb{E}_j [\|x_s^{t+1} - w^*\|_2^2] &= \|x_s^t - w^*\|_2^2 - 2\gamma \langle \nabla f(x_s^t), x_s^t - w^* \rangle + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2] \\
&\stackrel{\text{str. conv.}}{\leq} (1 - \mu\gamma) \|x_s^t - w^*\|_2^2 - 2\gamma (f(x_s^t) - f(w^*)) + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2]
\end{aligned}$$

Need to
bound this!

$$\mathbb{E}_j [\|g^t\|_2^2]$$

Smoothness Consequences I

Smoothness

$$f(w) \leq f(y) + \langle \nabla f(y), w - y \rangle + \frac{L}{2} \|w - y\|_2^2, \quad \text{for } i = 1, \dots, n$$

EXE: Lemma 1

$$f(y - \frac{1}{L} \nabla f(y)) - f(y) \leq -\frac{1}{2L} \|\nabla f(y)\|_2^2, \quad \forall y.$$

Proof:

Substituting $w = y - \frac{1}{L} \nabla f(y)$ into the smoothness inequality gives

$$\begin{aligned} f(y - \frac{1}{L} \nabla f(y)) - f(y) &\leq \langle \nabla f(y), -\frac{1}{L} \nabla f(y) \rangle + \frac{L}{2} \left\| -\frac{1}{L} \nabla f(y) \right\|_2^2 \\ &= -\frac{1}{2L} \|\nabla f(y)\|_2^2. \quad \blacksquare \end{aligned}$$

Smoothness Consequences II

Smoothness

$$f_i(w) \leq f_i(y) + \langle \nabla f_i(y), w - y \rangle + \frac{L_i}{2} \|w - y\|_2^2, \quad \text{for } i = 1, \dots, n$$

EXE: Lemma 2

$$\mathbb{E}[\|\nabla f_i(w) - \nabla f_i(w^*)\|_2^2] \leq 2L_{\max}(f(w) - f(w^*))$$

Proof: Let $g_i(w) = f_i(w) - f_i(w^*) - \langle \nabla f_i(w^*), w - w^* \rangle$ which is L_i -smooth.

Smoothness Consequences II

Smoothness

$$f_i(w) \leq f_i(y) + \langle \nabla f_i(y), w - y \rangle + \frac{L_i}{2} \|w - y\|_2^2, \quad \text{for } i = 1, \dots, n$$

EXE: Lemma 2

$$\mathbb{E}[\|\nabla f_i(w) - \nabla f_i(w^*)\|_2^2] \leq 2L_{\max}(f(w) - f(w^*))$$

Proof: Let $g_i(w) = f_i(w) - f_i(w^*) - \langle \nabla f_i(w^*), w - w^* \rangle$ which is L_i -smooth.

Smoothness Consequences II

Smoothness

$$f_i(w) \leq f_i(y) + \langle \nabla f_i(y), w - y \rangle + \frac{L_i}{2} \|w - y\|_2^2, \quad \text{for } i = 1, \dots, n$$

EXE: Lemma 2

$$\mathbb{E}[\|\nabla f_i(w) - \nabla f_i(w^*)\|_2^2] \leq 2L_{\max}(f(w) - f(w^*))$$

Proof: Let $g_i(w) = f_i(w) - f_i(w^*) - \langle \nabla f_i(w^*), w - w^* \rangle$ which is L_i -smooth.

Convexity of $f_i(w) \Rightarrow g_i(w) \geq 0$ for all w . From Lemma 1 we have

$$g_i(w) \geq g_i(w) - g_i(w - \frac{1}{L_i} \nabla g_i(w)) \geq \frac{1}{2L_i} \|\nabla g_i(w)\|_2^2 \geq \frac{1}{2L_{\max}} \|\nabla g_i(w)\|_2^2$$

Inserting definition of $g_i(w)$ we have

$$\frac{1}{2L_{\max}} \|\nabla f_i(w) - \nabla f_i(w^*)\|_2^2 \leq f_i(w) - f_i(w^*) - \langle \nabla f_i(w^*), w - w^* \rangle$$

Result follows by taking expectation of i .

Lemma 1

Bounding gradient estimate

$$g^t = \nabla f_i(x^t) - \nabla f_i(\tilde{w}) + \nabla f(\tilde{w})$$

EXE: Lemma 3

$$\mathbb{E}[\|g^t\|_2^2] \leq 4L_{\max}(f(x^t) - f(w^*)) + 4L_{\max}(f(\tilde{w}) - f(w^*))$$

Proof: Hint: use $\|a + b\|_2^2 \leq 2\|a\|_2^2 + 2\|b\|_2^2$ and Lemma 2

Bounding gradient estimate

$$g^t = \nabla f_i(x^t) - \nabla f_i(\tilde{w}) + \nabla f(\tilde{w})$$

EXE: Lemma 3

$$\mathbb{E}[\|g^t\|_2^2] \leq 4L_{\max}(f(x^t) - f(w^*)) + 4L_{\max}(f(\tilde{w}) - f(w^*))$$

Proof: Hint: use $\|a + b\|_2^2 \leq 2\|a\|_2^2 + 2\|b\|_2^2$ and Lemma 2

Bounding gradient estimate

$$g^t = \nabla f_i(x^t) - \nabla f_i(\tilde{w}) + \nabla f(\tilde{w})$$

EXE: Lemma 3

$$\mathbb{E}[\|g^t\|_2^2] \leq 4L_{\max}(f(x^t) - f(w^*)) + 4L_{\max}(f(\tilde{w}) - f(w^*))$$

Proof: Hint: use $\|a + b\|_2^2 \leq 2\|a\|_2^2 + 2\|b\|_2^2$ and Lemma 2

$$\begin{aligned} \mathbb{E}_j[\|g^t\|_2^2] &= \mathbb{E}_j[\|\nabla f_i(x^t) - \nabla f_i(w^*) + \nabla f_i(w^*) - \nabla f_i(\tilde{w}) + \nabla f(\tilde{w})\|_2^2] \\ &\leq 2\mathbb{E}_j[\|\nabla f_i(x^t) - \nabla f_i(w^*)\|_2^2] + 2\mathbb{E}_j[\|\nabla f_i(w^*) - \nabla f_i(\tilde{w}) + \nabla f(\tilde{w})\|_2^2] \\ &\leq 2\mathbb{E}_j[\|\nabla f_i(x^t) - \nabla f_i(w^*)\|_2^2] + 2\mathbb{E}_j[\|\nabla f_i(w^*) - \nabla f_i(\tilde{w})\|_2^2] \\ &= 4L_{\max}(f(x^t) - f(w^*) + f(\tilde{w}) - f(w^*)) \quad \blacksquare \end{aligned}$$

Lemma 2

Bounding gradient estimate

$$g^t = \nabla f_i(x^t) - \nabla f_i(\tilde{w}) + \nabla f(\tilde{w})$$

EXE: Lemma 3

$$\mathbb{E}[\|g^t\|_2^2] \leq 4L_{\max}(f(x^t) - f(w^*)) + 4L_{\max}(f(\tilde{w}) - f(w^*))$$

Proof: Hint: use $\|a + b\|_2^2 \leq 2\|a\|_2^2 + 2\|b\|_2^2$ and Lemma 2

$$\begin{aligned} \mathbb{E}_j[\|g^t\|_2^2] &= \mathbb{E}_j[\|\nabla f_i(x^t) - \nabla f_i(w^*) + \nabla f_i(w^*) - \nabla f_i(\tilde{w}) + \nabla f(\tilde{w})\|_2^2] \\ &\leq 2\mathbb{E}_j[\|\nabla f_i(x^t) - \nabla f_i(w^*)\|_2^2] + 2\mathbb{E}_j[\|\nabla f_i(w^*) - \nabla f_i(\tilde{w}) + \nabla f(\tilde{w})\|_2^2] \\ &\leq 2\mathbb{E}_j[\|\nabla f_i(x^t) - \nabla f_i(w^*)\|_2^2] + 2\mathbb{E}_j[\|\nabla f_i(w^*) - \nabla f_i(\tilde{w})\|_2^2] \\ &= 4L_{\max}(f(x^t) - f(w^*) + f(\tilde{w}) - f(w^*)) \quad \blacksquare \end{aligned}$$

Lemma 2

Bounding gradient estimate

$$g^t = \nabla f_i(x^t) - \nabla f_i(\tilde{w}) + \nabla f(\tilde{w})$$

EXE: Lemma 3

$$\mathbb{E}[\|g^t\|_2^2] \leq 4L_{\max}(f(x^t) - f(w^*)) + 4L_{\max}(f(\tilde{w}) - f(w^*))$$

Proof: Hint: use $\|a + b\|_2^2 \leq 2\|a\|_2^2 + 2\|b\|_2^2$ and Lemma 2

$$\begin{aligned} \mathbb{E}_j[\|g^t\|_2^2] &= \mathbb{E}_j[\|\nabla f_i(x^t) - \nabla f_i(w^*) + \nabla f_i(w^*) - \nabla f_i(\tilde{w}) + \nabla f(\tilde{w})\|_2^2] \\ &\leq 2\mathbb{E}_j[\|\nabla f_i(x^t) - \nabla f_i(w^*)\|_2^2] + 2\mathbb{E}_j[\|\nabla f_i(w^*) - \nabla f_i(\tilde{w}) + \nabla f(\tilde{w})\|_2^2] \\ &\leq 2\mathbb{E}_j[\|\nabla f_i(x^t) - \nabla f_i(w^*)\|_2^2] + 2\mathbb{E}_j[\|\nabla f_i(w^*) - \nabla f_i(\tilde{w})\|_2^2] \\ &= 4L_{\max}(f(x^t) - f(w^*) + f(\tilde{w}) - f(w^*)) \quad \blacksquare \end{aligned}$$

Lemma 2

Where we used in the first inequality that $\mathbb{E}[\|X - \mathbb{E}X\|_2^2] \leq \mathbb{E}[\|X\|_2^2]$ with $X = \nabla f_i(w^*) - \nabla f_i(\tilde{w})$ thus $\mathbb{E}[X] = -\nabla f(\tilde{w})$

Proof:

$$\begin{aligned}
\|x_s^{t+1} - w^*\|_2^2 &= \|x_s^t - w^* - \gamma g^t\|_2^2 \\
&= \|x_s^t - w^*\|_2^2 - 2\gamma \langle g^t, x_s^t - w^* \rangle + \gamma^2 \|g^t\|_2^2.
\end{aligned}$$

Taking expectation with respect to j

Unbiased estimator

$$\begin{aligned}
\mathbb{E}_j [\|x_s^{t+1} - w^*\|_2^2] &= \|x_s^t - w^*\|_2^2 - 2\gamma \langle \nabla f(x_s^t), x_s^t - w^* \rangle + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2] \\
&\stackrel{\text{str. conv.}}{\leq} (1 - \mu\gamma) \|x_s^t - w^*\|_2^2 - 2\gamma (f(x_s^t) - f(w^*)) + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2]
\end{aligned}$$

Need to
bound this!

$$\mathbb{E}_j [\|g^t\|_2^2]$$

Lemma 3 $g^t = \nabla f_{\mathbf{i}}(x_s^t) - \nabla f_{\mathbf{i}}(\tilde{w}_{s-1}) + \nabla f(\tilde{w}_{s-1})$

$$\mathbb{E}_j [\|g^t\|_2^2] \leq 4L_{\max}(f(x_s^t) - f(w^*)) + 4L_{\max}(f(\tilde{w}_{s-1}) - f(w^*))$$

Proof:

$$\begin{aligned}
\|x_s^{t+1} - w^*\|_2^2 &= \|x_s^t - w^* - \gamma g^t\|_2^2 \\
&= \|x_s^t - w^*\|_2^2 - 2\gamma \langle g^t, x_s^t - w^* \rangle + \gamma^2 \|g^t\|_2^2.
\end{aligned}$$

Taking expectation with respect to j

Unbiased estimator



$$\begin{aligned}
\mathbb{E}_j [\|x_s^{t+1} - w^*\|_2^2] &= \|x_s^t - w^*\|_2^2 - 2\gamma \langle \nabla f(x_s^t), x_s^t - w^* \rangle + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2] \\
&\stackrel{\text{str. conv.}}{\leq} (1 - \mu\gamma) \|x_s^t - w^*\|_2^2 - 2\gamma (f(x_s^t) - f(w^*)) + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2] \\
&\leq (1 - \mu\gamma) \|x_s^t - w^*\|_2^2 - 2\gamma (1 - 2\gamma L_{\max}) (f(x_s^t) - f(w^*)) \\
&\quad + 4\gamma^2 L_{\max} (f(w_{s-1}) - f(w^*))
\end{aligned}$$

Proof:

$$\begin{aligned}
\|x_s^{t+1} - w^*\|_2^2 &= \|x_s^t - w^* - \gamma g^t\|_2^2 \\
&= \|x_s^t - w^*\|_2^2 - 2\gamma \langle g^t, x_s^t - w^* \rangle + \gamma^2 \|g^t\|_2^2.
\end{aligned}$$

Taking expectation with respect to j

Unbiased estimator



$$\begin{aligned}
\mathbb{E}_j [\|x_s^{t+1} - w^*\|_2^2] &= \|x_s^t - w^*\|_2^2 - 2\gamma \langle \nabla f(x_s^t), x_s^t - w^* \rangle + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2] \\
&\stackrel{\text{str. conv.}}{\leq} (1 - \mu\gamma) \|x_s^t - w^*\|_2^2 - 2\gamma (f(x_s^t) - f(w^*)) + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2] \\
&\leq (1 - \mu\gamma) \|x_s^t - w^*\|_2^2 - 2\gamma (1 - 2\gamma L_{\max}) (f(x_s^t) - f(w^*)) \\
&\quad + 4\gamma^2 L_{\max} (f(w_{s-1}) - f(w^*))
\end{aligned}$$

Taking expectation and iterating from $t = 0, \dots, m-1$

Proof:

$$\begin{aligned}
\|x_s^{t+1} - w^*\|_2^2 &= \|x_s^t - w^* - \gamma g^t\|_2^2 \\
&= \|x_s^t - w^*\|_2^2 - 2\gamma \langle g^t, x_s^t - w^* \rangle + \gamma^2 \|g^t\|_2^2.
\end{aligned}$$

Taking expectation with respect to j

Unbiased estimator

$$\begin{aligned}
\mathbb{E}_j [\|x_s^{t+1} - w^*\|_2^2] &= \|x_s^t - w^*\|_2^2 - 2\gamma \langle \nabla f(x_s^t), x_s^t - w^* \rangle + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2] \\
&\stackrel{\text{str. conv.}}{\leq} (1 - \mu\gamma) \|x_s^t - w^*\|_2^2 - 2\gamma (f(x_s^t) - f(w^*)) + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2] \\
&\leq (1 - \mu\gamma) \|x_s^t - w^*\|_2^2 - 2\gamma (1 - 2\gamma L_{\max}) (f(x_s^t) - f(w^*)) \\
&\quad + 4\gamma^2 L_{\max} (f(w_{s-1}) - f(w^*))
\end{aligned}$$

Taking expectation and iterating from $t = 0, \dots, m-1$

$$\begin{aligned}
\mathbb{E}_j [\|x_s^m - w^*\|_2^2] &\leq (1 - \mu\gamma)^m \|x_s^0 - w^*\|_2^2 \\
&\quad - 2\gamma (1 - 2\gamma L_{\max}) S_m \sum_{t=0}^{m-1} \alpha_t (f(x_s^t) - f(w^*)) \\
&\quad + 4S_m \gamma^2 L_{\max} (f(w_{s-1}) - f(w^*))
\end{aligned}$$

$$\alpha_t := (1 - \mu\gamma)^{m-1-t}$$

Proof:

$$\begin{aligned}
\|x_s^{t+1} - w^*\|_2^2 &= \|x_s^t - w^* - \gamma g^t\|_2^2 \\
&= \|x_s^t - w^*\|_2^2 - 2\gamma \langle g^t, x_s^t - w^* \rangle + \gamma^2 \|g^t\|_2^2.
\end{aligned}$$

Unbiased estimator

Taking expectation with respect to j

$$\begin{aligned}
\mathbb{E}_j [\|x_s^{t+1} - w^*\|_2^2] &= \|x_s^t - w^*\|_2^2 - 2\gamma \langle \nabla f(x_s^t), x_s^t - w^* \rangle + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2] \\
&\stackrel{\text{str. conv.}}{\leq} (1 - \mu\gamma) \|x_s^t - w^*\|_2^2 - 2\gamma (f(x_s^t) - f(w^*)) + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2] \\
&\leq (1 - \mu\gamma) \|x_s^t - w^*\|_2^2 - 2\gamma (1 - 2\gamma L_{\max}) (f(x_s^t) - f(w^*)) \\
&\quad + 4\gamma^2 L_{\max} (f(w_{s-1}) - f(w^*))
\end{aligned}$$

Taking expectation and iterating from $t = 0, \dots, m-1$

$$\begin{aligned}
\mathbb{E}_j [\|x_s^m - w^*\|_2^2] &\leq (1 - \mu\gamma)^m \|x_s^0 - w^*\|_2^2 \\
&\quad - 2\gamma (1 - 2\gamma L_{\max}) S_m \sum_{t=0}^{m-1} \alpha_t (f(x_s^t) - f(w^*)) \\
&\quad + 4S_m \gamma^2 L_{\max} (f(w_{s-1}) - f(w^*))
\end{aligned}$$

$$S_m := \sum_{t=0}^{m-1} \alpha_t$$

$$\alpha_t := (1 - \mu\gamma)^{m-1-t}$$

Proof:

$$\begin{aligned}
\|x_s^{t+1} - w^*\|_2^2 &= \|x_s^t - w^* - \gamma g^t\|_2^2 \\
&= \|x_s^t - w^*\|_2^2 - 2\gamma \langle g^t, x_s^t - w^* \rangle + \gamma^2 \|g^t\|_2^2.
\end{aligned}$$

Unbiased estimator

Taking expectation with respect to j

$$\begin{aligned}
\mathbb{E}_j [\|x_s^{t+1} - w^*\|_2^2] &= \|x_s^t - w^*\|_2^2 - 2\gamma \langle \nabla f(x_s^t), x_s^t - w^* \rangle + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2] \\
&\stackrel{\text{str. conv.}}{\leq} (1 - \mu\gamma) \|x_s^t - w^*\|_2^2 - 2\gamma (f(x_s^t) - f(w^*)) + \gamma^2 \mathbb{E}_j [\|g^t\|_2^2] \\
&\leq (1 - \mu\gamma) \|x_s^t - w^*\|_2^2 - 2\gamma (1 - 2\gamma L_{\max}) (f(x_s^t) - f(w^*)) \\
&\quad + 4\gamma^2 L_{\max} (f(w_{s-1}) - f(w^*))
\end{aligned}$$

Taking expectation and iterating from $t = 0, \dots, m-1$

$$\mathbb{E}_j [\|x_s^m - w^*\|_2^2] \leq (1 - \mu\gamma)^m \|x_s^0 - w^*\|_2^2$$

$$\alpha_t := (1 - \mu\gamma)^{m-1-t}$$

$$S_m := \sum_{t=0}^{m-1} \alpha_t$$

$$\begin{aligned}
&-2\gamma (1 - 2\gamma L_{\max}) S_m \sum_{t=0}^{m-1} \alpha_t (f(x_s^t) - f(w^*)) \\
&+ 4S_m \gamma^2 L_{\max} (f(w_{s-1}) - f(w^*))
\end{aligned}$$

Rest on the board

Take for home Variance Reduction

- Variance reduced methods use only **one stochastic gradient per iteration** and converge linearly on strongly convex functions
- Choice of **fixed stepsize** possible
- **SAGA** only needs to know the smoothness parameter to work, but requires storing n past stochastic gradients
- **SVRG** only has $O(d)$ storage, but requires full gradient computations every so often. Has an extra “number of inner iterations” parameter to tune