

LARGE-SCALE CONVEX OPTIMIZATION: PARALLELIZATION AND VARIANCE REDUCTION

Ph.D. Thesis Defense

June 10th, 2024

Cheik Traoré¹

Supervisor: **Silvia Villa**¹

Jury: **Claudio Estatico**¹, **François Glineur**² and **Juan Peypouquet**³



TRAINING DATA-DRIVEN EXPERTS IN
OPTIMIZATION
MSCA-ITN 2019

Context

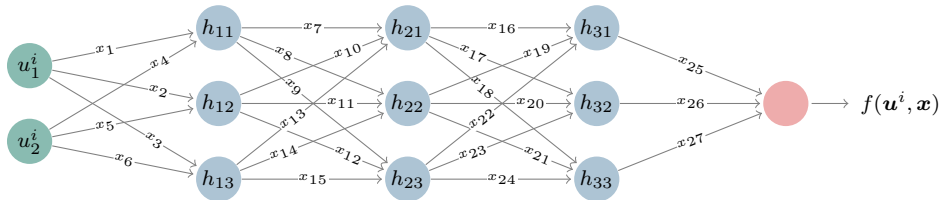
Training dataset: $(\mathbf{u}^i, y^i)_{i \in \{1, 2, \dots, n\}}$.

$$\underset{\mathbf{x} \in \mathbb{R}^m}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{u}^i; \mathbf{x}), y_i) + \lambda R(\mathbf{x}).$$

Context

Training dataset: $(\mathbf{u}^i, y^i)_{i \in \{1, 2, \dots, n\}}$.

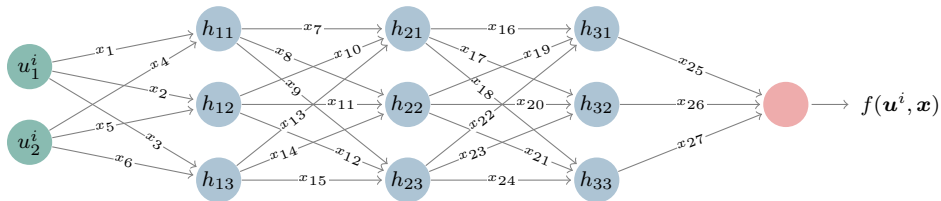
$$\underset{\mathbf{x} \in \mathbb{R}^m}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{u}^i; \mathbf{x}), y_i) + \lambda R(\mathbf{x}).$$



Context

Training dataset: $(\mathbf{u}^i, y^i)_{i \in \{1, 2, \dots, n\}}$.

$$\underset{\mathbf{x} \in \mathbb{R}^m}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{u}^i; \mathbf{x}), y_i) + \lambda R(\mathbf{x}).$$

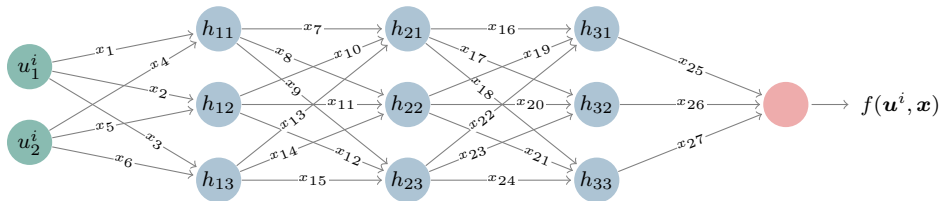


n and m can be very big: several BILLIONS!!!

Context

Training dataset: $(u^i, y^i)_{i \in \{1, 2, \dots, n\}}$.

$$\underset{\mathbf{x} \in \mathbb{R}^m}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \ell(f(u^i; \mathbf{x}), y_i) + \lambda R(\mathbf{x}).$$



n and m can be very big: several BILLIONS!!!

ChatGP4: **estimated** 1.76 trillion parameters (Georges Hotz).

Outline

General introduction

Asynchronous Forward-Backward

Variance reduction techniques for SPPA

Conclusion

Optimization for data science

Goal:

$$\underset{x \in \mathbf{H}}{\text{minimize}} \ F(x) = f(x) + g(x),$$

where \mathbf{H} is a separable real Hilbert space and $F, f, g: \mathbf{H} \rightarrow \mathbb{R}$.

Optimization for data science

Goal:

$$\underset{\boldsymbol{x} \in \mathbf{H}}{\text{minimize}} \ F(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x}),$$

where \mathbf{H} is a separable real Hilbert space and $F, f, g: \mathbf{H} \rightarrow \mathbb{R}$.

Example:

$$\underset{\boldsymbol{x} \in \mathbb{R}^m}{\text{minimize}} \ F(\boldsymbol{x}) = \underbrace{\frac{1}{n} \sum_{i=1}^n f_i(\boldsymbol{x})}_{f(\boldsymbol{x})} + \underbrace{\lambda R(\boldsymbol{x})}_{g(\boldsymbol{x})},$$

Most popular method

Gradient descent (GD):

$$\boldsymbol{x}^0 \in \mathbf{H}$$

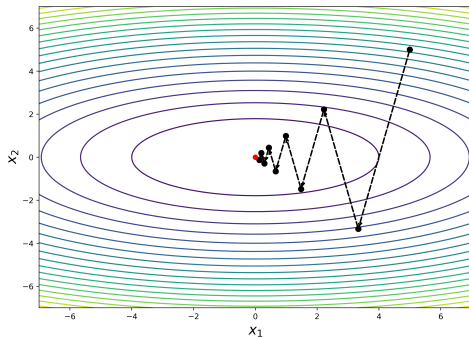
$$\boldsymbol{x}^{k+1} = \boldsymbol{x}^k - \gamma_k \nabla F(\boldsymbol{x}^k).$$

Most popular method

Gradient descent (GD):

$$\mathbf{x}^0 \in \mathbf{H}$$

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \nabla F(\mathbf{x}^k).$$



GD rates

F is μ -strongly convex for $\mu \geq 0$

GD rates

F is μ -strongly convex for $\mu \geq 0$:

$$(\forall \mathbf{x}, \mathbf{y} \in \mathbf{H}) \quad (t \in [0, 1]) \quad F(t\mathbf{x} + (1 - t)\mathbf{y}) \leq tF(\mathbf{x}) + (1 - t)F(\mathbf{y}) - \frac{\mu}{2}t(1 - t)\|\mathbf{x} - \mathbf{y}\|^2.$$

GD rates

F is μ –strongly convex for $\mu \geq 0$:

$$(\forall \mathbf{x}, \mathbf{y} \in \mathbf{H}) \quad (t \in [0, 1]) \quad F(t\mathbf{x} + (1 - t)\mathbf{y}) \leq tF(\mathbf{x}) + (1 - t)F(\mathbf{y}) - \frac{\mu}{2}t(1 - t)\|\mathbf{x} - \mathbf{y}\|^2.$$

F is L –smooth for $L \geq 0$:

$$(\forall \mathbf{x}, \mathbf{y} \in \mathbf{H}) \quad \|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$$

GD rates

Let $\gamma_k = \frac{1}{L}$.

GD rates

Let $\gamma_k = \frac{1}{L}$.

► If $\mu = 0$,

$$F(\mathbf{x}_k) - F_* \leq O(1/k) \quad (\text{sublinear rate}).$$

GD rates

Let $\gamma_k = \frac{1}{L}$.

► If $\mu = 0$,

$$F(\mathbf{x}_k) - F_* \leq O(1/k) \quad (\text{sublinear rate}).$$

► If $\mu > 0$,

$$F(\mathbf{x}_k) - F_* \leq O(\epsilon^k) \quad (\text{linear rate}),$$

with $\epsilon < 1$.

Sum of functions problem

First case: $g \equiv 0$ and $f(\boldsymbol{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\boldsymbol{x})$.

Sum of functions problem

First case: $g \equiv 0$ and $f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$.

The problem is now:

$$\underset{\mathbf{x} \in \mathbf{H}}{\text{minimize}} \ F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}).$$

Stochastic gradient descent (SGD)

SGD update is:

Select uniformly at random $i_k \in [n] := \{1, 2, \dots, n\}$ and do

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \nabla f_{i_k}(\mathbf{x}^k).$$

Non-smooth problem

Second case: g is not differentiable.

Non-smooth problem

Second case: g is not differentiable.

The problem is:

$$\underset{x \in \mathbf{H}}{\text{minimize}} \ F(x) = f(x) + g(x),$$

with f smooth.

Non-smooth problem

Second case: g is not differentiable.

The problem is:

$$\underset{\boldsymbol{x} \in \mathbf{H}}{\text{minimize}} \ F(\boldsymbol{x}) = f(\boldsymbol{x}) + g(\boldsymbol{x}),$$

with f smooth.

So F is non-smooth!!

Subgradient descent

Update:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \partial F(\mathbf{x}^k).$$

Subgradient descent

Update:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \partial F(\mathbf{x}^k).$$

Rates with $\gamma_k \rightarrow 0$ and bounded subgradient:

- Convex case: $F(\bar{\mathbf{x}}^k) - F_* \leq O\left(\frac{1}{\sqrt{k}}\right)$.
- Strongly convex case: $F(\bar{\mathbf{x}}^k) - F_* \leq O(1/k)$.

$\bar{\mathbf{x}}^k = \sum_{t=i}^k p_t \mathbf{x}^t$ with $\sum_{t=i}^k p_t = 1$ (ergodic rates).

Proximal Point Algorithm (PPA)

Update:

$$\begin{aligned}\mathbf{x}^{k+1} &= \text{prox}_{\gamma_k F}(\mathbf{x}^k) = \text{prox}_{\gamma_k(f+g)}(\mathbf{x}^k) \\ &= \underset{\mathbf{x}}{\text{argmin}} F(\mathbf{x}) + \frac{1}{2\gamma_k} \|\mathbf{x} - \mathbf{x}^k\|^2 \\ &= \underset{\mathbf{x}}{\text{argmin}} f(\mathbf{x}) + g(\mathbf{x}) + \frac{1}{2\gamma_k} \|\mathbf{x} - \mathbf{x}^k\|^2.\end{aligned}$$

Proximal Point Algorithm (PPA)

Update:

$$\begin{aligned}\mathbf{x}^{k+1} &= \text{prox}_{\gamma_k F}(\mathbf{x}^k) = \text{prox}_{\gamma_k(f+g)}(\mathbf{x}^k) \\ &= \underset{\mathbf{x}}{\operatorname{argmin}} F(\mathbf{x}) + \frac{1}{2\gamma_k} \|\mathbf{x} - \mathbf{x}^k\|^2 \\ &= \underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x}) + g(\mathbf{x}) + \frac{1}{2\gamma_k} \|\mathbf{x} - \mathbf{x}^k\|^2.\end{aligned}$$

Rates with $\gamma_k = \gamma \in \mathbb{R}_+$:

- Convex case: $F(\mathbf{x}_k) - F_* \leq O(1/k)$.
- Strongly convex case: $F(\mathbf{x}_k) - F_* \leq O(\varepsilon^k)$, with $\epsilon < 1$.

Forward-Backward

Proximal gradient (Forward-Backward) update:

$$\boldsymbol{x}^{k+1} = \text{prox}_{\gamma_k g} (\boldsymbol{x}^k - \gamma_k \nabla f(\boldsymbol{x}^k)) .$$

Forward-Backward

Proximal gradient (Forward-Backward) update:

$$\mathbf{x}^{k+1} = \text{prox}_{\gamma_k g} (\mathbf{x}^k - \gamma_k \nabla f(\mathbf{x}^k)) .$$

Rates with $\gamma_k = \frac{1}{L}$:

- Convex case: $F(\mathbf{x}_k) - F_* \leq O(1/k)$.
- Strongly convex case: $F(\mathbf{x}_k) - F_* \leq O(\epsilon^k)$, with $\epsilon < 1$.

Parallel Forward-Backward

Suppose that $\mathbf{H} = \bigoplus_{i=1}^m \mathbf{H}_i$, that we have a central server S and m machines: M_1, M_2, \dots, M_m .

Parallel Forward-Backward

Suppose that $\mathbf{H} = \bigoplus_{i=1}^m \mathbf{H}_i$, that we have a central server S and m machines: M_1, M_2, \dots, M_m .

- For $i \in \{1, 2, \dots, m\}$, machine M_i computes $\nabla_i f(\mathbf{x}^k)$.

Parallel Forward-Backward

Suppose that $\mathbf{H} = \bigoplus_{i=1}^m \mathbf{H}_i$, that we have a central server S and m machines: M_1, M_2, \dots, M_m .

- ▶ For $i \in \{1, 2, \dots, m\}$, machine M_i computes $\nabla_i f(\mathbf{x}^k)$.
- ▶ The server S reconstitutes $\nabla f(\mathbf{x}^k)$,

Parallel Forward-Backward

Suppose that $\mathbf{H} = \bigoplus_{i=1}^m \mathbf{H}_i$, that we have a central server S and m machines: M_1, M_2, \dots, M_m .

- ▶ For $i \in \{1, 2, \dots, m\}$, machine M_i computes $\nabla_i f(\mathbf{x}^k)$.
- ▶ The server S reconstitutes $\nabla f(\mathbf{x}^k)$,
- ▶ and makes the update:

$$\mathbf{x}^{k+1} = \text{prox}_{\gamma_k g} (\mathbf{x}^k - \gamma_k \nabla f(\mathbf{x}^k)) .$$

Synchronous Parallel Forward-Backward

Suppose that $\mathbf{H} = \bigoplus_{i=1}^m \mathbf{H}_i$, that we have a central server S and m machines: M_1, M_2, \dots, M_m .

- ▶ For $i \in \{1, 2, \dots, m\}$, machine M_i computes $\nabla_i f(\mathbf{x}^k)$.
- ▶ The server S reconstitutes $\nabla f(\mathbf{x}^k)$,
- ▶ and makes the update:

$$\mathbf{x}^{k+1} = \text{prox}_{\gamma_k g} (\mathbf{x}^k - \gamma_k \nabla f(\mathbf{x}^k)) .$$

Synchronous Parallel Forward-Backward

Suppose that $\mathbf{H} = \bigoplus_{i=1}^m \mathbf{H}_i$, that we have a central server S and m machines: M_1, M_2, \dots, M_m .

- ▶ For $i \in \{1, 2, \dots, m\}$, machine M_i computes $\nabla_i f(\mathbf{x}^k)$.
- ▶ The server S reconstitutes $\nabla f(\mathbf{x}^k)$,
- ▶ and makes the update:

$$\mathbf{x}^{k+1} = \text{prox}_{\gamma_k g} (\mathbf{x}^k - \gamma_k \nabla f(\mathbf{x}^k)) .$$

Synchronous algorithms are as slow as the slowest machine.

Outline

General introduction

Asynchronous Forward-Backward

Variance reduction techniques for SPPA

Conclusion

Problem setting

Consider:

$$\underset{\mathbf{x} \in \mathbf{H}}{\text{minimize}} \ F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x}), \quad g(\mathbf{x}) = \sum_{i=1}^m g_i(\mathbf{x}_i),$$

where $\mathbf{H} = \bigoplus_{i=1}^m \mathbf{H}_i$.

Coordinate Forward-Backward

Algorithm

Let $(\gamma_i)_{1 \leq i \leq m} \in \mathbb{R}_{++}^m$ and $\mathbf{x}^0 = (x_1^0, \dots, x_m^0) \in \mathbf{H}$. Iterate

for $k = 0, 1, \dots$

 choose i_k uniformly at random in $[m] := \{1, 2, \dots, m\}$
 for $i = 1, \dots, m$

$$x_i^{k+1} = \begin{cases} \text{prox}_{\gamma_i g_i}(x_i^k - \gamma_i \nabla_i f(\mathbf{x}^k)) & \text{if } i = i_k \\ x_i^k & \text{if } i \neq i_k. \end{cases}$$

Asynchronous Forward-Backward

Server S and m machines: M_1, M_2, \dots, M_m .

Asynchronous Forward-Backward

Server S and m machines: M_1, M_2, \dots, M_m .

- For $i \in \{1, 2, \dots, m\}$, machine M_i computes $\nabla_i f(\mathbf{x}^k)$.

Asynchronous Forward-Backward

Server S and m machines: M_1, M_2, \dots, M_m .

- ▶ For $i \in \{1, 2, \dots, m\}$, machine M_i computes $\nabla_i f(\mathbf{x}^k)$.
- ▶ The server S receives $\nabla_{i_k} f(\mathbf{x}^k)$ from machine M_{i_k} ,

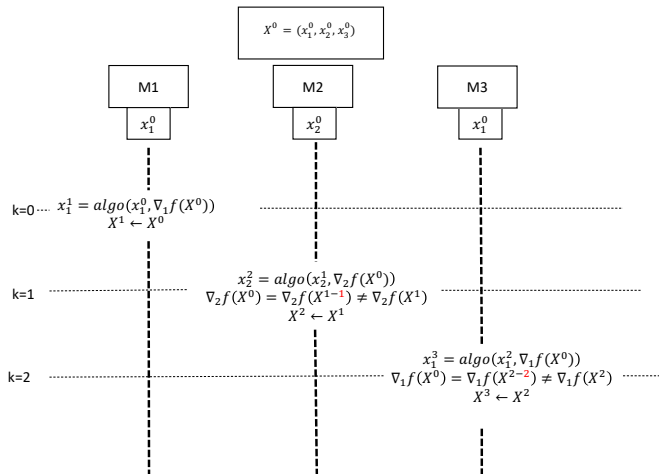
Asynchronous Forward-Backward

Server S and m machines: M_1, M_2, \dots, M_m .

- ▶ For $i \in \{1, 2, \dots, m\}$, machine M_i computes $\nabla_i f(\mathbf{x}^k)$.
- ▶ The server S receives $\nabla_{i_k} f(\mathbf{x}^k)$ from machine M_{i_k} ,
- ▶ and makes the update:

$$\mathbf{x}_{i_k}^{k+1} = \text{prox}_{\gamma_{i_k} g_{i_k}} \left(\mathbf{x}_{i_k}^k - \gamma_{i_k} \nabla_{i_k} f(\mathbf{x}^k) \right).$$

Delayed gradient in asynchronous setting



Asynchronous Forward-Backward

Algorithm

Let $(\gamma_i)_{1 \leq i \leq m} \subset \mathbb{R}_{++}^m$ and $\mathbf{x}^0 = (x_1^0, \dots, x_m^0) \in \mathbf{H}$. Iterate

for $k = 0, 1, \dots$

 choose i_k uniformly at random in $[m]$
 for $i = 1, \dots, m$

$$x_i^{k+1} = \begin{cases} \text{prox}_{\gamma_i g_i}(x_i^k - \gamma_i \nabla_i f(\mathbf{x}^{k-\mathbf{d}^k})) & \text{if } i = i_k \\ x_i^k & \text{if } i \neq i_k, \end{cases}$$

where $\mathbf{x}^{k-\mathbf{d}^k} = (x_1^{k-\mathbf{d}^k}, \dots, x_m^{k-\mathbf{d}^k})$.

Asynchronous Forward-Backward

Algorithm

Let $(\gamma_i)_{1 \leq i \leq m} \subset \mathbb{R}_{++}^m$ and $\mathbf{x}^0 = (x_1^0, \dots, x_m^0) \in \mathbf{H}$. Iterate

for $k = 0, 1, \dots$

choose i_k uniformly at random in $[m]$

for $i = 1, \dots, m$

$$x_i^{k+1} = \begin{cases} \text{prox}_{\gamma_i g_i}(x_i^k - \gamma_i \nabla_i f(\mathbf{x}^{k-\mathbf{d}^k})) & \text{if } i = i_k \\ x_i^k & \text{if } i \neq i_k, \end{cases}$$

where $\mathbf{x}^{k-\mathbf{d}^k} = (x_1^{k-\mathbf{d}^k}, \dots, x_m^{k-\mathbf{d}^k})$. $\mathbf{d}^k \in \mathbb{N}$.

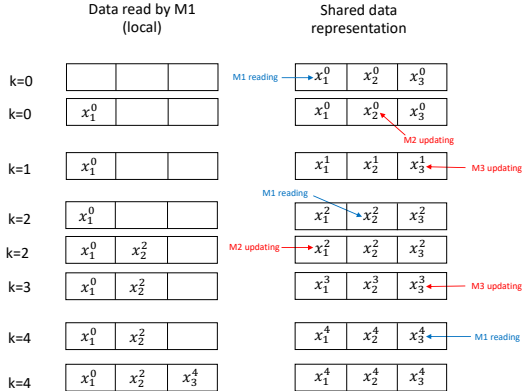
Read paradigm

- Consistent read

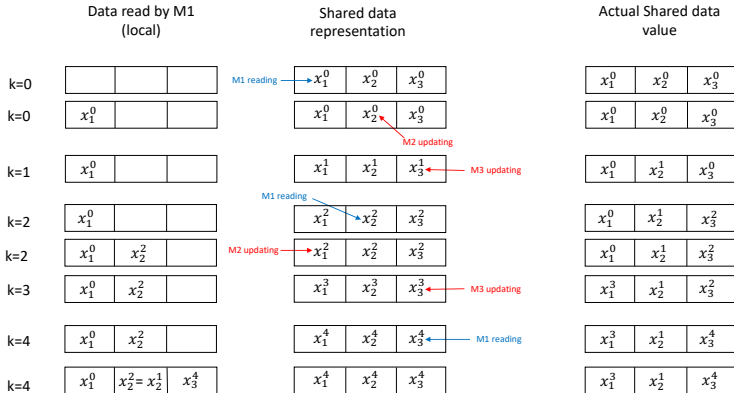
Read paradigm

- ▶ Consistent read
- ▶ Inconsistent read

Inconsistent read



Inconsistent read



Algorithm studied²

Algorithm

Let $(\gamma_i)_{1 \leq i \leq m} \subset \mathbb{R}_{++}^m$ and $\mathbf{x}^0 = (x_1^0, \dots, x_m^0) \in \mathbf{H}$. Iterate

for $k = 0, 1, \dots$

 choose randomly i_k in $[m]$ with probability p_{i_k}
 for $i = 1, \dots, m$

$$x_i^{k+1} = \begin{cases} \text{prox}_{\gamma_i g_i}(x_i^k - \gamma_i \nabla_i f(\mathbf{x}^{k-\mathbf{d}^k})) & \text{if } i = i_k \\ x_i^k & \text{if } i \neq i_k, \end{cases}$$

where $\mathbf{x}^{k-\mathbf{d}^k} = (x_1^{k-\mathbf{d}_1^k}, \dots, x_m^{k-\mathbf{d}_m^k})$.

²Traoré, Salzo, et al., "Convergence of an asynchronous block-coordinate forward-backward algorithm for convex composite optimization".

Algorithm studied²

Algorithm

Let $(\gamma_i)_{1 \leq i \leq m} \subset \mathbb{R}_{++}^m$ and $\mathbf{x}^0 = (x_1^0, \dots, x_m^0) \in \mathbf{H}$. Iterate

for $k = 0, 1, \dots$

 choose randomly i_k in $[m]$ with probability p_{i_k}
 for $i = 1, \dots, m$

$$x_i^{k+1} = \begin{cases} \text{prox}_{\gamma_i g_i}(x_i^k - \gamma_i \nabla_i f(\mathbf{x}^{k-\mathbf{d}^k})) & \text{if } i = i_k \\ x_i^k & \text{if } i \neq i_k, \end{cases}$$

where $\mathbf{x}^{k-\mathbf{d}^k} = (x_1^{k-\mathbf{d}_1^k}, \dots, x_m^{k-\mathbf{d}_m^k})$. $\mathbf{d}^k \in \mathbb{N}^m$.

²Traoré, Salzo, et al., "Convergence of an asynchronous block-coordinate forward-backward algorithm for convex composite optimization".

Algorithm studied²

Algorithm

Let $(\gamma_i)_{1 \leq i \leq m} \subset \mathbb{R}_{++}^m$ and $\mathbf{x}^0 = (x_1^0, \dots, x_m^0) \in \mathbf{H}$. Iterate

for $k = 0, 1, \dots$

 choose randomly i_k in $[m]$ with probability p_{i_k}
 for $i = 1, \dots, m$

$$x_i^{k+1} = \begin{cases} \text{prox}_{\gamma_i g_i}(x_i^k - \gamma_i \nabla_i f(\mathbf{x}^{k-\mathbf{d}^k})) & \text{if } i = i_k \\ x_i^k & \text{if } i \neq i_k, \end{cases}$$

where $\mathbf{x}^{k-\mathbf{d}^k} = (x_1^{k-\mathbf{d}_1^k}, \dots, x_m^{k-\mathbf{d}_m^k})$. $\mathbf{d}^k \in \mathbb{N}^m$. $p_{\max} = \max_i p_i$ and $p_{\min} = \min_i p_i > 0$.

²Traoré, Salzo, et al., "Convergence of an asynchronous block-coordinate forward-backward algorithm for convex composite optimization".

Assumptions

- ▶ $f: \mathbf{H} \rightarrow \mathbb{R}$ is convex and differentiable.
- ▶ For every $i \in \{1, \dots, m\}$, $g_i: H_i \rightarrow]-\infty, +\infty]$ is proper convex and lower semicontinuous.
- ▶ For all $\mathbf{x} \in \mathbf{H}$ and $i \in \{1, \dots, m\}$, the map $\nabla_i f(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \cdot, \mathbf{x}_{i+1}, \dots, \mathbf{x}_m): H_i \rightarrow H_i$ is Lipschitz continuous with constant L_i . Define $L_{\max} := \max_i L_i$ and $L_{\min} := \min_i L_i$.
- ▶ For all $\mathbf{x} \in \mathbf{H}$ and $i \in \{1, \dots, m\}$, the map $\nabla f(\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \cdot, \mathbf{x}_{i+1}, \dots, \mathbf{x}_m): H_i \rightarrow \mathbf{H}$ is Lipschitz continuous with constant $L_{\text{res}} > 0$. Note that $L_{\max} \leq L_{\text{res}}$.
- ▶ F attains its minimum $F_* := \min F$ on \mathbf{H} .

Assumptions

- ▶ We assume that the delays are deterministic and bounded by τ .
- ▶ The delay vector is independent of the coordinates.

Main result

Theorem (Convex case)

Assume that $\gamma_i < \frac{2}{L_i + 2\tau \frac{\rho_{\max}}{\sqrt{\rho_{\min}}}}$ for all $i \in [m]$.

Main result

Theorem (Convex case)

Assume that $\gamma_i < \frac{2}{L_i + 2\tau \frac{\rho_{\max}}{\sqrt{\rho_{\min}}}}$ for all $i \in [m]$. Then,

- $(x^k)_{k \in \mathbb{N}}$ converges weakly P-a.s. to x^* in $\operatorname{argmin} F$.

Main result

Theorem (Convex case)

Assume that $\gamma_i < \frac{2}{L_i + 2\tau \frac{p_{\max}}{\sqrt{p_{\min}}}}$ for all $i \in [m]$. Then,

- ▶ $(x^k)_{k \in \mathbb{N}}$ converges weakly P-a.s. to x^* in $\operatorname{argmin} F$.
- ▶ And

$$(\forall k \in \mathbb{N}) \quad \mathbb{E}[F(x^k) - F_*] \leq \frac{1}{k} \left(\frac{\operatorname{dist}_W^2(x^0, \operatorname{argmin} F)}{2} + C (F(x^0) - F_*) \right),$$

where $C = O(\tau)$ and $W = \bigoplus_{i=1}^m \frac{1}{\gamma_i p_i} \operatorname{Id}_i$.

Linear convergence

Luo-Tseng error bound condition:

$$(\forall \mathbf{x} \in X) \quad \text{dist}_{\Gamma^{-1}}(\mathbf{x}, \text{argmin } F) \leq C_{X, \Gamma^{-1}} \left\| \mathbf{x} - \text{prox}_g^{\Gamma^{-1}}(\mathbf{x} - \nabla^{\Gamma^{-1}} f(\mathbf{x})) \right\|_{\Gamma^{-1}},$$

where $\Gamma^{-1} = \bigoplus_{i=1}^m \frac{1}{\gamma_i} \text{Id}_i$.

Linear convergence

Luo-Tseng error bound condition:

$$(\forall \mathbf{x} \in \mathbf{X}) \quad \text{dist}_{\Gamma^{-1}}(\mathbf{x}, \argmin F) \leq C_{\mathbf{X}, \Gamma^{-1}} \left\| \mathbf{x} - \text{prox}_g^{\Gamma^{-1}}(\mathbf{x} - \nabla^{\Gamma^{-1}} f(\mathbf{x})) \right\|_{\Gamma^{-1}},$$

where $\Gamma^{-1} = \bigoplus_{i=1}^m \frac{1}{\gamma_i} \text{Id}_i$.

Equivalent to *quadratic growth*:

$$(\forall \mathbf{x} \in \mathbf{X}) \quad \frac{\text{func}(C_{\mathbf{X}, \Gamma^{-1}})}{2} \text{dist}_{\Gamma^{-1}}^2(\mathbf{x}, \argmin F) \leq F(\mathbf{x}_k) - F_*.$$

Linear convergence

Theorem (With error bound condition)

Suppose that $\gamma_i < \frac{2}{L_i + 2\tau \frac{\rho_{\max}}{\sqrt{\rho_{\min}}}}$ for all $i \in [m]$.

Linear convergence

Theorem (With error bound condition)

Suppose that $\gamma_i < \frac{2}{L_i + 2\tau \frac{\mathbf{p}_{\max}}{\sqrt{\mathbf{p}_{\min}}}}$ for all $i \in [m]$. Then,

$$\blacktriangleright (\forall k \in \mathbb{N}) \quad \mathbb{E}[F(\mathbf{x}^{k+1}) - F_*] \leq \left(1 - \frac{\mathbf{p}_{\min}}{\kappa + \theta}\right)^{\lfloor \frac{k+1}{\tau+1} \rfloor} \mathbb{E}[F(\mathbf{x}^0) - F_*],$$

where $\kappa \geq 1$ and $\theta > 0$.

Linear convergence

Theorem (With error bound condition)

Suppose that $\gamma_i < \frac{2}{L_i + 2\tau \frac{\mathbf{p}_{\max}}{\sqrt{\mathbf{p}_{\min}}}}$ for all $i \in [m]$. Then,

$$\blacktriangleright (\forall k \in \mathbb{N}) \quad \mathbb{E}[F(\mathbf{x}^{k+1}) - F_*] \leq \left(1 - \frac{\mathbf{p}_{\min}}{\kappa + \theta}\right)^{\lfloor \frac{k+1}{\tau+1} \rfloor} \mathbb{E}[F(\mathbf{x}^0) - F_*],$$

where $\kappa \geq 1$ and $\theta > 0$.

$\blacktriangleright (\mathbf{x}^k)_{k \in \mathbb{N}}$ converges strongly P-a.s. to $\mathbf{x}^* \in \operatorname{argmin} F$ and

$$(\forall k \in \mathbb{N}) \quad \mathbb{E}[\|\mathbf{x}^k - \mathbf{x}^*\|_{\Gamma^{-1}}] = \mathcal{O}\left((1 - \mathbf{p}_{\min}/(\kappa + \theta))^{\lfloor \frac{k}{\tau+1} \rfloor / 2}\right).$$

Related works

► Liu and Wright³

- constant stepsize
- uniform probability
- geometric (exponential) dependence of the stepsize on the delay.

³Liu et al., “Asynchronous stochastic coordinate descent: Parallelism and convergence properties”.

Related works

- ▶ Liu and Wright³
 - constant stepsize
 - uniform probability
 - geometric (exponential) dependence of the stepsize on the delay.
- ▶ Salzo and Villa⁴
 - No delay consideration.

³Liu et al., “Asynchronous stochastic coordinate descent: Parallelism and convergence properties”.

⁴Salzo et al., “Parallel random block-coordinate forward–backward algorithm: a unified convergence analysis”.

Related works

- ▶ Liu and Wright³
 - constant stepsize
 - uniform probability
 - geometric (exponential) dependence of the stepsize on the delay.
- ▶ Salzo and Villa⁴
 - No delay consideration.
- ▶ Other works: Davis⁵, Cannelli et al.⁶

³Liu et al., “Asynchronous stochastic coordinate descent: Parallelism and convergence properties”.

⁴Salzo et al., “Parallel random block-coordinate forward–backward algorithm: a unified convergence analysis”.

⁵Davis, “The asynchronous palm algorithm for nonsmooth nonconvex problems”.

⁶Cannelli et al., “Asynchronous parallel algorithms for nonconvex optimization”.

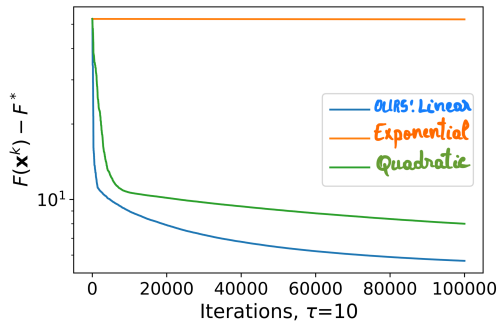
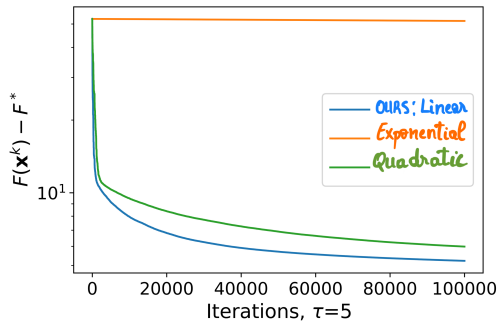
Experiments

$A \in \mathbb{R}^{n \times m}$ and $b \in \mathbb{R}^n$,

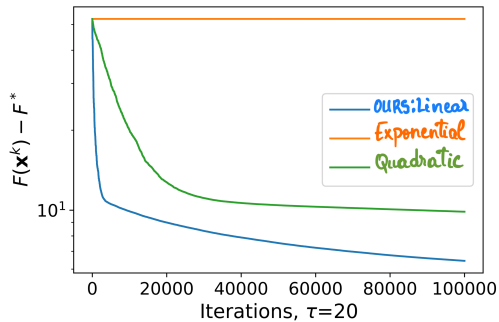
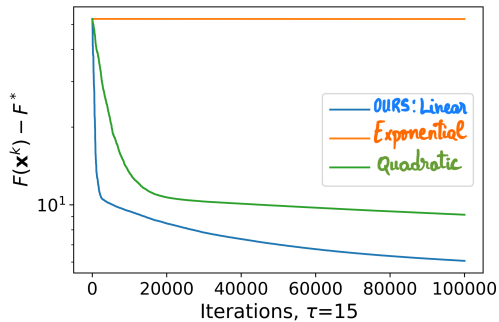
$$\underset{\mathbf{x} \in \mathbb{R}^m}{\text{minimize}} \quad \frac{1}{2} \|A\mathbf{x} - b\|_2^2 + \lambda \|\mathbf{x}\|_1 \quad (\lambda > 0).$$

Here, $f(\mathbf{x}) = (1/2) \|A\mathbf{x} - b\|_2^2$ and $g_i(x_i) = \lambda |x_i|$.

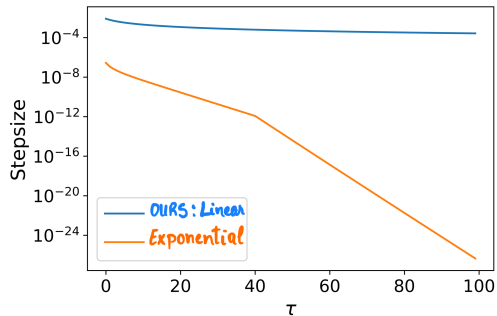
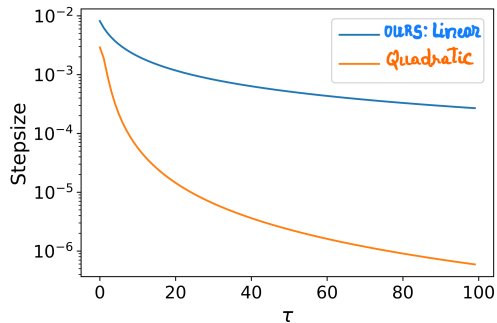
Comparing to existing works



Comparing to existing works



Comparing to existing works



Remarks on the delay vector

- ▶ inconsistent read

Remarks on the delay vector

- ▶ inconsistent read
- ▶ Need τ to fix the stepsize.

Remarks on the delay vector

- ▶ inconsistent read
- ▶ Need τ to fix the stepsize.
- ▶ independence of the coordinates

Outline

General introduction

Asynchronous Forward-Backward

Variance reduction techniques for SPPA

Conclusion

Problem

Recall the first case:

$$\underset{\mathbf{x} \in \mathbf{H}}{\text{minimize}} \quad F(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}),$$

where for all $i \in [n]$, $f_i : \mathbf{H} \rightarrow \mathbb{R}$.

SGD

For all $k \in \mathbb{N}$,

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \nabla f_{i_k}(\mathbf{x}^k)$$

SGD

For all $k \in \mathbb{N}$,

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \nabla f_{i_k}(\mathbf{x}^k)$$

► $\gamma_k \rightarrow 0$.

SGD

For all $k \in \mathbb{N}$,

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \nabla f_{i_k}(\mathbf{x}^k)$$

► $\gamma_k \rightarrow 0$. For instance, $\gamma_k = \frac{\gamma_0}{k^\beta}$ with $\beta \in [1/2, 1]$.

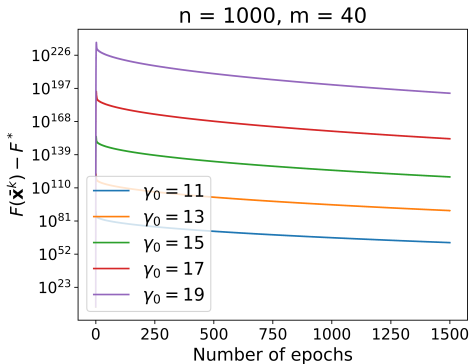
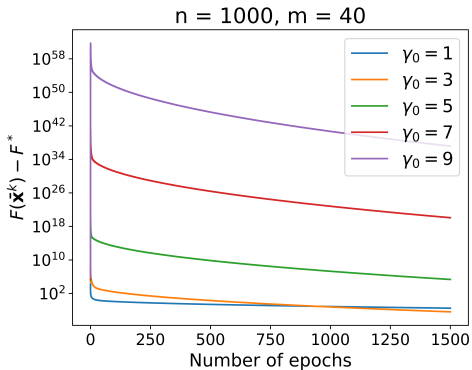
SGD

For all $k \in \mathbb{N}$,

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma_k \nabla f_{i_k}(\mathbf{x}^k)$$

- ▶ $\gamma_k \rightarrow 0$. For instance, $\gamma_k = \frac{\gamma_0}{k^\beta}$ with $\beta \in [1/2, 1]$.
- ▶ Rates in expectation: $O\left(\frac{1}{\sqrt{k}}\right)$ for the convex case and $O(1/k)$ for the strongly convex case.

SGD unstable w.r.t. γ_0



How can we alleviate the instability with respect to γ_0 ?

How can we alleviate the instability with respect to γ_0 ?

By using, for example, stochastic proximal point algorithm (SPPA)!

SPPA

For all $k \in \mathbb{N}$,

$$\boldsymbol{x}^{k+1} = \text{prox}_{\gamma_k f_{i_k}}(\boldsymbol{x}^k).$$

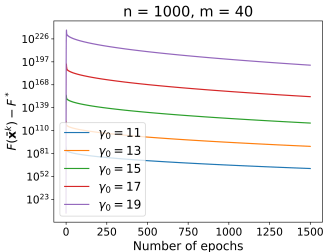
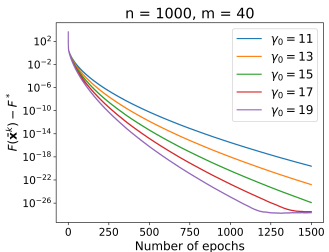
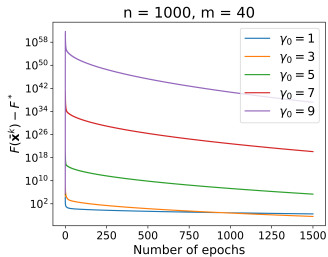
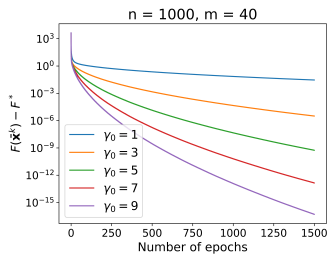
SPPA

For all $k \in \mathbb{N}$,

$$\mathbf{x}^{k+1} = \text{prox}_{\gamma_k f_{i_k}}(\mathbf{x}^k).$$

Converges with same stepsize rule as SGD!!

SPPA more stable



(a) SPPA

(b) SGD

SPPA more stable

From Asi and Duchi⁷:

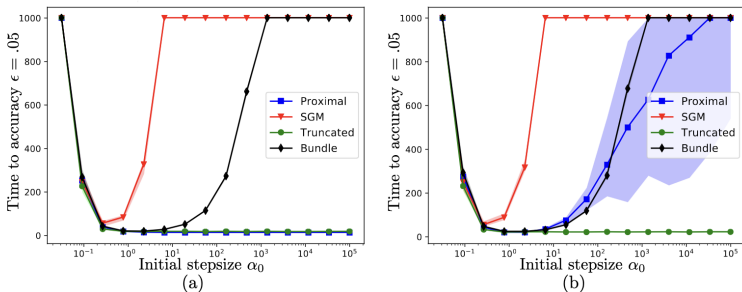


Figure 2. The number of iterations to achieve ϵ -accuracy versus initial stepsize α_0 for linear regression with $m = 1000$, $n = 40$, and condition number $\kappa(A) = 1$. (a) The noiseless setting with $\sigma = 0$. (b) Noisy setting with $\sigma = \frac{1}{2}$.

⁷Asi et al., “Stochastic (approximate) proximal point methods: convergence, optimality, and adaptivity”.

What about the rates ?

What about the rates ?

Does SPPA help recover the full GD rates ?

What about the rates ?

Does SPPA help recover the full GD rates ?

No! Same as SGD.

One problem with SPPA⁸ and SGD⁹ bounds

$\forall i \in [n]$ f_i convex and L -smooth:

$$\mathbb{E}[F(\bar{x}^k) - F_*] \leq \frac{\text{dist}(x^0, \text{argmin } F)^2}{\sum_{t=0}^{k-1} \gamma_t} + 2\sigma^2 \frac{\sum_{t=0}^{k-1} \gamma_t^2}{\sum_{t=0}^{k-1} \gamma_t},$$

⁸Traoré, Apidopoulos, et al., “Variance reduction techniques for stochastic proximal point algorithms”.

⁹Garrigos et al., “Handbook of convergence theorems for (stochastic) gradient methods”.

One problem with SPPA⁸ and SGD⁹ bounds

$\forall i \in [n]$ f_i convex and L -smooth:

$$\mathbb{E}[F(\bar{\mathbf{x}}^k) - F_*] \leq \frac{\text{dist}(\mathbf{x}^0, \text{argmin } F)^2}{\sum_{t=0}^{k-1} \gamma_t} + 2\sigma^2 \frac{\sum_{t=0}^{k-1} \gamma_t^2}{\sum_{t=0}^{k-1} \gamma_t},$$

$$\sigma^2 := \sup_{\mathbf{x}^* \in \text{argmin } F} \mathbb{E} \|\nabla f_i(\mathbf{x}^*) - \mathbb{E}[\nabla f_i(\mathbf{x}^*)]\|^2.$$

⁸Traoré, Apidopoulos, et al., “Variance reduction techniques for stochastic proximal point algorithms”.

⁹Garrigos et al., “Handbook of convergence theorems for (stochastic) gradient methods”.

"Can we do better ?", Dr. Cris Vega, modern day philosopher.

"Can we do better ?", Dr. Cris Vega, modern day philosopher.

Yes, by using variance reduction techniques!

A way to reduce the variance of a R.V.

- ▶ X a random variable (R.V.).

A way to reduce the variance of a R.V.

- ▶ X a random variable (R.V.).
 - Goal: reduce variance of X .

A way to reduce the variance of a R.V.

- ▶ X a random variable (R.V.).
 - Goal: reduce variance of X .
- ▶ Given Z , easy-to-compute $E[Z]$.

A way to reduce the variance of a R.V.

- ▶ X a random variable (R.V.).
 - Goal: reduce variance of X .
- ▶ Given Z , easy-to-compute $E[Z]$.
- ▶ Define

$$X_Z := X - Z + E[Z].$$

A way to reduce the variance of a R.V.

- ▶ X a random variable (R.V.).
 - Goal: reduce variance of X .
- ▶ Given Z , easy-to-compute $E[Z]$.
- ▶ Define

$$X_Z := X - Z + E[Z].$$

- ▶ We want

$$\text{Var}(X_Z) < \text{Var}(X).$$

A way to reduce the variance of a R.V.

- ▶ X a random variable (R.V.).
 - Goal: reduce variance of X .
- ▶ Given Z , easy-to-compute $E[Z]$.

- ▶ Define

$$X_Z := X - Z + E[Z].$$

- ▶ We want

$$\text{Var}(X_Z) < \text{Var}(X).$$

- ▶ We know

$$\text{Var}(X_Z) = \text{Var}(X) + \text{Var}(Z) - 2\text{Cov}(X, Z).$$

A way to reduce the variance of a R.V.

- ▶ X a random variable (R.V.).
 - Goal: reduce variance of X .
- ▶ Given Z , easy-to-compute $E[Z]$.

- ▶ Define

$$X_Z := X - Z + E[Z].$$

- ▶ We want

$$\text{Var}(X_Z) < \text{Var}(X).$$

- ▶ We know

$$\text{Var}(X_Z) = \text{Var}(X) + \text{Var}(Z) - 2\text{Cov}(X, Z).$$

- ▶ It is sufficient to have $\text{Cov}(X, Z) > \frac{1}{2}\text{Var}(Z)$.

General variance reduction scheme for SGD

- At each iteration k , $X = \nabla f_{i_k}(\mathbf{x}^k)$ and $\mathbb{E}[X|\mathbf{x}^k] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}^k) = \nabla F(\mathbf{x}^k)$.

General variance reduction scheme for SGD

- At each iteration k , $X = \nabla f_{i_k}(\mathbf{x}^k)$ and $\mathbb{E}[X|\mathbf{x}^k] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}^k) = \nabla F(\mathbf{x}^k)$.
- Goal: find a good Z and replace $X = \nabla f_{i_k}(\mathbf{x}^k)$ by

$$X_Z := \nabla f_{i_k}(\mathbf{x}^k) - Z + \mathbb{E}[Z|\mathbf{x}^k].$$

General variance reduction scheme for SGD

- ▶ At each iteration k , $X = \nabla f_{i_k}(\mathbf{x}^k)$ and $\mathbb{E}[X|\mathbf{x}^k] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}^k) = \nabla F(\mathbf{x}^k)$.
- ▶ Goal: find a good Z and replace $X = \nabla f_{i_k}(\mathbf{x}^k)$ by

$$X_Z := \nabla f_{i_k}(\mathbf{x}^k) - Z + \mathbb{E}[Z|\mathbf{x}^k].$$

- ▶ Depending on the choice of Z and how it is used, we get different algorithms.

At iteration k , $Z = \nabla f_{i_k}(\tilde{\mathbf{x}} = \mathbf{x}^{k-d})$. Then $X_Z := \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \mathbb{E}[\nabla f_{i_k}(\tilde{\mathbf{x}}) | \mathbf{x}^k]$.

Stochastic variance reduced gradient (SVRG)¹⁰

At iteration k , $Z = \nabla f_{i_k}(\tilde{\mathbf{x}} = \mathbf{x}^{k-d})$. Then $X_Z := \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}) + \mathbb{E}[\nabla f_{i_k}(\tilde{\mathbf{x}}) | \mathbf{x}^k]$.

Algorithm

Let $\gamma > 0$ and set $\tilde{\mathbf{x}}^0 \in \mathcal{H}$. Then

for $s = 0, 1, \dots$

$\mathbf{x}^0 = \tilde{\mathbf{x}}^s$, compute $\nabla F(\tilde{\mathbf{x}}^s)$

for $k = 0, \dots, \ell - 1$

choose i_k uniformly at random in $[n]$

$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma (\nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\tilde{\mathbf{x}}^s) + \nabla F(\tilde{\mathbf{x}}^s))$

choose ξ_s uniformly at random in $\{0, 1, \dots, \ell - 1\}$

$\tilde{\mathbf{x}}^{s+1} = \sum_{k=0}^{\ell-1} \delta_{k, \xi_s} \mathbf{x}^k$,

where $\delta_{k,h}$ is the Kronecker symbol.

¹⁰Johnson et al., "Accelerating stochastic gradient descent using predictive variance reduction".

At iteration k , $Z = \nabla f_{i_k}(\phi_{i_k}^k = \mathbf{x}^{k-d})$, with $f_{i_{k-d}} = f_{i_k}$. Then

$$X_Z := \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\phi_{i_k}^k) + \mathbb{E}[\nabla f_{i_k}(\phi_{i_k}^k) \mid \mathbf{x}^k].$$

Stochastic average gradient algorithm (SAGA)¹¹

At iteration k , $Z = \nabla f_{i_k}(\phi_{i_k}^k = \mathbf{x}^{k-d})$, with $f_{i_{k-d}} = f_{i_k}$. Then

$$X_Z := \nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\phi_{i_k}^k) + \mathbb{E}[\nabla f_{i_k}(\phi_{i_k}^k) | \mathbf{x}^k].$$

Algorithm

Let $\gamma > 0$. Set $\mathbf{x}^0 \in \mathcal{H}$ and, $\forall i \in [n]$, $\phi_i^0 = \mathbf{x}^0$. Then

for $k = 0, 1, \dots$

$$\left[\begin{array}{l} \text{choose } i_k \text{ uniformly at random in } [n] \\ \mathbf{x}^{k+1} = \mathbf{x}^k - \gamma \left(\nabla f_{i_k}(\mathbf{x}^k) - \nabla f_{i_k}(\phi_{i_k}^k) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\phi_i^k) \right), \\ \forall i \in [n]: \phi_i^{k+1} = \phi_i^k + \delta_{i,i_k}(\mathbf{x}^k - \phi_i^k), \end{array} \right.$$

where $\delta_{i,j}$ is the Kronecker symbol.

¹¹Defazio et al., "SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives".

Variance reduction worked for SGD

Algorithm	Rates for strongly convex F	Cost of 1 Iteration
GD	$O(\varepsilon^k)$	$O(n)$
SGD	$O(1/k)$	$O(1)$
V.R.	$O(\varepsilon^k)$	$O(1)$

$O(1)$ = cost of $\nabla f_i(\mathbf{x})$.

Our goal: variance reduction for SPPA.

Algorithm proposed¹²

Algorithm (Generic)

Let $\gamma > 0$ and $x^0 \in H$. Then

for $k = 0, 1, \dots$

 | choose i_k uniformly at random in $[n]$
 | $x^{k+1} = \text{prox}_{\gamma f_{i_k}}(x^k + \gamma e^k)$.

¹²Traoré, Apidopoulos, et al., “Variance reduction techniques for stochastic proximal point algorithms”.

Algorithm proposed¹²

Algorithm (Generic)

Let $\gamma > 0$ and $x^0 \in H$. Then

for $k = 0, 1, \dots$

$\left[\begin{array}{l} \text{choose } i_k \text{ uniformly at random in } [n] \\ x^{k+1} = \text{prox}_{\gamma f_{i_k}}(x^k + \gamma e^k). \end{array} \right.$

Let $w^k = \nabla f_{i_k}(x^{k+1}) - e^k$.

¹²Traoré, Apidopoulos, et al., “Variance reduction techniques for stochastic proximal point algorithms”.

Algorithm proposed¹²

Algorithm (Generic)

Let $\gamma > 0$ and $\mathbf{x}^0 \in \mathcal{H}$. Then

for $k = 0, 1, \dots$

 | choose i_k uniformly at random in $[n]$
 | $\mathbf{x}^{k+1} = \text{prox}_{\gamma f_{i_k}}(\mathbf{x}^k + \gamma \mathbf{e}^k)$.

Let $\mathbf{w}^k = \nabla f_{i_k}(\mathbf{x}^{k+1}) - \mathbf{e}^k$.

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma [\nabla f_{i_k}(\mathbf{x}^{k+1}) - \mathbf{e}^k] = \mathbf{x}^k - \gamma \mathbf{w}^k.$$

¹²Traoré, Apidopoulos, et al., “Variance reduction techniques for stochastic proximal point algorithms”.

Algorithm proposed¹²

Algorithm (Generic)

Let $\gamma > 0$ and $\mathbf{x}^0 \in \mathcal{H}$. Then

for $k = 0, 1, \dots$

 | choose i_k uniformly at random in $[n]$
 | $\mathbf{x}^{k+1} = \text{prox}_{\gamma f_{i_k}}(\mathbf{x}^k + \gamma \mathbf{e}^k)$.

Let $\mathbf{w}^k = \nabla f_{i_k}(\mathbf{x}^{k+1}) - \mathbf{e}^k$.

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \gamma [\nabla f_{i_k}(\mathbf{x}^{k+1}) - \mathbf{e}^k] = \mathbf{x}^k - \gamma \mathbf{w}^k.$$

For the analysis we consider

$$\mathbf{v}^k := \nabla f_{i_k}(\mathbf{x}^k) - \mathbf{e}^k.$$

¹²Traoré, Apidopoulos, et al., “Variance reduction techniques for stochastic proximal point algorithms”.

Assumptions

Let $A, B, D \in \mathbb{R}_+$ and $\rho \in [0, 1]$ and a real-valued random variable C such that, for every $k \in \mathbb{N}$,

1. $E[e^k | \mathfrak{F}_k] = 0$ a.s.,

Assumptions

Let $A, B, D \in \mathbb{R}_+$ and $\rho \in [0, 1]$ and a real-valued random variable C such that, for every $k \in \mathbb{N}$,

1. $\mathbb{E}[e^k \mid \mathfrak{F}_k] = 0$ a.s.,
2. $\mathbb{E}[\|\mathbf{v}^k\|^2 \mid \mathfrak{F}_k] \leq 2A(F(\mathbf{x}^k) - F_*) + B\sigma_k^2 + C$ a.s.,

Assumptions

Let $A, B, D \in \mathbb{R}_+$ and $\rho \in [0, 1]$ and a real-valued random variable C such that, for every $k \in \mathbb{N}$,

1. $\mathbb{E}[e^k \mid \mathfrak{F}_k] = 0$ a.s.,
2. $\mathbb{E}[\|\mathbf{v}^k\|^2 \mid \mathfrak{F}_k] \leq 2A(F(\mathbf{x}^k) - F_*) + B\sigma_k^2 + C$ a.s.,
3. $\mathbb{E}[\sigma_{k+1}^2] \leq (1 - \rho)\mathbb{E}[\sigma_k^2] + 2D\mathbb{E}[F(\mathbf{x}^k) - F_*]$.

Assumptions

1. $\operatorname{argmin} F \neq \emptyset$.
2. For all $i \in [n]$, f_i is convex and, moreover, L -smooth, i.e., differentiable and such that

$$(\forall \mathbf{x}, \mathbf{y} \in \mathbf{H}) \quad \|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$$

for some $L > 0$. As a consequence, F is convex and L -smooth.

3. F satisfies the PL condition with constant $\mu > 0$, i.e.,

$$(\forall \mathbf{x} \in \mathbf{H}) \quad F(\mathbf{x}) - F_* \leq \frac{1}{2\mu} \|\nabla F(\mathbf{x})\|^2.$$

Common result

Proposition

Let $M > 0$. Then, for all $k \in \mathbb{N}$,

$$\begin{aligned} \mathbb{E}[\text{dist}(\mathbf{x}^{k+1}, \text{argmin } F)^2] + \gamma^2 M \mathbb{E}[\sigma_{k+1}^2] &\leq \mathbb{E}[\text{dist}(\mathbf{x}^k, \text{argmin } F)^2] + \gamma^2 [M + B - \rho M] \mathbb{E}[\sigma_k^2] \\ &\quad - 2\gamma [1 - \gamma(A + MD)] \mathbb{E}[F(\mathbf{x}^k) - F_*] \\ &\quad + \gamma^2 \mathbb{E}[C]. \end{aligned}$$

Convex F

Theorem

Let $M > 0$ and $\gamma > 0$ be such that $M \geq B/\rho$ and $\gamma < 1/(A + MD)$. Then,

$$(\forall k \in \mathbb{N}) \quad \mathbb{E}[F(\bar{\mathbf{x}}^k) - F_*] \leq \frac{\text{dist}(\mathbf{x}^0, \text{argmin } F)^2 + \gamma^2 M \mathbb{E}[\sigma_0^2]}{2\gamma k [1 - \gamma(A + MD)]},$$

with $\bar{\mathbf{x}}^k = \frac{1}{k} \sum_{t=0}^{k-1} \mathbf{x}^t$.

F satisfying PL

Theorem

Let M be such that $M > B/\rho$ and $\gamma > 0$ such that $\gamma < 1/(A + MD)$.

Set

$$q := \max \left\{ 1 - \gamma\mu(1 - \gamma(A + MD)), 1 + \frac{B}{M} - \rho \right\}.$$

Then $q \in]0, 1[$ and

$$(\forall k \in \mathbb{N}) \quad \mathbb{E}[\text{dist}(\mathbf{x}^k, \text{argmin } F)^2] \leq q^k (\text{dist}(\mathbf{x}^0, \text{argmin } F)^2 + \gamma^2 M \mathbb{E}[\sigma_0^2]),$$

$$\mathbb{E}[F(\mathbf{x}^k) - F_*] \leq \frac{q^k L}{2} (\text{dist}(\mathbf{x}^0, \text{argmin } F)^2 + \gamma^2 M \mathbb{E}[\sigma_0^2]).$$

SVRP

Algorithm (SVRP)

Let $\gamma > 0$ and $\tilde{\mathbf{x}}^0 \in \mathcal{H}$. Then

for $s = 0, 1, \dots$

$$\left[\begin{array}{l} \mathbf{x}^0 = \tilde{\mathbf{x}}^s, \text{ compute } \nabla F(\tilde{\mathbf{x}}^s) \\ \text{for } k = 0, \dots, \ell - 1 \\ \quad \left[\begin{array}{l} \text{choose } i_k \text{ uniformly at random in } [n] \\ \mathbf{x}^{k+1} = \text{prox}_{\gamma f_{i_k}} (\mathbf{x}^k + \gamma \nabla f_{i_k}(\tilde{\mathbf{x}}^s) - \gamma \nabla F(\tilde{\mathbf{x}}^s)) \end{array} \right. \\ \text{choose } \xi_s \text{ uniformly at random in } \{0, 1, \dots, \ell - 1\} \\ \left. \tilde{\mathbf{x}}^{s+1} = \sum_{k=0}^{\ell-1} \delta_{k, \xi_s} \mathbf{x}^k, \right. \end{array} \right.$$

where $\delta_{k,h}$ is the Kronecker symbol.

SVRP results

Theorem (PL case)

Suppose that

$$0 < \gamma < \frac{1}{2(2L - \mu)} \quad \text{and} \quad \ell > \frac{1}{\mu\gamma(1 - 2\gamma(2L - \mu))}.$$

Then

$$(\forall s \in \mathbb{N}) \quad \mathbb{E} [F(\tilde{\mathbf{x}}^{s+1}) - F_*] \leq q^s (F(\mathbf{x}^0) - F_*),$$

with $q := \left(\frac{1}{\mu\gamma(1 - 2L\gamma)\ell} + \frac{2\gamma(L - \mu)}{1 - 2L\gamma} \right) < 1$.

L-SVRP

Algorithm (L-SVRP)

Let $\gamma > 0$ and set $\mathbf{x}^0 = \mathbf{u}^0 \in \mathbf{H}$. Then

for $k = 0, 1, \dots$

$$\left[\begin{array}{l} \text{choose } i_k \text{ uniformly at random in } [n] \\ \mathbf{x}^{k+1} = \text{prox}_{\gamma f_{i_k}} (\mathbf{x}^k + \gamma \nabla f_{i_k}(\mathbf{u}^k) - \gamma \nabla F(\mathbf{u}^k)) \\ \varepsilon^k \text{ Bernoulli r.v. with } P(\varepsilon^k = 1) = p \in]0, 1] \\ \mathbf{u}^{k+1} = (1 - \varepsilon^k) \mathbf{u}^k + \varepsilon^k \mathbf{x}^k, \end{array} \right.$$

L-SVRP results

Corollary (Convex case)

Let $M \geq \frac{2}{p}$ and $\gamma < \frac{1}{L(2+pM)}$. Then

$$(\forall k \in \mathbb{N}) \quad \mathbb{E}[F(\bar{\mathbf{x}}^k) - F_*] \leq \frac{\text{dist}(\mathbf{x}^0, \text{argmin } F)^2 + \gamma^2 M \mathbb{E}[\sigma_0^2]}{2\gamma k [1 - \gamma L(2 + pM)]},$$

with $\bar{\mathbf{x}}^k = \frac{1}{k} \sum_{t=0}^{k-1} \mathbf{x}^t$.

L-SVRP results

Corollary (PL case)

Let $M > 2/p$ and $\gamma < \frac{1}{L(2+pM)}$. Then

$$(\forall k \in \mathbb{N}) \quad \mathbb{E}[\text{dist}(\mathbf{x}^k, \text{argmin } F)]^2 \leq q^k (\text{dist}(\mathbf{x}^0, \text{argmin } F)^2 + \gamma^2 M \mathbb{E}[\sigma_0^2]) ,$$

$$\mathbb{E}[F(\mathbf{x}^k) - F_*] \leq \frac{q^k L}{2} (\text{dist}(\mathbf{x}^0, \text{argmin } F)^2 + \gamma^2 M \mathbb{E}[\sigma_0^2]) ,$$

with $0 < q < 1$.

SAPA

Algorithm (SAPA)

Let $\gamma > 0$. Set $\mathbf{x}^0 \in \mathbf{H}$ and, $\forall i \in [n]$, $\phi_i^0 = \mathbf{x}^0$. Then

for $k = 0, 1, \dots$

$$\left[\begin{array}{l} \text{choose } i_k \text{ uniformly at random in } [n] \\ \mathbf{x}^{k+1} = \text{prox}_{\gamma f_{i_k}} \left(\mathbf{x}^k + \gamma \nabla f_{i_k}(\phi_{i_k}^k) - \frac{\gamma}{n} \sum_{i=1}^n \nabla f_i(\phi_i^k) \right) \\ \forall i \in [n]: \phi_i^{k+1} = \phi_i^k + \delta_{i,i_k} (\mathbf{x}^k - \phi_i^k), \end{array} \right.$$

where $\delta_{i,j}$ is the Kronecker symbol.

SAPA results

Corollary (Convex case)

Let $M \geq 2n$ and $\gamma < \frac{1}{L(2+M/n)}$. Then

$$(\forall k \in \mathbb{N}) \quad \mathbb{E}[F(\bar{\mathbf{x}}^k) - F_*] \leq \frac{\text{dist}(\mathbf{x}^0, \text{argmin } F)^2 + \gamma^2 M \mathbb{E}[\sigma_0^2]}{2\gamma k [1 - \gamma L(2 + M/n)]},$$

with $\bar{\mathbf{x}}^k = \frac{1}{k} \sum_{t=0}^{k-1} \mathbf{x}^t$.

SAPA results

Corollary (PL case)

Let $M > 2n$ and $\gamma < \frac{1}{L(2+M/n)}$. Then

$$(\forall k \in \mathbb{N}) \quad \mathbb{E}[\text{dist}(\mathbf{x}^k, \text{argmin } F)]^2 \leq q^k (\text{dist}(\mathbf{x}^0, \text{argmin } F)^2 + \gamma^2 M \mathbb{E}[\sigma_0^2]) ,$$

$$\mathbb{E}[F(\mathbf{x}^k) - F_*] \leq \frac{q^k L}{2} (\text{dist}(\mathbf{x}^0, \text{argmin } F)^2 + \gamma^2 M \mathbb{E}[\sigma_0^2]) ,$$

with $0 < q < 1$.

Related works

	Algorithm	Smooth + convex	Smooth + SC	Smooth + PL	Non-smooth + SC
Defazio ¹³	Point-Saga	NA	$O(\varepsilon^k)$	NA	$O(1/k)$
Khaled et al. ¹⁴	L-SVRP	NA	$O(\varepsilon^k)$	NA	NA
Milzarek et al. ¹⁵	SNSPP	NA	$O(\varepsilon^k)$	NA	NA
Traoré et al. ¹⁶	Unified	$O(1/k)$ (not for SVRP)	$O(\varepsilon^k)$	$O(\varepsilon^k)$	NA

Table: Comparison to related works.

¹³Defazio, “A simple practical accelerated method for finite sums”.

¹⁴Khaled et al., “Faster federated optimization under second-order similarity”.

¹⁵Milzarek et al., “A semismooth Newton stochastic proximal point algorithm with variance reduction”.

¹⁶Traoré, Apidopoulos, et al., “Variance reduction techniques for stochastic proximal point algorithms”.

Experiments

Ordinary least squares (OLS):

$$\underset{\mathbf{x} \in \mathbb{R}^m}{\text{minimize}} \quad F(\mathbf{x}) = \frac{1}{2n} \|A\mathbf{x} - \mathbf{b}\|^2 = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (\langle \mathbf{a}_i, \mathbf{x} \rangle - b_i)^2,$$

where \mathbf{a}_i is the i^{th} row of the matrix $A \in \mathbb{R}^{n \times m}$ and $b_i \in \mathbb{R}$ for all $i \in [n]$

SAPA/SVRP/SPPA

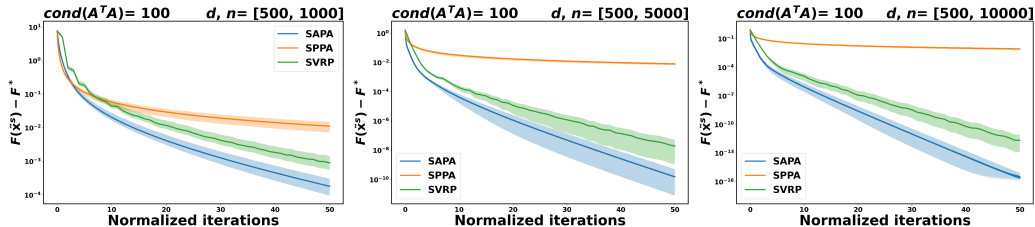


Figure: SAPA (blue) and SVRP (green) compared to SPPA (orange).

SAGA/SAPA

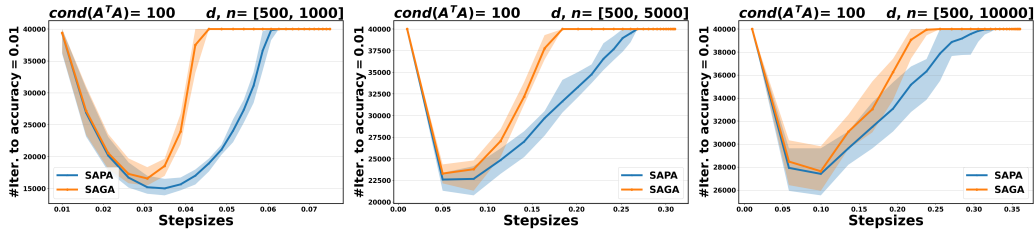


Figure: SAPA (in blue) compared to SAGA (in orange).

SVRG/SVRP

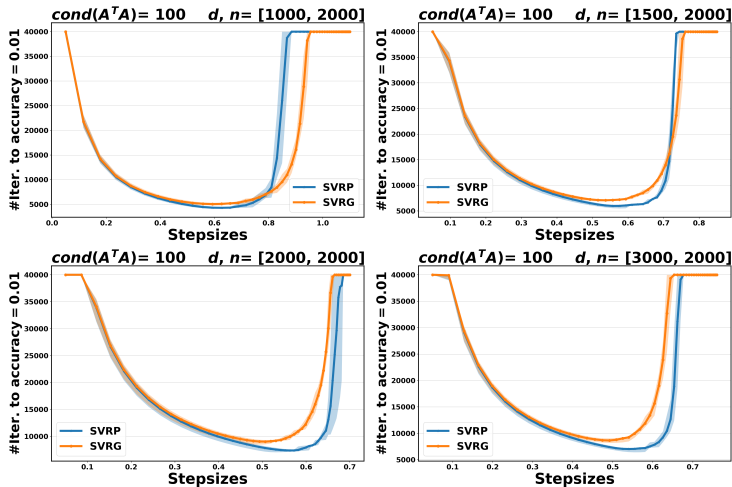


Figure: SVRP (in blue) compared SVRG (in orange).

Remarks on the results

- ▶ We only have results for the smooth case.

Remarks on the results

- ▶ We only have results for the smooth case.
- ▶ For SVRP, no results for the convex case.

Remarks on the results

- ▶ We only have results for the smooth case.
- ▶ For SVRP, no results for the convex case.
- ▶ We have derived results for SPPA from the generic algorithm.

Outline

General introduction

Asynchronous Forward-Backward

Variance reduction techniques for SPPA

Conclusion

Conclusion

Asynchronous Forward-Backward

Variance reduction for SPPA

Conclusion

Asynchronous Forward-Backward

Summary

- ▶ considered with abstract probability and coordinate-wise adaptive stepsize
- ▶ provided convergence of the iterates
- ▶ provided standard convergence rates for convex and error bound cases
- ▶ results depend linearly on τ

Variance reduction for SPPA

Summary

- ▶ made a unified study
- ▶ proved standard rates for convex and PL cases
- ▶ recovered proximal version of standard gradient variance reduced algorithms
- ▶ empirical comparison

Conclusion

Asynchronous Forward-Backward

Summary

- ▶ considered with abstract probability and coordinate-wise adaptive stepsize
- ▶ provided convergence of the iterates
- ▶ provided standard convergence rates for convex and error bound cases
- ▶ results depend linearly on τ

Future directions

- ▶ results for delay-wise adaptive stepsize
- ▶ considered coordinates dependent delay
- ▶ prove the tightest delay dependence

Variance reduction for SPPA

Summary

- ▶ made a unified study
- ▶ proved standard rates for convex and PL cases
- ▶ recovered proximal version of standard gradient variance reduced algorithms
- ▶ empirical comparison

Future directions

- ▶ convex case for SVRP
- ▶ non-smooth analysis
- ▶ Bregman and/or zeroth-order versions

Conclusion

Asynchronous Forward-Backward

Summary

- ▶ considered with abstract probability and coordinate-wise adaptive stepsize
- ▶ provided convergence of the iterates
- ▶ provided standard convergence rates for convex and error bound cases
- ▶ results depend linearly on τ

Future directions

- ▶ results for delay-wise adaptive stepsize
- ▶ considered coordinates dependent delay
- ▶ prove the tightest delay dependence

Variance reduction for SPPA

Summary

- ▶ made a unified study
- ▶ proved standard rates for convex and PL cases
- ▶ recovered proximal version of standard gradient variance reduced algorithms
- ▶ empirical comparison

Future directions

- ▶ convex case for SVRP
- ▶ non-smooth analysis
- ▶ Bregman and/or zeroth-order versions

Conclusion

Asynchronous Forward-Backward

Summary

- ▶ considered with abstract probability and coordinate-wise adaptive stepsize
- ▶ provided convergence of the iterates
- ▶ provided standard convergence rates for convex and error bound cases
- ▶ results depend linearly on τ

Future directions

- ▶ results for delay-wise adaptive stepsize
- ▶ considered coordinates dependent delay
- ▶ prove the tightest delay dependence

Variance reduction for SPPA

Summary

- ▶ made a unified study
- ▶ proved standard rates for convex and PL cases
- ▶ recovered proximal version of standard gradient variance reduced algorithms
- ▶ empirical comparison

Future directions

- ▶ convex case for SVRP
- ▶ non-smooth analysis
- ▶ Bregman and/or zeroth-order versions

Thank you for your attention!!