SARAH Algorithm

Lam M. Nguyen, Jie Liu, Katya Scheinberg, Martin Takáč

INFORMS Annual Meeting October 24, 2017



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

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

Training set: $\{(x_i, y_i)\}_{i=1}^n$ with $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$

 f_i - strongly convex: linear regression, binary classification

 ℓ_2 -regularized least squares regression: $f_i(w) = (x_i^T w - y_i)^2 + \frac{\lambda}{2} ||w||^2$

 ℓ_2 -regularized logistic regression: $f_i(w) = \log(1 + \exp(-y_i x_i^T w)) + \frac{\lambda}{2} ||w||^2$

 f_i - **nonconvex**: neural networks

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

Training set: $\{(x_i, y_i)\}_{i=1}^n$ with $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$

 f_i - strongly convex: linear regression, binary classification

 ℓ_2 -regularized least squares regression: $f_i(w) = (x_i^T w - y_i)^2 + \frac{\lambda}{2} ||w||^2$

 ℓ_2 -regularized logistic regression: $f_i(w) = \log(1 + \exp(-y_i x_i^T w)) + \frac{\lambda}{2} ||w||^2$

 f_i - **nonconvex**: neural networks

HOW WE COULD SOLVE THIS OPTIMIZATION PROBLEM?

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

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

P is *L*-smooth and μ -strongly convex

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

P is *L*-smooth and μ -strongly convex

Iterative methods (using gradient)

Given initial point w_0 . Update: $w_{t+1} = w_t - \eta_t v_t$, t = 0,1,2,...

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

P is *L*-smooth and μ -strongly convex

Iterative methods (using gradient)

Given initial point w_0 . Update: $w_{t+1} = w_t - \eta_t v_t$, t = 0,1,2,...

Goal: achieve ϵ -accurate solution w_T such that

- $P(w_T) P(w^*) \le \epsilon$
- $||w_T w^*||^2 \le \epsilon$
- $\|\nabla P(w_T)\|^2 \le \epsilon$

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

P is *L*-smooth and μ -strongly convex

Iterative methods (using gradient)

Given initial point w_0 . Update: $w_{t+1} = w_t - \eta_t v_t$, t = 0,1,2,...

Goal: achieve ϵ -accurate solution w_T such that

- $P(w_T) P(w^*) \le \epsilon$
- $||w_T w^*||^2 \le \epsilon$
- $\|\nabla P(w_T)\|^2 \le \epsilon$

Gradient Descent: $w_{t+1} = w_t - \eta \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$

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

P is *L*-smooth and μ -strongly convex

Iterative methods (using gradient)

Given initial point w_0 . Update: $w_{t+1} = w_t - \eta_t v_t$, t = 0,1,2,...

Goal: achieve ϵ -accurate solution w_T such that

• $P(w_T) - P(w^*) \le \epsilon$

•
$$||w_T - w^*||^2 \le \epsilon$$

• $\|\nabla P(w_T)\|^2 \le \epsilon$

Gradient Descent: $w_{t+1} = w_t - \eta \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$

The total work complexity: Number of component gradient evaluations

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

P is *L*-smooth and μ -strongly convex

Iterative methods (using gradient)

Given initial point w_0 . Update: $w_{t+1} = w_t - \eta_t v_t$, t = 0,1,2,...

Goal: achieve ϵ -accurate solution w_T such that

•
$$P(w_T) - P(w^*) \le \epsilon$$

•
$$||w_T - w^*||^2 \le \epsilon$$

• $\|\nabla P(w_T)\|^2 \le \epsilon$

Need "*n*" work per iteration

Gradient Descent: $w_{t+1} = w_t - \eta \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$

The total work complexity: Number of component gradient evaluations

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

P is *L*-smooth and μ -strongly convex

Iterative methods (using gradient)

Given initial point w_0 . Update: $w_{t+1} = w_t - \eta_t v_t$, t = 0,1,2,...

Goal: achieve ϵ -accurate solution w_T such that



per iteration

The total work complexity: Number of component gradient evaluations

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

P is *L*-smooth and μ -strongly convex

Iterative methods (using gradient)

Given initial point w_0 . Update: $w_{t+1} = w_t - \eta_t v_t$, t = 0,1,2,...

Goal: achieve ϵ -accurate solution w_T such that

•
$$P(w_T) - P(w^*) \le \epsilon$$

•
$$\|w_T - w^*\|^2 \le \epsilon$$

•
$$\|\nabla P(w_T)\|^2 \le \epsilon$$

Need "*n*" work **per iteration**

Gradient Descent: $w_{t+1} = w_t - \eta \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$

The total work complexity: Number of component gradient evaluations Machine learning and Big Data applications $\Rightarrow n \gg 1$ ("*n*" is very large)

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

P is *L*-smooth and μ -strongly convex

Iterative methods (using gradient)

Given initial point w_0 . Update: $w_{t+1} = w_t - \eta_t v_t$, t = 0,1,2,...

Goal: achieve ϵ -accurate solution w_T such that

•
$$P(w_T) - P(w^*) \le \epsilon$$

•
$$||w_T - w^*||^2 \le \epsilon$$

•
$$\|\nabla P(w_T)\|^2 \le \epsilon$$

Need "*n*" work per iteration

Gradient Descent: $w_{t+1} = w_t - \eta \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$

The total work complexity: Number of component gradient evaluations Machine learning and Big Data applications $\Rightarrow n \gg 1$ ("*n*" is very large)

TOO MUCH WORK !!!

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

P is *L*-smooth and μ -strongly convex

Iterative methods (using gradient)

Given initial point w_0 . Update: $w_{t+1} = w_t - \eta_t v_t$, t = 0,1,2,...

Goal: achieve ϵ -accurate solution w_T such that

• $P(w_T) - P(w^*) \le \epsilon$

•
$$||w_T - w^*||^2 \le \epsilon$$

• $\|\nabla P(w_T)\|^2 \le \epsilon$

Need "*n*" work **per iteration**

Gradient Descent: $w_{t+1} = w_t - \eta \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$ **Newton Method**: $w_{t+1} = w_t - [\mathbf{H} P(w_t)]^{-1} \nabla P(w_t)$

The total work complexity: Number of component gradient evaluations Machine learning and Big Data applications $\Rightarrow n \gg 1$ ("*n*" is very large)

TOO MUCH WORK !!!

Computing "less" work per iteration

Stochastic Gradient Descent (SGD)

- 1. Choose initial point w_0
- 2. For t = 0, 1, 2, ...
- 3. $w_{t+1} = w_t \eta_t \nabla f_{i_t}(w_t), i_t \in \{1, \dots, n\}$

- 1. Choose initial point w_0
- 2. For t = 0, 1, 2, ...

Need only "1" work per iteration

3. $w_{t+1} = w_t - \eta_t \nabla f_{i_t}(w_t), i_t \in \{1, ..., n\}$

1. Choose initial point w_0

Need only "1" work per iteration

2. For t = 0, 1, 2, ...3. $w_{t+1} = w_t - \eta_t \nabla f_{i_t}(w_t), i_t \in \{1, ..., n\}$

To guarantee convergence:

$$\sum_{t=0}^{\infty} \eta_t = \infty \quad \text{and} \quad \sum_{t=0}^{\infty} \eta_t^2 < \infty$$

- 1. Choose initial point w_0
- 2. For t = 0, 1, 2, ...
- 3. $w_{t+1} = w_t \eta_t \nabla f_{i_t}(w_t), i_t \in \{1, \dots, n\}$

Need only "1" work per iteration

To guarantee convergence: $\sum_{t=0}^{\infty} \eta_t = \infty \quad \text{and} \quad \sum_{t=0}^{\infty} \eta_t^2 < \infty$ If $\eta_t = \frac{d}{\gamma + t}$ then $\mathbb{E}[P(w_t) - P(w^*)] \le \frac{c}{\gamma + t}$

- 1. Choose initial point w_0
- 2. For t = 0, 1, 2, ...
- 3. $w_{t+1} = w_t \eta_t \nabla f_{i_t}(w_t), i_t \in \{1, ..., n\}$

Need only "1" work per iteration

To guarantee convergence: $\sum_{t=0}^{\infty} \eta_t = \infty \quad \text{and} \quad \sum_{t=0}^{\infty} \eta_t^2 < \infty$ If $\eta_t = \frac{d}{\gamma + t}$ then $\mathbb{E}[P(w_t) - P(w^*)] \le \frac{c}{\gamma + t}$

=> Require $O\left(\frac{1}{\epsilon}\right)$ total work to achieve ϵ -accurate solution

- 1. Choose initial point w_0
- 2. For t = 0, 1, 2, ...
- 3. $w_{t+1} = w_t \eta_t \nabla f_{i_t}(w_t), i_t \in \{1, ..., n\}$

Need only "1" work per iteration

To guarantee convergence: $\sum_{t=0}^{\infty} \eta_t = \infty \quad \text{and} \quad \sum_{t=0}^{\infty} \eta_t^2 < \infty$ If $\eta_t = \frac{d}{\gamma + t}$ then $\mathbb{E}[P(w_t) - P(w^*)] \le \frac{c}{\gamma + t}$

=> Require $O\left(\frac{1}{\epsilon}\right)$ total work to achieve ϵ -accurate solution

Pros:

• Each iteration is independent on "*n*"

- 1. Choose initial point w_0
- 2. For t = 0, 1, 2, ...

To guarantee convergence.

3. $w_{t+1} = w_t - \eta_t \nabla f_{i_t}(w_t), i_t \in \{1, ..., n\}$

Need only "1" work per iteration

If
$$\eta_t = \frac{d}{\gamma + t}$$
 then $\mathbb{E}[P(w_t) - P(w^*)] \le \frac{c}{\gamma + t}$

=> Require $O\left(\frac{1}{\epsilon}\right)$ total work to achieve ϵ -accurate solution

Pros:

• Each iteration is independent on "*n*"

Cons:

• Sublinear convergence rate

- 1. Choose initial point w_0
- 2. For t = 0, 1, 2, ...

To guarantee convergence.

3. $w_{t+1} = w_t - \eta_t \nabla f_{i_t}(w_t), i_t \in \{1, ..., n\}$

Need only "1" work per iteration

If
$$\eta_t = \frac{d}{\gamma + t}$$
 then $\mathbb{E}[P(w_t) - P(w^*)] \le \frac{c}{\gamma + t}$

=> Require $O\left(\frac{1}{\epsilon}\right)$ total work to achieve ϵ -accurate solution

Pros:

• Each iteration is independent on "*n*"

Cons:

• Sublinear convergence rate

Can we get a linear convergence rate?

- 1. Choose initial point w_0
- 2. For t = 0, 1, 2, ...

To guarantee convergence.

3. $w_{t+1} = w_t - \eta_t \nabla f_{i_t}(w_t), i_t \in \{1, ..., n\}$

Need only "1" work per iteration

If
$$\eta_t = \frac{d}{\gamma + t}$$
 then $\mathbb{E}[P(w_t) - P(w^*)] \le \frac{c}{\gamma + t}$

=> Require $O\left(\frac{1}{\epsilon}\right)$ total work to achieve ϵ -accurate solution

Pros:

• Each iteration is independent on "*n*"

Cons:

• Sublinear convergence rate



Can we get a linear convergence rate?

H. Robbins and S. Monro. A Stochastic Approximation Method, 1951 Léon Bottou, Frank E Curtis, Jorge Nocedal. Optimization methods for large-scale machine learning, 2016

SAG [M. Schmidt et. al., 2013] and SAGA [A. Defazio et. al., 2014]

- Compute the full gradient at the initial point
- Keep a table of "past" gradients
- In each iteration, update one "gradient" in the table

SAG [M. Schmidt et. al., 2013] and SAGA [A. Defazio et. al., 2014]

- Compute the full gradient at the initial point
- Keep a table of "past" gradients
- In each iteration, update one "gradient" in the table

(SAG)
$$w_{t+1} = w_t - \eta_t \cdot \frac{1}{n} \sum_{i=1}^n y_{i,t} \qquad y_{i,t} = \begin{cases} \nabla f_i(w_t), & \text{if } i_t = i \\ y_{i,t-1}, & \text{otherwise} \end{cases}$$

(SAGA)
$$w_{t+1} = w_t - \eta \left(\nabla f_{i_t}(w_t) - y_{i_t,t-1} + \frac{1}{n} \sum_{i=1}^n y_{i,t-1} \right)$$

SAG [M. Schmidt et. al., 2013] and SAGA [A. Defazio et. al., 2014]

- Compute the full gradient at the initial point
- Keep a table of "past" gradients
- In each iteration, update one "gradient" in the table

(SAG)
$$w_{t+1} = w_t - \eta_t \cdot \frac{1}{n} \sum_{i=1}^n y_{i,t} \qquad y_{i,t} = \begin{cases} \nabla f_i(w_t), & \text{if } i_t = i \\ y_{i,t-1}, & \text{otherwise} \end{cases}$$

(SAGA)
$$w_{t+1} = w_t - \eta \left(\nabla f_{i_t}(w_t) - y_{i_t,t-1} + \frac{1}{n} \sum_{i=1}^n y_{i,t-1} \right)$$

Pros:

• Linear convergence rate

SAG [M. Schmidt et. al., 2013] and SAGA [A. Defazio et. al., 2014]

- Compute the full gradient at the initial point
- Keep a table of "past" gradients
- In each iteration, update one "gradient" in the table

(SAG)
$$w_{t+1} = w_t - \eta_t \cdot \frac{1}{n} \sum_{i=1}^n y_{i,t} \qquad y_{i,t} = \begin{cases} \nabla f_i(w_t), & \text{if } i_t = i \\ y_{i,t-1}, & \text{otherwise} \end{cases}$$

(SAGA)
$$w_{t+1} = w_t - \eta \left(\nabla f_{i_t}(w_t) - y_{i_t,t-1} + \frac{1}{n} \sum_{i=1}^n y_{i,t-1} \right)$$

Pros:

Cons:

• Linear convergence rate

• Extra storage! Need to store "*n*" gradients

M. Schmidt, N. Le Roux, F. Bach. Minimizing Finite Sums with the Stochastic Average Gradient, 2013 A. Defazio, F. Bach, S. Lacoste-Julien. SAGA: A Fast Incremental Gradient Method With Support for Non-Strongly Convex Composite Objectives, 2014

SAG [M. Schmidt et. al., 2013] and SAGA [A. Defazio et. al., 2014]

- Compute the full gradient at the initial point
- Keep a table of "past" gradients
- In each iteration, update one "gradient" in the table

(SAG)
$$w_{t+1} = w_t - \eta_t \cdot \frac{1}{n} \sum_{i=1}^n y_{i,t} \qquad y_{i,t} = \begin{cases} \nabla f_i(w_t), & \text{if } i_t = i \\ y_{i,t-1}, & \text{otherwise} \end{cases}$$

(SAGA)
$$w_{t+1} = w_t - \eta \left(\nabla f_{i_t}(w_t) - y_{i_t,t-1} + \frac{1}{n} \sum_{i=1}^n y_{i,t-1} \right)$$

Pros:

• Linear convergence rate

Cons:

• Extra storage! Need to store "*n*" gradients

Can we eliminate the extra storage and get a linear convergence rate?

M. Schmidt, N. Le Roux, F. Bach. Minimizing Finite Sums with the Stochastic Average Gradient, 2013

A. Defazio, F. Bach, S. Lacoste-Julien. SAGA: A Fast Incremental Gradient Method With Support for Non-Strongly Convex Composite Objectives, 2014

SVRG [R. Johnson & T. Zhang, 2013]

- Modify stochastic gradient
 - 1. Choose initial point w_0
 - 2. Set $\widetilde{w} = w_0$
 - 3. For t = 0, 1, 2, ..., m

4.
$$w_{t+1} = w_t - \eta \underbrace{\left(\nabla f_{i_t}(w_t) - \nabla f_{i_t}(\widetilde{w}) + \nabla P(\widetilde{w})\right)}_{v_t}, i_t \in \{1, \dots, n\}$$

SVRG [R. Johnson & T. Zhang, 2013]

- Modify stochastic gradient
 - 1. Choose initial point w_0
 - 2. Set $\widetilde{w} = w_0$
 - 3. For t = 0, 1, 2, ..., m

4.
$$w_{t+1} = w_t - \eta \underbrace{\left(\nabla f_{i_t}(w_t) - \nabla f_{i_t}(\widetilde{w}) + \nabla P(\widetilde{w})\right)}_{v_t}, i_t \in \{1, \dots, n\}$$

• Let
$$\widetilde{w}^+ \in \{w_0, w_1, \dots, w_{m-1}\}$$

• Choose $\eta < \frac{1}{4L}$ and *m* such that $\alpha \coloneqq \frac{1}{\mu \eta (1 - 2L\eta)m} + \frac{2L\eta}{1 - 2L\eta} < 1$

SVRG [R. Johnson & T. Zhang, 2013]

- Modify stochastic gradient
 - 1. Choose initial point w_0
 - 2. Set $\widetilde{w} = w_0$
 - 3. For t = 0, 1, 2, ..., m

4.
$$w_{t+1} = w_t - \eta \underbrace{\left(\nabla f_{i_t}(w_t) - \nabla f_{i_t}(\widetilde{w}) + \nabla P(\widetilde{w})\right)}_{v_t}, i_t \in \{1, \dots, n\}$$

• Let
$$\widetilde{w}^+ \in \{w_0, w_1, \dots, w_{m-1}\}$$

• Choose $\eta < \frac{1}{4L}$ and *m* such that $\alpha \coloneqq \frac{1}{\mu n(1 - 2Ln)m} + \frac{2L\eta}{1 - 2Ln} < 1$

Then $\mathbb{E}[P(\widetilde{w}^+) - P(w^*)] \le \alpha \cdot \mathbb{E}[P(\widetilde{w}) - P(w^*)]$

SVRG [R. Johnson & T. Zhang, 2013]

- Modify stochastic gradient
 - 1. Choose initial point w_0
 - 2. Set $\widetilde{w} = w_0$
 - 3. For t = 0, 1, 2, ..., m

4.
$$w_{t+1} = w_t - \eta \underbrace{\left(\nabla f_{i_t}(w_t) - \nabla f_{i_t}(\widetilde{w}) + \nabla P(\widetilde{w})\right)}_{v_t}, i_t \in \{1, \dots, n\}$$

• Let
$$\widetilde{w}^+ \in \{w_0, w_1, \dots, w_{m-1}\}$$

• Choose $\eta < \frac{1}{4L}$ and *m* such that $\alpha \coloneqq \frac{1}{\mu\eta(1 - 2L\eta)m} + \frac{2L\eta}{1 - 2L\eta} < 1$

Then $\mathbb{E}[P(\widetilde{w}^+) - P(w^*)] \le \alpha \cdot \mathbb{E}[P(\widetilde{w}) - P(w^*)]$

Note: For fixed η , it would not converge to the optimal solution!
SVRG [R. Johnson & T. Zhang, 2013]

- Modify stochastic gradient
 - 1. Choose initial point w_0
 - 2. Set $\widetilde{w} = w_0$
 - 3. For t = 0, 1, 2, ..., m

4.
$$w_{t+1} = w_t - \eta \underbrace{\left(\nabla f_{i_t}(w_t) - \nabla f_{i_t}(\widetilde{w}) + \nabla P(\widetilde{w})\right)}_{v_t}, i_t \in \{1, \dots, n\}$$

• Let
$$\widetilde{w}^+ \in \{w_0, w_1, \dots, w_{m-1}\}$$

• Choose $\eta < \frac{1}{4L}$ and *m* such that $\alpha \coloneqq \frac{1}{\mu\eta(1 - 2L\eta)m} + \frac{2L\eta}{1 - 2L\eta} < 1$

Then $\mathbb{E}[P(\widetilde{w}^+) - P(w^*)] \le \alpha \cdot \mathbb{E}[P(\widetilde{w}) - P(w^*)]$

Note: For fixed η , it would not converge to the optimal solution!



SVRG [R. Johnson & T. Zhang, 2013]

- Modify stochastic gradient
 - 1. Choose initial point w_0
 - 2. Set $\widetilde{w} = w_0$
 - 3. For t = 0, 1, 2, ..., m

4.
$$w_{t+1} = w_t - \eta \underbrace{\left(\nabla f_{i_t}(w_t) - \nabla f_{i_t}(\widetilde{w}) + \nabla P(\widetilde{w})\right)}_{v_t}, i_t \in \{1, \dots, n\}$$

• Let
$$\widetilde{w}^+ \in \{w_0, w_1, \dots, w_{m-1}\}$$

• Choose $\eta < \frac{1}{4L}$ and *m* such that $\alpha \coloneqq \frac{1}{\mu\eta(1 - 2L\eta)m} + \frac{2L\eta}{1 - 2L\eta} < 1$

Then $\mathbb{E}[P(\widetilde{w}^+) - P(w^*)] \le \alpha \cdot \mathbb{E}[P(\widetilde{w}) - P(w^*)]$

Note: For fixed η , it would not converge to the optimal solution!

$$\widetilde{w}^{(0)} \to \widetilde{w}^{(1)} \to \cdots \to \widetilde{w}^{(s)}$$

SVRG [R. Johnson & T. Zhang, 2013]

- Modify stochastic gradient
 - 1. Choose initial point w_0
 - 2. Set $\widetilde{w} = w_0$
 - 3. For t = 0, 1, 2, ..., m

4.
$$w_{t+1} = w_t - \eta \underbrace{\left(\nabla f_{i_t}(w_t) - \nabla f_{i_t}(\widetilde{w}) + \nabla P(\widetilde{w})\right)}_{v_t}, i_t \in \{1, \dots, n\}$$

• Let
$$\widetilde{w}^+ \in \{w_0, w_1, \dots, w_{m-1}\}$$

• Choose $\eta < \frac{1}{4L}$ and *m* such that $\alpha \coloneqq \frac{1}{\mu\eta(1 - 2L\eta)m} + \frac{2L\eta}{1 - 2L\eta} < 1$

Then $\mathbb{E}[P(\widetilde{w}^+) - P(w^*)] \le \alpha \cdot \mathbb{E}[P(\widetilde{w}) - P(w^*)]$

Note: For fixed η , it would not converge to the optimal solution!

$$\widetilde{w}^{(0)} \to \widetilde{w}^{(1)} \to \dots \to \widetilde{w}^{(s)}$$

Hence, $\mathbb{E}\left[P(\widetilde{w}^{(s)}) - P(w^*)\right] \le \alpha^s \cdot \left[P(\widetilde{w}^{(0)}) - P(w^*)\right]$



SVRG [R. Johnson & T. Zhang, 2013]

- Modify stochastic gradient
 - 1. Choose initial point w_0
 - 2. Set $\widetilde{w} = w_0$
 - 3. For t = 0, 1, 2, ..., m

4.
$$w_{t+1} = w_t - \eta \underbrace{\left(\nabla f_{i_t}(w_t) - \nabla f_{i_t}(\widetilde{w}) + \nabla P(\widetilde{w})\right)}_{v_t}, i_t \in \{1, \dots, n\}$$

• Let
$$\widetilde{w}^+ \in \{w_0, w_1, \dots, w_{m-1}\}$$

• Choose $\eta < \frac{1}{4L}$ and *m* such that $\alpha \coloneqq \frac{1}{\mu\eta(1 - 2L\eta)m} + \frac{2L\eta}{1 - 2L\eta} < 1$

Then $\mathbb{E}[P(\widetilde{w}^+) - P(w^*)] \le \alpha \cdot \mathbb{E}[P(\widetilde{w}) - P(w^*)]$

Note: For fixed η , it would not converge to the optimal solution!

$$\widetilde{w}^{(0)} \to \widetilde{w}^{(1)} \to \cdots \to \widetilde{w}^{(s)}$$

Hence, $\mathbb{E}[P(\widetilde{w}^{(s)}) - P(w^*)] \le \alpha^s \cdot [P(\widetilde{w}^{(0)}) - P(w^*)]$

No storage is required!

SVRG One Outer Loop Behavior

An issue:









• How to choose "*m*" in algorithm?



The trajectory for one outer loop is very "unstable"

SARAH

- L. Nguyen, J. Liu, K. Scheinberg, and M. Takac. SARAH: A Novel Method for Machine Learning Problems Using Stochastic Recursive Gradient, 2017
- L. Nguyen, J. Liu, K. Scheinberg, and M. Takac. Stochastic Recursive Gradient Algorithm for Nonconvex Optimization, 2017

- It also does **restarting** as SVRG
- It takes **recursive** gradient estimator

- It also does **restarting** as SVRG
- It takes **recursive** gradient estimator

Parameters: the learning rate $\eta > 0$ and the inner loop size m. Initialize: \tilde{w}_0 Iterate: for s = 1, 2, ... do $w_0 = \tilde{w}_{s-1}$ $v_0 = \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(w_0)$ $w_1 = w_0 - \eta v_0$ Iterate: for t = 1, ..., m - 1 do Sample i_t uniformly at random from [n] $v_t = \nabla f_{i_t}(w_t) - \nabla f_{i_t}(w_{t-1}) + v_{t-1}$ $w_{t+1} = w_t - \eta v_t$ end for Set $\tilde{w}_s = w_t$ with t chosen uniformly at random from $\{0, 1, \dots, m\}$ end for

- It also does **restarting** as SVRG
- It takes **recursive** gradient estimator

Parameters: the learning rate $\eta > 0$ and the inner loop size m. Initialize: \tilde{w}_0 Iterate: for s = 1, 2, ... do $w_0 = \tilde{w}_{s-1}$ $v_0 = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_0)$ $w_1 = w_0 - \eta v_0$ Iterate: for t = 1, ..., m - 1 do Sample i_t uniformly at random from [n] $v_{t} = \nabla f_{i_{t}}(w_{t}) - \nabla f_{i_{t}}(w_{t-1}) + v_{t-1}$ $w_{t+1} = w_t - \eta v_t$ end for Set $\tilde{w}_s = w_t$ with t chosen uniformly at random from $\{0, 1, \ldots, m\}$ end for

Outer loop

- It also does **restarting** as SVRG
- It takes **recursive** gradient estimator



- It also does **restarting** as SVRG
- It takes **recursive** gradient estimator



- It also does **restarting** as SVRG
- It takes **recursive** gradient estimator



SARAH update (stochastic gradient computing)

- It also does **restarting** as SVRG
- It takes **recursive** gradient estimator



SARAH update (stochastic gradient computing) No storage is required!

Derivation of SARAH update

• Try to approximate **Gradient Descent**

Recall the update: $w_{t+1} = w_t - \eta v_t$

We want: $v_t \approx \nabla P(w_t)$

Recall the update: $w_{t+1} = w_t - \eta v_t$

We want: $\boldsymbol{v}_t \approx \boldsymbol{\nabla} \boldsymbol{P}(\boldsymbol{w}_t)$

According to L-Lipschitz smooth property, we have

$$\begin{aligned} ||\nabla P(w_t) - \nabla P(w_{t-1})|| &\leq L ||w_t - w_{t-1}|| = L\eta ||v_{t-1}|| \\ ||\nabla f_i(w_t) - \nabla f_i(w_{t-1})|| &\leq L ||w_t - w_{t-1}|| = L\eta ||v_{t-1}||, \forall i \in \{1, \dots, n\} \end{aligned}$$

Recall the update: $w_{t+1} = w_t - \eta v_t$

We want: $\boldsymbol{v}_t \approx \boldsymbol{\nabla} \boldsymbol{P}(\boldsymbol{w}_t)$

According to L-Lipschitz smooth property, we have

$$\begin{aligned} ||\nabla P(w_t) - \nabla P(w_{t-1})|| &\leq L ||w_t - w_{t-1}|| = L\eta ||v_{t-1}|| \\ ||\nabla f_i(w_t) - \nabla f_i(w_{t-1})|| &\leq L ||w_t - w_{t-1}|| = L\eta ||v_{t-1}||, \forall i \in \{1, \dots, n\} \end{aligned}$$

When η is small enough, we have

 $\nabla P(w_t) \approx \nabla f_i(w_t) - \nabla f_i(w_{t-1}) + \nabla P(w_{t-1})$

Recall the update: $w_{t+1} = w_t - \eta v_t$

We want: $\boldsymbol{v}_t \approx \boldsymbol{\nabla} \boldsymbol{P}(\boldsymbol{w}_t)$

According to L-Lipschitz smooth property, we have

$$\begin{aligned} ||\nabla P(w_t) - \nabla P(w_{t-1})|| &\leq L ||w_t - w_{t-1}|| = L\eta ||v_{t-1}|| \\ ||\nabla f_i(w_t) - \nabla f_i(w_{t-1})|| &\leq L ||w_t - w_{t-1}|| = L\eta ||v_{t-1}||, \forall i \in \{1, \dots, n\} \end{aligned}$$

When η is small enough, we have

$$\nabla P(w_t) \approx \nabla f_i(w_t) - \nabla f_i(w_{t-1}) + \nabla P(w_{t-1})$$
$$\Rightarrow v_t \approx \nabla f_i(w_t) - \nabla f_i(w_{t-1}) + v_{t-1}$$

 $\mathbb{E}[v_t | \mathcal{F}_t] = \nabla P(w_t)$

 $\mathbb{E}[v_t | \mathcal{F}_t] = \nabla P(w_t)$

• SARAH is conditionally biased

 $\mathbb{E}[v_t | \mathcal{F}_t] = \nabla P(w_t)$

• SARAH is conditionally biased

Recall: $v_t = \nabla f_{i_t}(w_t) - \nabla f_{i_t}(w_{t-1}) + v_{t-1}$

We have

$$\mathbb{E}[v_t | \mathcal{F}_t] = \nabla P(w_t) - \nabla P(w_{t-1}) + v_{t-1} \neq \nabla P(w_t)$$

Conditioned on $\{w_0, i_0, i_1, \dots, i_{t-1}\}$

 $\mathbb{E}[v_t | \mathcal{F}_t] = \nabla P(w_t)$

• SARAH is conditionally biased

Recall: $v_t = \nabla f_{i_t}(w_t) - \nabla f_{i_t}(w_{t-1}) + v_{t-1}$

We have

$$\mathbb{E}[v_t | \mathcal{F}_t] = \nabla P(w_t) - \nabla P(w_{t-1}) + v_{t-1} \neq \nabla P(w_t)$$

Conditioned on $\{w_0, i_0, i_1, \dots, i_{t-1}\}$

However,

 $\mathbb{E}[v_t] = \mathbb{E}[\nabla P(w_t)]$

• Choose $\eta < \frac{1}{L}$ and *m* such that

$$\sigma\coloneqq \frac{1}{\mu\eta(m+1)} + \frac{L\eta}{2-L\eta} < 1$$

• Choose $\eta < \frac{1}{L}$ and *m* such that

$$\sigma \coloneqq \frac{1}{\mu\eta(m+1)} + \frac{L\eta}{2 - L\eta} < 1$$

Then, $\mathbb{E}[||\nabla P(\widetilde{w}^{(s)})||^2] \leq \sigma^s \cdot ||\nabla P(\widetilde{w}^{(0)})||^2$

• Choose $\eta < \frac{1}{L}$ and *m* such that

$$\sigma \coloneqq \frac{1}{\mu\eta(m+1)} + \frac{L\eta}{2 - L\eta} < 1$$

Then, $\mathbb{E}[||\nabla P(\widetilde{w}^{(s)})||^2] \leq \sigma^s \cdot ||\nabla P(\widetilde{w}^{(0)})||^2$

It is a little bit better than SVRG (since SARAH could use the fixed learning rate, whose size is **larger than** that of SVRG and $\sigma < \alpha$ with the same η and m).

• Choose $\eta < \frac{1}{r}$ and *m* such that

$$\sigma \coloneqq \frac{1}{\mu\eta(m+1)} + \frac{L\eta}{2 - L\eta} < 1$$

Then, $\mathbb{E}[||\nabla P(\widetilde{w}^{(s)})||^2] \leq \sigma^s \cdot ||\nabla P(\widetilde{w}^{(0)})||^2$

It is a little bit better than SVRG (since SARAH could use the fixed learning rate, whose size is **larger than** that of SVRG and $\sigma < \alpha$ with the same η and m).

But, they are still considered as the same rate of convergence (linear)

• Choose $\eta < \frac{1}{r}$ and *m* such that

$$\sigma \coloneqq \frac{1}{\mu\eta(m+1)} + \frac{L\eta}{2 - L\eta} < 1$$

Then, $\mathbb{E}[||\nabla P(\widetilde{w}^{(s)})||^2] \leq \sigma^s \cdot ||\nabla P(\widetilde{w}^{(0)})||^2$

It is a little bit better than SVRG (since SARAH could use the fixed learning rate, whose size is **larger than** that of SVRG and $\sigma < \alpha$ with the same η and m).

But, they are still considered as the same rate of convergence (linear)

What is the main difference between SARAH and SVRG?

SARAH One Outer Loop

Recall the update: $w_{t+1} = w_t - \eta v_t$

Recall the update: $w_{t+1} = w_t - \eta v_t$

• *P* is *L*-smooth and μ -strongly convex $\mathbb{E}[||v_t||^2] \le \rho^t \cdot \mathbb{E}[||\nabla P(w_0)||^2]$ $\rho = 1 - \left(\frac{2}{\eta L} - 1\right)\mu^2\eta^2 < 1, \qquad \eta < \frac{2}{L}$ Recall the update: $w_{t+1} = w_t - \eta v_t$

- *P* is *L*-smooth and μ -strongly convex $\mathbb{E}[||v_t||^2] \le \rho^t \cdot \mathbb{E}[||\nabla P(w_0)||^2]$ $\rho = 1 - \left(\frac{2}{\eta L} - 1\right)\mu^2\eta^2 < 1, \qquad \eta < \frac{2}{L}$
- Each f_i , $\forall i$, is *L*-smooth and μ -strongly convex

$$\begin{split} \mathbb{E}[||v_t||^2] &\leq \rho^t \cdot \mathbb{E}[||\nabla P(w_0)||^2] \\ \rho &= 1 - \frac{2\mu L\eta}{\mu + L} < 1, \qquad \qquad \eta \leq \frac{2}{L + \mu} \end{split}$$
Recall the update: $w_{t+1} = w_t - \eta v_t$

E

- *P* is *L*-smooth and μ -strongly convex $\mathbb{E}[||v_t||^2] \le \rho^t \cdot \mathbb{E}[||\nabla P(w_0)||^2]$ $\rho = 1 - \left(\frac{2}{\eta L} - 1\right)\mu^2\eta^2 < 1, \qquad \eta < \frac{2}{L}$
- Each f_i , $\forall i$, is *L*-smooth and μ -strongly convex

$$\begin{split} [||v_t||^2] &\leq \rho^t \cdot \mathbb{E}[||\nabla P(w_0)||^2] \\ \rho &= 1 - \frac{2\mu L\eta}{\mu + L} < 1, \qquad \qquad \eta \leq \frac{2}{L + \mu} \end{split}$$

Hence,

$$\mathbb{E}\big[||\boldsymbol{v}_t||^2\big] \to \mathbf{0} \Rightarrow \mathbb{E}\big[||\boldsymbol{w}_{t+1} - \boldsymbol{w}_t||^2\big] \to \mathbf{0}$$

Recall the update: $w_{t+1} = w_t - \eta v_t$

E

- *P* is *L*-smooth and μ -strongly convex $\mathbb{E}[||v_t||^2] \le \rho^t \cdot \mathbb{E}[||\nabla P(w_0)||^2]$ $\rho = 1 - \left(\frac{2}{\eta L} - 1\right)\mu^2\eta^2 < 1, \qquad \eta < \frac{2}{L}$
- Each f_i , $\forall i$, is *L*-smooth and μ -strongly convex

$$\begin{split} [||v_t||^2] &\leq \rho^t \cdot \mathbb{E}[||\nabla P(w_0)||^2] \\ \rho &= 1 - \frac{2\mu L\eta}{\mu + L} < 1, \qquad \qquad \eta \leq \frac{2}{L + \mu} \end{split}$$

Hence,

$$\mathbb{E}\big[||\boldsymbol{v}_t||^2\big] \to \mathbf{0} \Rightarrow \mathbb{E}\big[||\boldsymbol{w}_{t+1} - \boldsymbol{w}_t||^2\big] \to \mathbf{0}$$

SARAH is converging (somewhere) within a single outer loop with fixed "large" learning rate

Recall the update: $w_{t+1} = w_t - \eta v_t$

- *P* is *L*-smooth and μ -strongly convex $\mathbb{E}[||v_t||^2] \le \rho^t \cdot \mathbb{E}[||\nabla P(w_0)||^2]$ $\rho = 1 - \left(\frac{2}{\eta L} - 1\right)\mu^2\eta^2 < 1, \qquad \eta < \frac{2}{L}$
- Each f_i , $\forall i$, is *L*-smooth and μ -strongly convex

$$\begin{split} \mathbb{E}[||v_t||^2] &\leq \rho^t \cdot \mathbb{E}[||\nabla P(w_0)||^2] \\ \rho &= 1 - \frac{2\mu L\eta}{\mu + L} < 1, \qquad \qquad \eta \leq \frac{2}{L + \mu} \end{split}$$

Hence,

$$\mathbb{E}\big[||\boldsymbol{v}_t||^2\big] \to \mathbf{0} \Rightarrow \mathbb{E}\big[||\boldsymbol{w}_{t+1} - \boldsymbol{w}_t||^2\big] \to \mathbf{0}$$

SARAH is converging (somewhere) within a single outer loop with fixed "large" learning rate

SARAH converges to ϵ -accurate solution within a single outer loop with fixed "small" learning rate for general convex and nonconvex cases (Results and Proofs in the papers)









One Outer Loop Behavior



One Outer Loop Behavior



SARAH is more stable than SVRG!

Sensitivity of SVRG and SARAH on "m"



Sensitivity of SVRG and SARAH on "m"



SARAH has a similar behavior!





SARAH+

• L. Nguyen, J. Liu, K. Scheinberg, and M. Takac. SARAH: A Novel Method for Machine Learning Problems Using Stochastic Recursive Gradient, 2017

Fact #1: Size of update is shrinking It doesn't make sense to do many tiny steps!

Fact #1: Size of update is shrinking It doesn't make sense to do many tiny steps! Heuristic: Restart algorithm when $||v_t||^2 \le \gamma ||v_0||^2$

Fact #1:Size of update is shrinkingIt doesn't make sense to do many tiny steps!

Heuristic: Restart algorithm when $||v_t||^2 \leq \gamma ||v_0||^2$



Fact #1:Size of update is shrinkingIt doesn't make sense to do many tiny steps!

Heuristic: Restart algorithm when $||v_t||^2 \leq \gamma ||v_0||^2$



 $\gamma \approx 1/10$

good performance across many datasets

Datasat	SARAH	SVRG	SAG	SGD+	FISTA
Dataset	(m^*,η^*)	(m^*,η^*)	(η^*)	(η^*)	(η^*)
covtype	(2n, 0.9/L)	(n, 0.8/L)	0.3/L	0.06/L	50/L
ijcnn1	(0.5n, 0.8/L)	(n, 0.5/L)	0.7/L	0.1/L	90/L
news20	(0.5n, 0.9/L)	(n, 0.5/L)	0.1/L	0.2/L	30/L
rcv1	(0.7n, 0.7/L)	(0.5n, 0.9/L)	0.1/L	0.1/L	120/L



Dataset	$egin{array}{c} { m SARAH} \ (m^*,\eta^*) \end{array}$	${{ m SVRG} \over (m^*,\eta^*)}$	$SAG \ (\eta^*)$	SGD+ (η^*)	FISTA (η^*)
covtype	(2n, 0.9/L)	(n, 0.8/L)	0.3/L	0.06/L	50/L
ijcnn1	(0.5n, 0.8/L)	(n, 0.5/L)	0.7/L	0.1/L	90/L
news20	(0.5n, 0.9/L)	(n, 0.5/L)	0.1/L	0.2/L	30/L
rcv1	(0.7n, 0.7/L)	(0.5n, 0.9/L)	0.1/L	0.1/L	120/L



Dataset	SARAH	SVRG	SAG	SGD+	FISTA
Dataset	(m^*,η^*)	(m^*,η^*)	(η^*)	(η^*)	(η^*)
covtype	(2n, 0.9/L)	(n, 0.8/L)	0.3/L	0.06/L	50/L
ijcnn1	(0.5n, 0.8/L)	(n, 0.5/L)	0.7/L	0.1/L	90/L
news20	(0.5n, 0.9/L)	(n, 0.5/L)	0.1/L	0.2/L	30/L
rcv1	(0.7n, 0.7/L)	(0.5n, 0.9/L)	0.1/L	0.1/L	120/L





Strongly convex case: $\kappa = L/\mu$ is a condition number

		Fixed	Low
Method	Complexity	Learning	Storage
		Rate	Cost
GD	$\mathcal{O}\left(n\kappa\log\left(1/\epsilon\right)\right)$	1	1
SGD	$\mathcal{O}\left(1/\epsilon ight)$	×	✓
SVRG	$\mathcal{O}\left((n+\kappa)\log\left(1/\epsilon\right)\right)$	✓	1
SAG/SAGA	$\mathcal{O}\left(\left(n+\kappa\right)\log\left(1/\epsilon\right)\right)$	1	×
SARAH	$\mathcal{O}\left((n+\kappa)\log\left(1/\epsilon\right)\right)$	✓	1

Strongly convex case: $\kappa = L/\mu$ is a condition number

		Fixed	Low
Method	Complexity	Learning	Storage
		Rate	Cost
GD	$\mathcal{O}\left(n\kappa\log\left(1/\epsilon\right)\right)$	1	1
SGD	$\mathcal{O}\left(1/\epsilon ight)$	×	~
SVRG	$\mathcal{O}\left((n+\kappa)\log\left(1/\epsilon\right)\right)$	1	~
SAG/SAGA	$\mathcal{O}\left(\left(n+\kappa\right)\log\left(1/\epsilon\right)\right)$	1	×
SARAH	$\mathcal{O}\left(\left(n+\kappa\right)\log\left(1/\epsilon\right)\right)$	1	1
		•••••	

Practical variant available

	Method	Complexity
	GD	$\mathcal{O}\left(n/\epsilon ight)$
For smooth	SGD	$\mathcal{O}\left(1/\epsilon^2 ight)$
(general) convex	SVRG	$\mathcal{O}\left(n + (\sqrt{n}/\epsilon)\right)$
functions	SAGA	$\mathcal{O}\left(n+(n/\epsilon) ight)$
	SARAH	$\mathcal{O}\left(\left(n+(1/\epsilon)\right)\log(1/\epsilon)\right)$

	Method	Complexity
	GD	$\mathcal{O}\left(\frac{n}{\epsilon}\right)$
For smooth	SGD	$\mathcal{O}\left(\frac{\sigma^2}{\epsilon^2}\right)$
nonconvex functions	SVRG	$\mathcal{O}\left(n + \frac{n^{2/3}}{\nu\epsilon}\right)$
	SARAH	$\mathcal{O}\left(n+\frac{1}{\epsilon^2}\right)$



	Method	Complexity
	GD	$\mathcal{O}\left(n/\epsilon ight)$
For smooth	SGD	$\mathcal{O}\left(1/\epsilon^2\right)$
(general) convex	SVRG	$\mathcal{O}\left(n + (\sqrt{n}/\epsilon)\right)$
functions	SAGA	$\mathcal{O}\left(n+(n/\epsilon)\right)$
	SARAH	$\mathcal{O}\left(\left(n+(1/\epsilon)\right)\log(1/\epsilon)\right)$

	Method	Complexity
	GD	$\mathcal{O}\left(\frac{n}{\epsilon}\right)$
For smooth	SGD	$\mathcal{O}\left(\frac{\sigma^2}{\epsilon^2}\right)$
nonconvex functions	SVRG	$\mathcal{O}\left(n + \frac{n^{2/3}}{\nu\epsilon}\right)$
Tunctions	SARAH	$\mathcal{O}\left(n+\frac{1}{\epsilon^2}\right)$

SARAH converges to ϵ -accurate solution within a single outer loop with fixed "small" learning rate for nonconvex case (Results and Proofs in the papers)

	Method	Complexity
	GD	$\mathcal{O}\left(n/\epsilon ight)$
For smooth	SGD	$\mathcal{O}\left(1/\epsilon^2\right)$
(general) convex	SVRG	$\mathcal{O}\left(n + (\sqrt{n}/\epsilon)\right)$
functions	SAGA	$\mathcal{O}\left(n+(n/\epsilon)\right)$
	SARAH	$\mathcal{O}\left(\left(n+(1/\epsilon)\right)\log(1/\epsilon)\right)$

	Method	Complexity
	GD	$\mathcal{O}\left(\frac{n}{\epsilon}\right)$
For smooth	SGD	$\mathcal{O}\left(\frac{\sigma^2}{\epsilon^2}\right)$
nonconvex functions	SVRG	$\mathcal{O}\left(n + \frac{n^{2/3}}{\nu\epsilon}\right)$
Tunctions	SARAH	$\mathcal{O}\left(n+\frac{1}{\epsilon^2}\right)$

SARAH converges to ϵ -accurate solution within a single outer loop with fixed "small" learning rate for nonconvex case (Results and Proofs in the papers)

Can we get rid of dependence on "n"?

	Method	Complexity
	GD	$\mathcal{O}\left(n/\epsilon ight)$
For smooth	SGD	$\mathcal{O}\left(1/\epsilon^2 ight)$
(general) convex	SVRG	$\mathcal{O}\left(n + (\sqrt{n}/\epsilon)\right)$
functions	SAGA	$\mathcal{O}\left(n+(n/\epsilon)\right)$
	SARAH	$\mathcal{O}\left(\left(n+(1/\epsilon)\right)\log(1/\epsilon)\right)$

	Method	Complexity
	GD	$\mathcal{O}\left(\frac{n}{\epsilon}\right)$
For smooth	SGD	$\mathcal{O}\left(\frac{\sigma^2}{\epsilon^2}\right)$
nonconvex functions	SVRG	$\mathcal{O}\left(n + \frac{n^{2/3}}{\nu\epsilon}\right)$
	SARAH	$\mathcal{O}\left(n+\frac{1}{\epsilon^2}\right)$

SARAH converges to ϵ -accurate solution within a single outer loop with fixed "small" learning rate for nonconvex case (Results and Proofs in the papers)

Can we get rid of dependence on "n"?



iSARAH

• L. Nguyen, K. Scheinberg, and M. Takac. **Inexact SARAH for Large Scale Machine Learning Problems**. *In preparation*.

```
Parameters: the learning rate \eta > 0 and the inner loop size m.
Initialize: \tilde{w}_0
Iterate:
for s = 1, 2, ... do
  w_0 = \tilde{w}_{s-1}
  Choose a subset I \subseteq [n] of size b uniformly at random (without replacement)
  v_0 = \frac{1}{b} \sum_{i \in I} \nabla f_i(w_0)
  w_1 = w_0 - \eta v_0
  Iterate:
  for t = 1, ..., m - 1 do
      Sample i_t uniformly at random from [n]
      v_t = \nabla f_{i_t}(w_t) - \nabla f_{i_t}(w_{t-1}) + v_{t-1}
      w_{t+1} = w_t - \eta v_t
  end for
   Set \tilde{w} = w_t with t chosen uniformly at random from \{0, 1, \dots, m\}
  Output: \tilde{w}_s = \tilde{w}
end for
```

Parameters: the learning rate $\eta > 0$ and the inner loop size m. Initialize: \tilde{w}_0 Iterate: for s = 1, 2, ... do **NOT computing Full gradient** $w_0 = \tilde{w}_{s-1}$ Choose a subset $I \subseteq [n]$ is size b uniformly at random (without replacement) $v_0 = \frac{1}{b} \sum_{i \in I} \nabla f_i(w_0)$ $w_1 = w_0 - \eta v_0$ Iterate: for t = 1, ..., m - 1 do Sample i_t uniformly at random from [n] $v_t = \nabla f_{i_t}(w_t) - \nabla f_{i_t}(w_{t-1}) + v_{t-1}$ $w_{t+1} = w_t - \eta v_t$ end for Set $\tilde{w} = w_t$ with t chosen uniformly at random from $\{0, 1, \dots, m\}$ **Output:** $\tilde{w}_s = \tilde{w}$ end for

Parameters: the learning rate $\eta > 0$ and the inner loop size m. Initialize: \tilde{w}_0 Iterate: for s = 1, 2, ... do **NOT** computing Full gradient $w_0 = \tilde{w}_{s-1}$ Choose a subset $I \subseteq [n]$ is size b uniformly at random (without replacement) $v_0 = \frac{1}{b} \sum_{i \in I} \nabla f_i(w_0)$ $w_1 = w_0 - \eta v_0$ Iterate: for t = 1, ..., m - 1 do Sample i_t uniformly at random from [n] $v_t = \nabla f_{i_t}(w_t) - \nabla f_{i_t}(w_{t-1}) + v_{t-1}$ $w_{t+1} = w_t - \eta v_t$ end for Set $\tilde{w} = w_t$ with t chosen uniformly at random from $\{0, 1, \dots, m\}$ **Output:** $\tilde{w}_s = \tilde{w}$ end for

- For smooth general convex functions: b = m
- For smooth nonconvex functions: $b = \sqrt{m}$



	Method	Complexity
	GD	$\mathcal{O}\left(n/\epsilon ight)$
For smooth	SGD	$\mathcal{O}\left(1/\epsilon^2 ight)$
	SVRG	$\mathcal{O}\left(n + (\sqrt{n}/\epsilon)\right)$
(general) convex	SAGA	$\mathcal{O}\left(n+(n/\epsilon) ight)$
functions	SARAH	$\mathcal{O}\left(\left(n+(1/\epsilon)\right)\log(1/\epsilon) ight)$
	iSARAH	$\mathcal{O}\left((1/\epsilon)\log(1/\epsilon) ight)$

	Method	Complexity
	GD	$\mathcal{O}\left(\frac{n}{\epsilon}\right)$
For smooth	SGD	$\mathcal{O}\left(\frac{\sigma^2}{\epsilon^2}\right)$
nonconvex functions	SVRG	$\mathcal{O}\left(n + \frac{n^{2/3}}{\nu\epsilon}\right)$
	SARAH	$\mathcal{O}\left(n+\frac{1}{\epsilon^2}\right)$
	iSARAH	$\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

	Method	Complexity
	GD	$\mathcal{O}\left(n/\epsilon ight)$
For smooth	SGD	$\mathcal{O}\left(1/\epsilon^2\right)$
	SVRG	$\mathcal{O}\left(n + (\sqrt{n}/\epsilon)\right)$
(general) convex	SAGA	$\mathcal{O}\left(n+(n/\epsilon) ight)$
functions	SARAH	$\mathcal{O}\left(\left(n+(1/\epsilon)\right)\log(1/\epsilon)\right)$
	iSARAH	$\mathcal{O}\left((1/\epsilon)\log(1/\epsilon) ight)$

	Method	Complexity
	GD	$\mathcal{O}\left(\frac{n}{\epsilon}\right)$
For smooth	SGD	$\mathcal{O}\left(\frac{\sigma^2}{\epsilon^2}\right)$
nonconvex functions	SVRG	$\mathcal{O}\left(n + \frac{n^{2/3}}{\nu\epsilon}\right)$
	SARAH	$\mathcal{O}\left(n+\frac{1}{\epsilon^2}\right)$
	iSARAH	$\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

Total complexity of iSARAH does not depend on "n"
Convergence Rates Comparisons

	Method	Complexity
	GD	$\mathcal{O}\left(n/\epsilon ight)$
For smooth (general) convex functions	SGD	$\mathcal{O}\left(1/\epsilon^2 ight)$
	SVRG	$\mathcal{O}\left(n + (\sqrt{n}/\epsilon)\right)$
	SAGA	$\mathcal{O}\left(n+(n/\epsilon) ight)$
	SARAH	$\mathcal{O}\left(\left(n+(1/\epsilon)\right)\log(1/\epsilon)\right)$
	iSARAH	$\mathcal{O}\left((1/\epsilon)\log(1/\epsilon) ight)$

	Method	Complexity
	GD	$\mathcal{O}\left(\frac{n}{\epsilon}\right)$
For smooth	SGD	$\mathcal{O}\left(\frac{\sigma^2}{\epsilon^2}\right)$
nonconvex functions	SVRG	$\mathcal{O}\left(n + \frac{n^{2/3}}{\nu\epsilon}\right)$
	SARAH	$\mathcal{O}\left(n+\frac{1}{\epsilon^2}\right)$
	iSARAH	$\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

Total complexity of iSARAH does not depend on "*n*" VERY USEFUL for large scale machine learning problems!!!



SARAH

THANK YOU !!!

LamNguyen.MLTD@gmail.com

Lam M. Nguyen – Lehigh University http://coral.ise.lehigh.edu/lmn214/