

# ICML | 2019

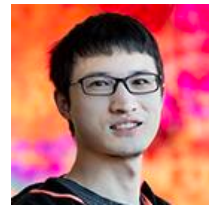
Thirty-sixth International Conference on  
Machine Learning



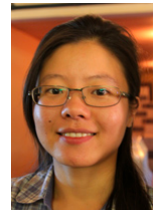
香港大學  
THE UNIVERSITY OF HONG KONG



## SAGA with Arbitrary Sampling



Xun Qian



Zheng Qu



Peter Richtárik



# The Problem

# The Problem: Regularized Empirical Risk Minimization

$$\min_{x \in \mathbb{R}^d} P(x) \stackrel{\text{def}}{=} \left( \sum_{i=1}^n \lambda_i f_i(x) \right) + \psi(x)$$

# The Problem: Regularized Empirical Risk Minimization

$$\min_{x \in \mathbb{R}^d} P(x) \stackrel{\text{def}}{=} \underbrace{\left( \sum_{i=1}^n \lambda_i f_i(x) \right)}_{f(x)} + \psi(x)$$

# The Problem: Regularized Empirical Risk Minimization

$$\min_{x \in \mathbb{R}^d} P(x) \stackrel{\text{def}}{=} \underbrace{\left( \sum_{i=1}^n \lambda_i f_i(x) \right)}_{f(x)} + \psi(x)$$

Regularizer

# The Problem: Regularized Empirical Risk Minimization

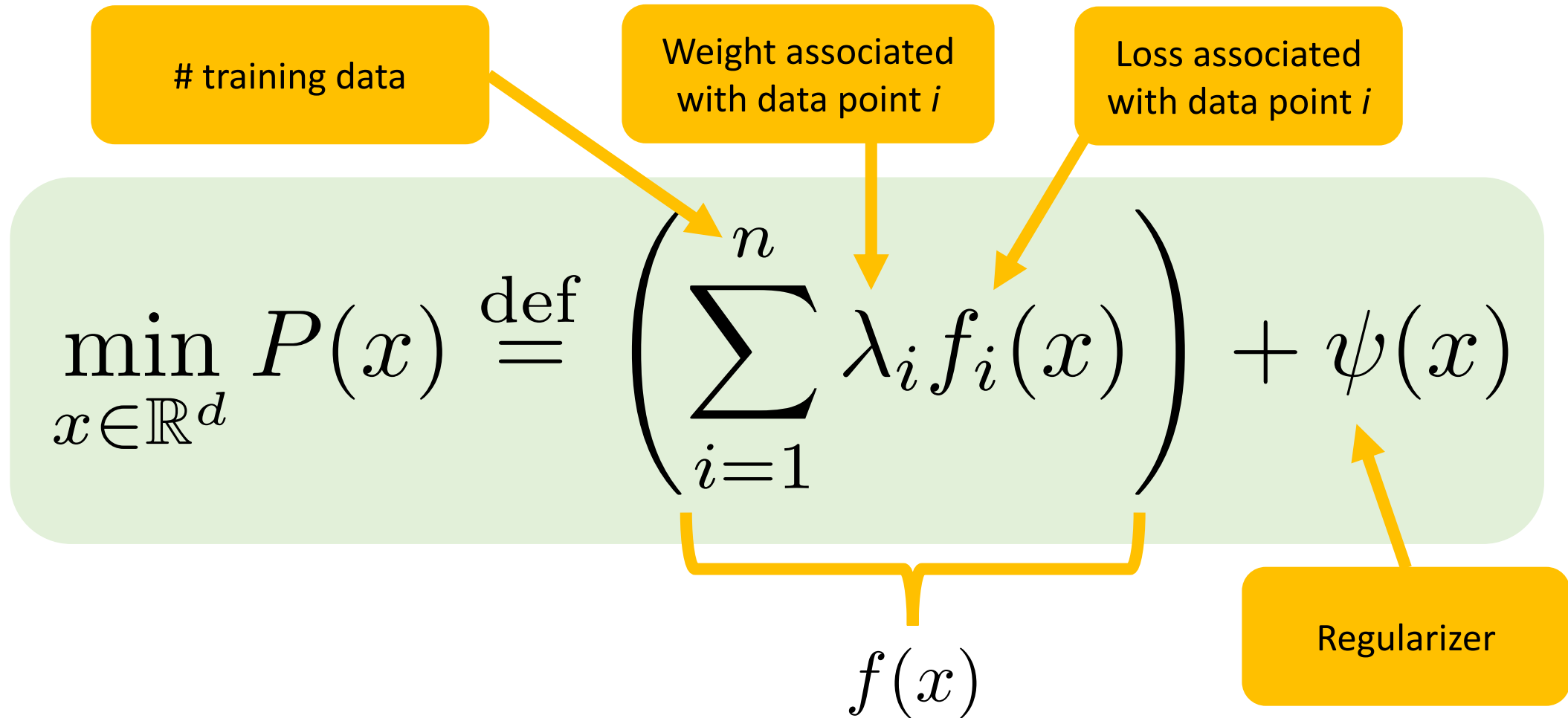
# training data

$$\min_{x \in \mathbb{R}^d} P(x) \stackrel{\text{def}}{=} \left( \sum_{i=1}^n \lambda_i f_i(x) \right) + \psi(x)$$

$f(x)$

Regularizer

# The Problem: Regularized Empirical Risk Minimization



# The Problem: Regularized Empirical Risk Minimization

# training data

Weight associated  
with data point  $i$

Loss associated  
with data point  $i$

$$\min_{x \in \mathbb{R}^d} P(x) \stackrel{\text{def}}{=} \left( \sum_{i=1}^n \lambda_i f_i(x) \right) + \psi(x)$$

Parameters  
describing the model

$f(x)$

Regularizer



# Arbitrary Sampling

# SGD with Arbitrary Sampling

# SGD with Arbitrary Sampling

- 1 In iteration  $k$ , we have  $x^k$  available

# SGD with Arbitrary Sampling

- 1 In iteration  $k$ , we have  $x^k$  available
- 2 Sample a random set  $S_k \subseteq \{1, 2, \dots, n\}$

# SGD with Arbitrary Sampling

- 1 In iteration  $k$ , we have  $x^k$  available
- 2 Sample a random set  $S_k \subseteq \{1, 2, \dots, n\}$
- 3 Compute the gradients  $\nabla f_i(x^k)$  for  $i \in S_k$  only

# SGD with Arbitrary Sampling

- 1 In iteration  $k$ , we have  $x^k$  available
- 2 Sample a random set  $S_k \subseteq \{1, 2, \dots, n\}$
- 3 Compute the gradients  $\nabla f_i(x^k)$  for  $i \in S_k$  only
- 4 Approximate the gradient  $\nabla f(x^k)$  using  $\{\nabla f_i(x^k) : i \in S_k\}$

# SGD with Arbitrary Sampling

- 1 In iteration  $k$ , we have  $x^k$  available
- 2 Sample a random set  $S_k \subseteq \{1, 2, \dots, n\}$
- 3 Compute the gradients  $\nabla f_i(x^k)$  for  $i \in S_k$  only
- 4 Approximate the gradient  $\nabla f(x^k)$  using  $\{\nabla f_i(x^k) : i \in S_k\}$
- 5 Take a stochastic gradient descent step to obtain  $x^{k+1}$

# SGD with Arbitrary Sampling

- 1 In iteration  $k$ , we have  $x^k$  available
- 2 Sample a random set  $S_k \subseteq \{1, 2, \dots, n\}$
- 3 Compute the gradients  $\nabla f_i(x^k)$  for  $i \in S_k$  only
- 4 Approximate the gradient  $\nabla f(x^k)$  using  $\{\nabla f_i(x^k) : i \in S_k\}$
- 5 Take a stochastic gradient descent step to obtain  $x^{k+1}$

**Arbitrary sampling paradigm** (R. & Takáč 2013): want to be able to sample from **any** distribution over all  $2^n$  subsets of  $\{1, 2, \dots, n\}$

$$p_i \stackrel{\text{def}}{=} \text{Prob}(i \in S_k)$$

$$p_i > 0 \text{ for all } i = 1, 2, \dots, n$$



# Arbitrary Sampling: Examples for $n = 3$

GD

$$S_k = \{1, 2, 3\} \text{ with prob } 1$$

# Arbitrary Sampling: Examples for $n = 3$

GD

$$S_k = \{1, 2, 3\} \text{ with prob } 1$$

SAGA

$$S_k = \{1\} \text{ with prob } 1/3$$

$$S_k = \{2\} \text{ with prob } 1/3$$

$$S_k = \{3\} \text{ with prob } 1/3$$

# Arbitrary Sampling: Examples for $n = 3$

GD

$$S_k = \{1, 2, 3\} \text{ with prob } 1$$

SAGA

$$S_k = \{1\} \text{ with prob } 1/3$$

$$S_k = \{2\} \text{ with prob } 1/3$$

$$S_k = \{3\} \text{ with prob } 1/3$$

SAGA with nonuniform sampling

$$S_k = \{1\} \text{ with prob } p_1$$

$$S_k = \{2\} \text{ with prob } p_2$$

$$S_k = \{3\} \text{ with prob } p_3$$

# Arbitrary Sampling: Examples for $n = 3$

GD

$$S_k = \{1, 2, 3\} \text{ with prob } 1$$

SAGA

$$S_k = \{1\} \text{ with prob } 1/3$$

$$S_k = \{2\} \text{ with prob } 1/3$$

$$S_k = \{3\} \text{ with prob } 1/3$$

SAGA with nonuniform sampling

$$S_k = \{1\} \text{ with prob } p_1$$

$$S_k = \{2\} \text{ with prob } p_2$$

$$S_k = \{3\} \text{ with prob } p_3$$

Minibatch SAGA (with 2-nice sampling)

$$S_k = \{1, 2\} \text{ with prob } 1/3$$

$$S_k = \{2, 3\} \text{ with prob } 1/3$$

$$S_k = \{3, 1\} \text{ with prob } 1/3$$

# Arbitrary Sampling: Examples for $n = 3$

GD

$$S_k = \{1, 2, 3\} \text{ with prob } 1$$

SAGA

$$S_k = \{1\} \text{ with prob } 1/3$$

$$S_k = \{2\} \text{ with prob } 1/3$$

$$S_k = \{3\} \text{ with prob } 1/3$$

SAGA with nonuniform sampling

$$S_k = \{1\} \text{ with prob } p_1$$

$$S_k = \{2\} \text{ with prob } p_2$$

$$S_k = \{3\} \text{ with prob } p_3$$

Minibatch SAGA (with 2-nice sampling)

$$S_k = \{1, 2\} \text{ with prob } 1/3$$

$$S_k = \{2, 3\} \text{ with prob } 1/3$$

$$S_k = \{3, 1\} \text{ with prob } 1/3$$

Interpolation between GD and SAGA

$$S_k = \{1, 2, 3\} \text{ with prob } 1/2$$

$$S_k = \{1\} \text{ with prob } 1/6$$

$$S_k = \{2\} \text{ with prob } 1/6$$

$$S_k = \{3\} \text{ with prob } 1/6$$

# **A Brief History of Arbitrary Sampling**

| #  | Paper   | Algorithm         | Comment  |
|----|---|-------------------|--|
| 1  | <b>R. &amp; Takáč (OL 2016; arXiv 2013)</b><br>On optimal probabilities in stochastic coordinate descent methods  | NSync             | <b>Arbitrary sampling (AS) first introduced</b><br>Analysis of coordinate descent under strong convexity   |
| 2  | <b>Qu, R. &amp; Zhang (NeurIPS 2015)</b><br>Quartz: Randomized dual coordinate ascent with arbitrary sampling   | QUARTZ            | <b>First AS SGD method for min <math>P</math></b><br>Primal-dual stochastic fixed point method; variance reduced   |
| 3  | <b>Csiba &amp; R. (arXiv 2015)</b><br>Primal method for ERM with flexible mini-batching schemes and non-convex losses   | Dual-free SDCA    | <b>First primal-only AS SGD method for min <math>P</math></b><br>Variance-reduced  |
| 4  | <b>Qu &amp; R. (OMS 2016)</b><br>Coordinate descent with arbitrary sampling I: algorithms and complexity  | ALPHA             | <b>First accelerated coordinate descent method with AS</b><br>Analysis for smooth convex functions   |
| 5  | <b>Qu &amp; R. (OMS 2016)</b><br>Coordinate descent with arbitrary sampling II: expected separable overapproximation  |                   | <b>First dedicated study of ESO inequalities needed for analysis of AS methods</b> $\mathbb{E}_S \left[ \left\  \sum_{i \in S} A_i h_i \right\ ^2 \right] \leq \sum_{i=1}^n p_i v_i \ h_i\ ^2$ |
| 6  | <b>Chambolle, Ehrhardt, R. &amp; Schoenlieb (SIOPT 2018)</b><br>Stochastic primal-dual hybrid gradient algorithm with arbitrary sampling and imaging applications | SPDHGM            | <b>Chambolle-Pock method with AS</b>   |
| 7  | <b>Hanzely, Mishchenko &amp; R. (NeurIPS 2018)</b><br>SEGA: Variance reduction via gradient sketching   | SEGA              | <b>Variance-reduce coordinate descent with AS</b>  |
| 8  | <b>Hanzely &amp; R. (AISTATS 2019)</b><br>Accelerated coordinate descent with arbitrary sampling and best rates for minibatches                                   | ACD               | <b>First accelerated coordinate descent method with AS</b><br>Analysis for smooth strongly convex functions<br>Importance sampling for minibatches   |
| 9  | <b>Horváth &amp; R. (ICML 2019)</b><br>Nonconvex variance reduced optimization with arbitrary sampling  | SARAH, SVRG, SAGA | <b>First non-convex analysis of an AS method</b><br>First optimal mini-batch sampling  |
| 10 | <b>Gower, Loizou, Qian, Sailanbayev, Shulgin &amp; R. (ICML 2019)</b><br>SGD: general analysis and improved rates   | SGD-AS            | <b>First AS variant of SGD (without variance reduction)</b><br>Optimal minibatch size  |
| 11 | <b>Qian, Qu &amp; R. (ICML 2019)</b><br>SAGA with arbitrary sampling  | <b>SAGA-AS</b>    | <b>First AS variant of SAGA</b>  |

# **The Algorithm**



# New Method: SAGA-AS (high level)

## The Problem

$$\min_{x \in \mathbb{R}^d} P(x) \stackrel{\text{def}}{=} \left( \sum_{i=1}^n \lambda_i f_i(x) \right) + \psi(x)$$

# New Method: SAGA-AS (high level)

Arbitrary Sampling

1

Sample fresh  $S_k \subseteq \{1, 2, \dots, n\}$

The Problem

$$\min_{x \in \mathbb{R}^d} P(x) \stackrel{\text{def}}{=} \left( \sum_{i=1}^n \lambda_i f_i(x) \right) + \psi(x)$$

# New Method: SAGA-AS (high level)

## The Problem

$$\min_{x \in \mathbb{R}^d} P(x) \stackrel{\text{def}}{=} \left( \sum_{i=1}^n \lambda_i f_i(x) \right) + \psi(x)$$

Arbitrary Sampling

1 Sample fresh  $S_k \subseteq \{1, 2, \dots, n\}$

$$2 \quad \mathbf{J}_{:i}^{k+1} = \begin{cases} \nabla f_i(x^k) & i \in S_k \\ \mathbf{J}_{:i}^k & i \notin S_k \end{cases}$$

Jacobian Sketch, i.e., a random matrix approximating the Jacobian:

$$\mathbf{J}^{k+1} \approx \mathbf{G}(x^k) \stackrel{\text{def}}{=} [\nabla f_1(x^k), \dots, \nabla f_n(x^k)] \in \mathbb{R}^{d \times n}$$

# New Method: SAGA-AS (high level)

## The Problem

$$\min_{x \in \mathbb{R}^d} P(x) \stackrel{\text{def}}{=} \left( \sum_{i=1}^n \lambda_i f_i(x) \right) + \psi(x)$$

Arbitrary Sampling

1 Sample fresh  $S_k \subseteq \{1, 2, \dots, n\}$

$$2 \quad \mathbf{J}_{:i}^{k+1} = \begin{cases} \nabla f_i(x^k) & i \in S_k \\ \mathbf{J}_{:i}^k & i \notin S_k \end{cases}$$

Jacobian Sketch, i.e., a random matrix approximating the Jacobian:

$$\mathbf{J}^{k+1} \approx \mathbf{G}(x^k) \stackrel{\text{def}}{=} [\nabla f_1(x^k), \dots, \nabla f_n(x^k)] \in \mathbb{R}^{d \times n}$$

3 Use  $\mathbf{J}^{k+1}$ ,  $\mathbf{J}^k$  to build an unbiased estimator  $g^k$  of  $\nabla f(x^k)$

# New Method: SAGA-AS (high level)

## The Problem

$$\min_{x \in \mathbb{R}^d} P(x) \stackrel{\text{def}}{=} \left( \sum_{i=1}^n \lambda_i f_i(x) \right) + \psi(x)$$

Arbitrary Sampling

1 Sample fresh  $S_k \subseteq \{1, 2, \dots, n\}$

$$2 \quad \mathbf{J}_{:i}^{k+1} = \begin{cases} \nabla f_i(x^k) & i \in S_k \\ \mathbf{J}_{:i}^k & i \notin S_k \end{cases}$$

Jacobian Sketch, i.e., a random matrix approximating the Jacobian:

$$\mathbf{J}^{k+1} \approx \mathbf{G}(x^k) \stackrel{\text{def}}{=} [\nabla f_1(x^k), \dots, \nabla f_n(x^k)] \in \mathbb{R}^{d \times n}$$

3 Use  $\mathbf{J}^{k+1}$ ,  $\mathbf{J}^k$  to build an unbiased estimator  $g^k$  of  $\nabla f(x^k)$

$$4 \quad x^{k+1} = \text{prox}_{\alpha\psi} (x^k - \alpha g^k)$$

Proximal SGD step with fixed step size

$$\text{prox}_{\psi}(x) \stackrel{\text{def}}{=} \arg \min_y \left\{ \frac{1}{2} \|y - x\|^2 + \psi(y) \right\}$$

# Convergence Theory

# Convergence Theory

$$\mathbb{E}_S \left[ \|\mathbf{M} \text{Diag}(\theta_S) \mathbf{I}_S \lambda\|^2 \right] \leq \sum_{i=1}^n \mathcal{A}_i \lambda_i^2 \|\mathbf{M}_{:i}\|^2 + \mathcal{B} \|\mathbf{M} \lambda\|^2$$

**Lyapunov function:**

$$\Psi^k \stackrel{\text{def}}{=} \|x^k - x^*\|^2 + 2\alpha \sum_{i=1}^n \sigma_i \mathcal{A}_i \lambda_i^2 \|\mathbf{J}_{:i}^k - \nabla f_i(x^*)\|^2$$

| Regime   | Arbitrary sampling  | Thm |
|--|---|-----|
| <b>Smooth</b><br>$\psi \equiv 0$<br>$f_i$ is $L_i$ -smooth, $f$ is $\mu$ -strongly convex  | $\max \left\{ \max_{1 \leq i \leq n} \left\{ \frac{1}{p_i} + \frac{4(1+\mathcal{B})L_i \mathcal{A}_i \lambda_i}{\mu} \right\}, \frac{2\mathcal{B}(1+1/\mathcal{B})L}{\mu} \right\} \log\left(\frac{1}{\epsilon}\right)$ | 3.3 |
| <b>Nonsmooth</b><br>$P$ satisfies $\mu$ -growth condition (19) and Assumption 4.3<br>$f_i(x) = \phi_i(\mathbf{A}_i^\top x)$ , $\phi_i$ is $1/\gamma$ -smooth, $f$ is $L$ -smooth | $\left( 2 + \max \left\{ \frac{6L}{\mu}, 3 \max_{1 \leq i \leq n} \left\{ \frac{1}{p_i} + \frac{4v_i \lambda_i}{p_i \mu \gamma} \right\} \right\} \right) \log\left(\frac{1}{\epsilon}\right)$                          | 4.4 |
| <b>Nonsmooth</b><br>$\psi$ is $\mu$ -strongly convex<br>$f_i(x) = \phi_i(\mathbf{A}_i^\top x)$ , $\phi_i$ is $1/\gamma$ -smooth  | $\max_{1 \leq i \leq n} \left\{ 1 + \frac{1}{p_i} + \frac{3v_i \lambda_i}{p_i \mu \gamma} \right\} \log\left(\frac{1}{\epsilon}\right)$   | 4.5 |

Table 1. Iteration complexity results for SAGA-AS. We have  $p_i := \mathbb{P}(i \in S)$ , where  $S$  is a sampling of subsets of  $[n]$  utilized by SAGA-AS. The key complexity parameters  $\mathcal{A}_i$ ,  $\mathcal{B}$ , and  $v_i$  are defined in the sections containing the theorems.

**Expected Separable Over-approximation (ESO):**

$$\mathbb{E}_S \left[ \left\| \sum_{i \in S} \mathbf{A}_i h_i \right\|^2 \right] \leq \sum_{i=1}^n p_i v_i \|h_i\|^2 \quad p_i \stackrel{\text{def}}{=} \text{Prob}(i \in S_k)$$

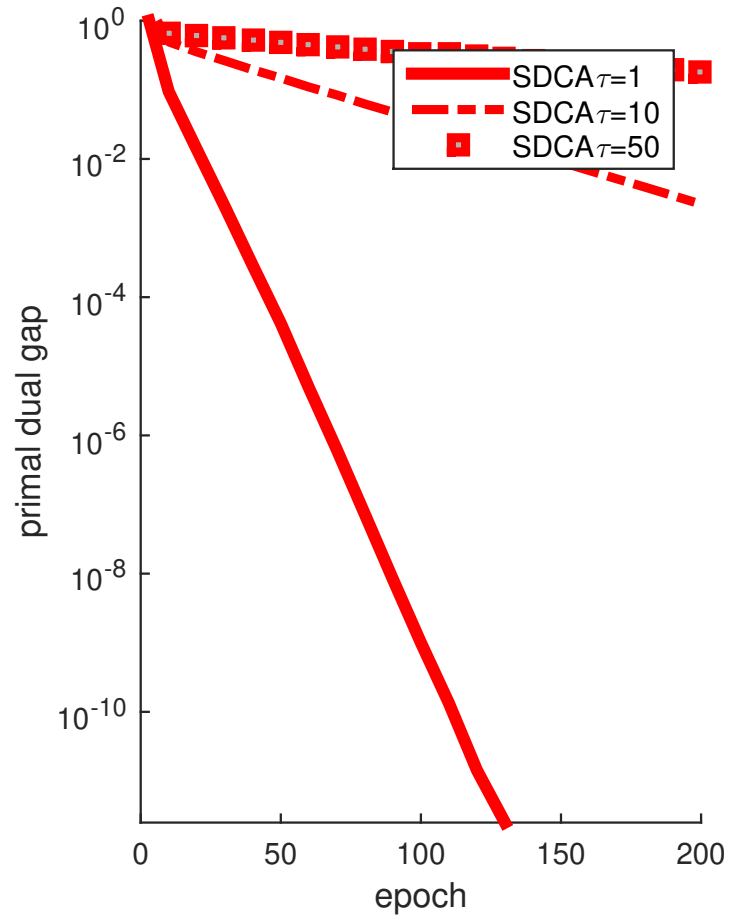
# Contributions



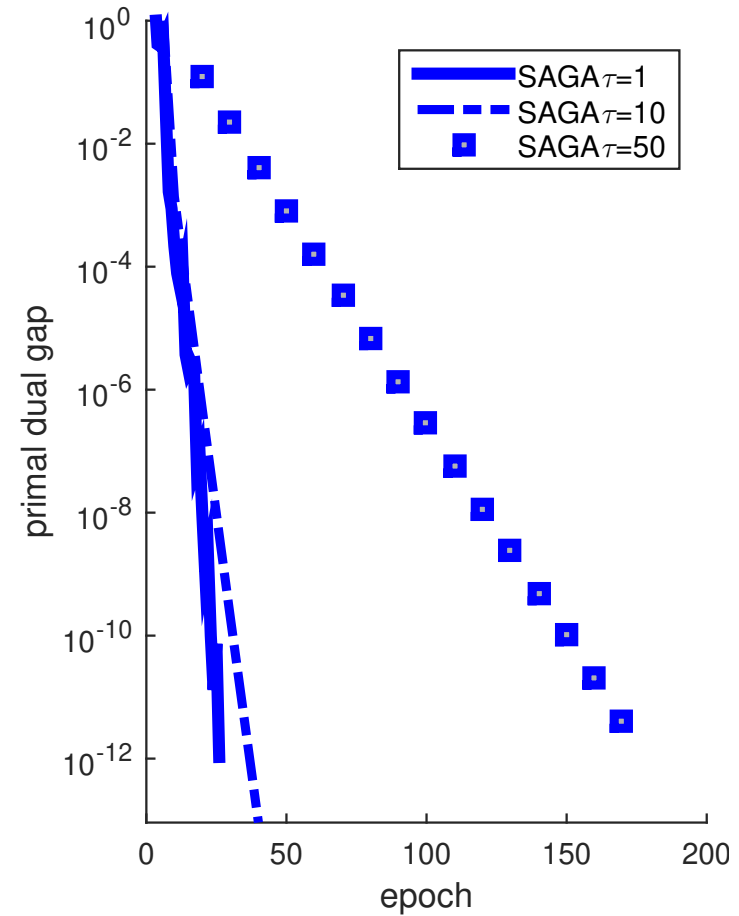
|                         | SAGA<br>(Defazio et al 2014)                 | QUARTZ<br>(Qu et al 2015)                               | JacSketch<br>(Gower et al 2018)  | SAGA-AS<br>(THIS WORK)                |
|-------------------------|--|---|--|---------------------------------------|
| PRIMAL / DUAL           | Primal                                       | Primal-dual   | Primal   | Primal                                |
| SAMPLING                | Uniform sampling of<br>of single data points | Arbitrary sampling<br>(first AS method for<br>min $P$ ) | A general sketching<br>mechanism, but does not<br>cover arbitrary sampling | Arbitrary sampling                    |
| IMPORTANCE<br>SAMPLING? | NO   | YES   | YES<br>(first SAGA-IS, but not for<br>minibatches)                         | YES<br>(also for minibatches)         |
| REGULARIZER             | Support for any<br>convex regularizer        | Support for strongly<br>convex regularizer              | No support for a<br>regularizer  | Support for any convex<br>regularizer |
| RATE                    | Linear                                       | Linear  | Linear   | Linear<br>(same or better)            |
| ASSUMPTIONS             | Each $f_i$ strongly<br>convex                | strongly convex<br>regularizer                          | Each $f_i$ strongly convex   | $P$ satisfying quadratic<br>growth    |
| HANDLING BIAS           | Scaling                                      | Built in  | Bias-correcting random<br>variable   | Bias-correcting random<br>vector      |

# Experiments

# SDCA vs SAGA

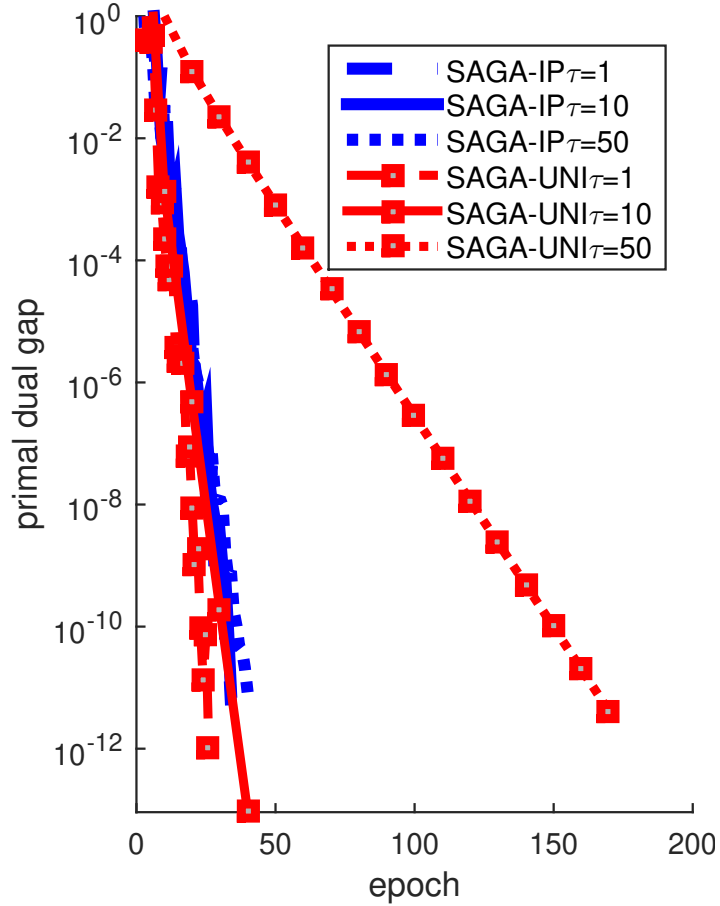


ijcnn1

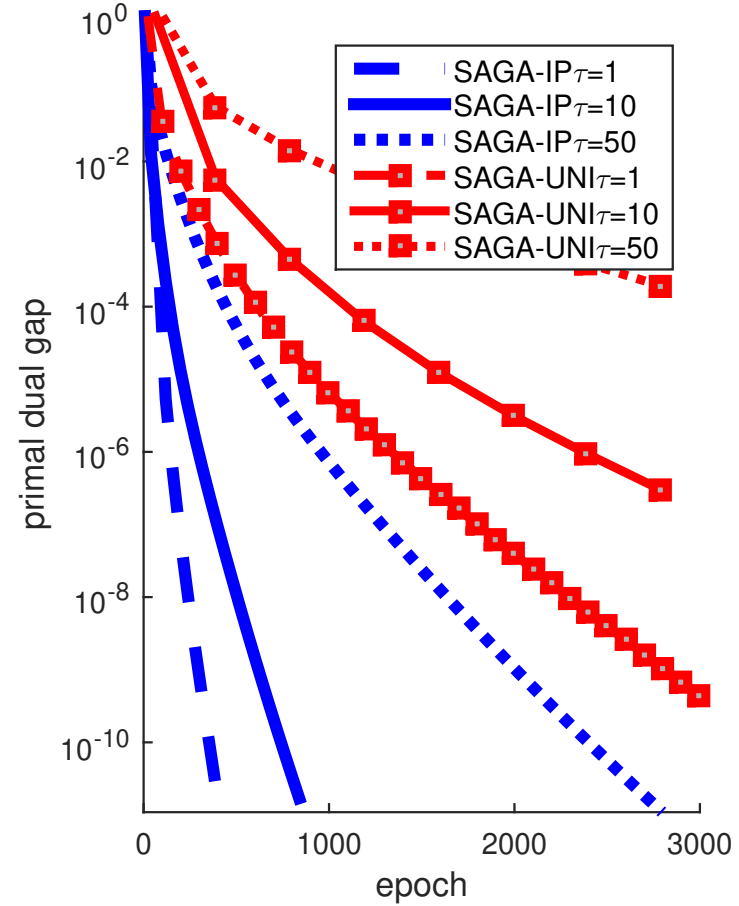


ijcnn1

# Uniform vs Importance Sampling



ijcnn1



w8a

**What's Next?**

### The Problem

$$\min_{x \in \mathbb{R}^d} P(x) \stackrel{\text{def}}{=} \left( \sum_{i=1}^n \lambda_i f_i(x) \right) + \psi(x), \quad (1)$$

where  $f \stackrel{\text{def}}{=} \sum_{i=1}^n \lambda_i f_i(x)$ ,  $f_i$  are smooth and convex,  $\lambda_i > 0$  are weights, and  $\psi: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  is closed and convex.

### Sampling

**Sampling:** A random set valued mapping  $S$  with values being subsets of  $\{1, \dots, n\}$ . A sampling is uniquely defined by assigning probabilities to all  $2^n$  subsets of  $\{1, \dots, n\}$ . Let  $\tau \stackrel{\text{def}}{=} \mathbb{E}[|S|]$  be the expected size of  $S$ , and define

$$p_i \stackrel{\text{def}}{=} \text{Prob}(i \in S), \quad i \in \{1, \dots, n\}.$$

A sampling is called **proper** if  $p_i > 0$  for all  $i$ . For  $C \subseteq \{1, \dots, n\}$ , let

$$p_C \stackrel{\text{def}}{=} \text{Prob}(S = C).$$

**Bias-correcting random vector:** vector  $\theta_S = (\theta_S^1, \dots, \theta_S^n) \in \mathbb{R}^n$  with the property

$$\mathbb{E}[\text{Diag}(\theta_S) \mathbf{I}_S e] = e, \quad \text{i.e.,} \quad \mathbb{E}[\theta_S^i \mathbf{1}_{i \in S}] = 1, \quad \forall i, \quad (2)$$

where

- $e$ :  $n \times 1$  vector of all ones
- $\mathbf{I}$ :  $n \times n$  identity matrix
- $\mathbf{I}_S$ :  $n \times n$  matrix with ones in places  $(i, i)$  for  $i \in S$
- $\mathbf{1}_{i \in S}$ : indicator random variable of the event  $i \in S$ , i.e.,  $\mathbf{1}_{i \in S} = 1$  if  $i \in S$  and  $\mathbf{1}_{i \in S} = 0$  if  $i \notin S$

### Algorithm

**Prox operator:**  $\text{prox}_\psi^{\alpha}(x) \stackrel{\text{def}}{=} \arg \min \left\{ \frac{1}{2\alpha} \|x - y\|^2 + \psi(y) \right\}$

**Gradient matrix:**  $\mathbf{G}(x) \stackrel{\text{def}}{=} [\nabla f_1(x), \dots, \nabla f_n(x)]$

**Algorithm 1: SAGA with Arbitrary Sampling**

*Initialize:*  $x^0 \in \mathbb{R}^d$ ,  $\mathbf{J}^0 \in \mathbb{R}^{n \times n}$

*Parameters:* arbitrary  $\theta_S$ , stepsize  $\alpha > 0$

**for**  $k = 1, 2, \dots, \mathbf{c}$

Sample fresh  $S_k \subseteq \{1, \dots, n\}$

$\mathbf{J}^{k+1} = \mathbf{J}^k + (\mathbf{G}(x^k) - \mathbf{J}^k) \mathbf{I}_{S_k}$

$g^k = \mathbf{J}^k \lambda + (\mathbf{G}(x^k) - \mathbf{J}^k) \lambda$

$x^{k+1} = \text{prox}_\psi^{\alpha}(x^k - \alpha g^k)$

**end**

### Smooth

**Assumptions:**

- $f_i$  is convex and  $L_i$ -smooth
- $f$  is  $\mu$ -strongly convex and  $L$ -smooth
- There exist constants  $\mathcal{A}_i \geq 0$  and  $0 \leq \mathcal{B} \leq 1$  such that for any matrix  $\mathbf{M} \in \mathbb{R}^{d \times n}$

$$\mathbb{E}[\| \text{MDiag}(\theta_S) \mathbf{I}_S \mathbf{M} \|^2] \leq \sum_{i=1}^n \mathcal{A}_i \lambda_i^2 \|\mathbf{M}_i\|^2 + \mathcal{B} \|\mathbf{M}\|^2$$

**Lyapunov function:**

$$\Psi^k \stackrel{\text{def}}{=} \|x^k - x^*\|^2 + 2\alpha \sum_{i=1}^n \sigma_i \lambda_i^2 \|\mathbf{J}_i^k - \nabla f_i(x^*)\|^2,$$

where  $\sigma_i = \frac{1}{4(1+\mathcal{B})L_i \mathcal{A}_i \lambda_i}$  and  $x^*$  is a solution of (1).

### Convergence Result ( $\mathbb{E}[\Psi^k] \leq \epsilon \cdot \mathbb{E}[\Psi^0]$ )

**$\mu$  is known:**  $\alpha = \min_i \left\{ \frac{\mu}{\mu + 4(1+\mathcal{B})L_i \mathcal{A}_i \lambda_i}, \frac{\mathcal{B}^{-1}}{2(1+\mathcal{B})L} \right\}$

$$k \geq \max_i \left\{ \frac{1}{p_i} + \frac{4(1+\mathcal{B})L_i \mathcal{A}_i \lambda_i}{\mu}, \frac{2\mathcal{B}(1+\frac{1}{\mathcal{B}})L}{\mu} \right\} \log \left( \frac{1}{\epsilon} \right).$$

**$\mu$  is unknown:**  $\alpha = \min_i \left\{ \frac{\mu}{8(1+\mathcal{B})L_i \mathcal{A}_i \lambda_i}, \frac{\mathcal{B}^{-1}}{2(1+\mathcal{B})L} \right\}$

$$k \geq \max_i \left\{ \frac{2}{p_i}, \frac{8(1+\mathcal{B})L_i \mathcal{A}_i \lambda_i}{\mu}, \frac{2\mathcal{B}(1+\frac{1}{\mathcal{B}})L}{\mu} \right\} \log \left( \frac{1}{\epsilon} \right).$$

### Interface For Sampling

- Proper sampling:  $\mathcal{A}_i = \beta_i \stackrel{\text{def}}{=} \sum_{C \subseteq [n]: i \in C} p_C |\mathcal{C}| (\theta_C^i)^2$ ,  $\mathcal{B} = 0$ .
- $\tau$ -nice sampling ( $\theta_S^i = \frac{1}{p_i}$ ):  $\mathcal{A}_i = \frac{\tau}{p_i} - \frac{\tau-1}{n-p_i}$ ,  $\mathcal{B} = \frac{n(\tau-1)}{\tau(n-1)}$ .
- Independent sampling ( $\theta_S^i = \frac{1}{p_i}$ ):  $\mathcal{A}_i = \frac{1}{p_i} - 1$ ,  $\mathcal{B} = 1$ .

### Optimal Bias-Correcting Random Vector

Let  $\Theta(S)$  be the collection of all bias-correcting random vectors associated with sampling  $S$ , i.e.,  $\mathbb{E}[\theta_S \mathbf{I}_S e] = e$ . Let  $\mathbb{F}^{\Theta(S)}$

### Nonsmooth Case (strongly convex)

**Assumptions:**

- $f_i(x) = \phi_i(\mathbf{A}_i^\top x)$
- $\phi$  is  $1/\gamma$ -smooth and convex
- $\psi_i$  is  $\mu$ -strongly convex
- Choose  $\theta_S^i = 1/p_i$
- Let  $v_i$  satisfy the **ESO inequality**:

$$\mathbb{E}_S \left[ \left\| \sum_{i \in S} \mathbf{A}_i / h_i \right\|^2 \right] \leq \sum_{i=1}^n p_i v_i \|h_i\|^2.$$

**Lyapunov function:**

$$\Psi^k \stackrel{\text{def}}{=} \|x^k - x^*\|^2 + \alpha \sum_{i=1}^n \sigma_i \lambda_i^2 \|\mathbf{J}_i^k - \nabla \phi_i(\mathbf{A}_i^\top x^*)\|^2,$$

### Strongly convex

and convex

$\sigma_i = \frac{1}{2v_i \lambda_i}$

**ESO inequality**

**Nullspace consistency:** For any  $x^*, y^* \in \mathcal{X}^*$  we have

$$\mathbf{A}_i^\top x^* = \mathbf{A}_i^\top y^*, \quad \forall i \in [n],$$

where  $\mathcal{X}^* \stackrel{\text{def}}{=} \arg \min \{P(x) : x \in \mathbb{R}^d\}$ .

**Quadratic functional growth condition:** there is a constant  $\mu > 0$  such that

$$P(x^k) - P^* \geq \frac{\mu}{2} \|x^k - [x^k]^*\|^2, \quad w.p.1, \quad \forall k \geq 1,$$

where  $[x]^* = \arg \min \{\|x - y\| : y \in \mathcal{X}^*\}$ , for the sequence  $\{x^k\}$  produced by the Algorithm.

**Lyapunov function:**

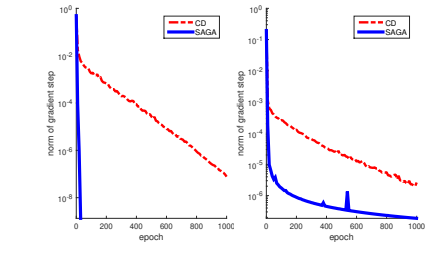
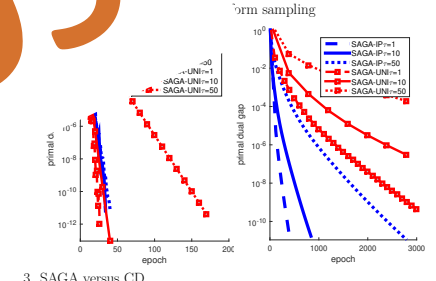
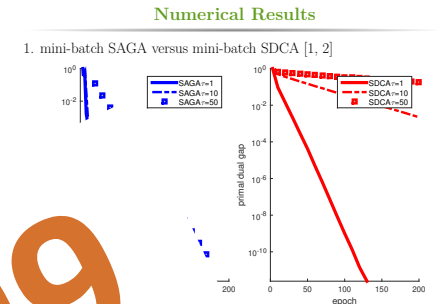
$$\Psi^k \stackrel{\text{def}}{=} \|x^k - [x^k]^*\|^2 + \alpha \sum_{i=1}^n \sigma_i \lambda_i^2 \|\alpha_i^k - \nabla \phi_i(\mathbf{A}_i^\top x^*)\|^2,$$

### Convergence Result ( $\mathbb{E}[\Psi^k] \leq \epsilon \cdot \mathbb{E}[\Psi^0]$ )

**$\mu$  is known:**  $\alpha = \min \left\{ \frac{1}{3} \min_{i \leq i \leq n} \frac{\mu}{\mu + 4(1+\mathcal{B})L_i \gamma}, \frac{1}{3L} \right\}$

$$k \geq \left( 2 + \max \left\{ \frac{6L}{\mu}, 3 \max_i \left( \frac{1}{p_i} + \frac{4v_i \lambda_i}{\mu \gamma} \right) \right\} \right) \log \left( \frac{1}{\epsilon} \right).$$

**$\mu$  is unknown:**  $\alpha = \min \left\{ \min_{i \leq i \leq n} \frac{\mu}{12v_i \lambda_i \gamma}, \frac{1}{3L} \right\}$

$$k \geq \left( 2 + \max \left\{ \frac{6L}{\mu}, \max_i \left\{ \frac{24v_i \lambda_i}{\mu p_i \gamma}, \frac{2}{p_i} \right\} \right\} \right) \log \left( \frac{1}{\epsilon} \right).$$


### References

- Shai Shalev-Shwartz and Tong Zhang. Stochastic dual coordinate ascent methods for regularized loss. *Journal of Machine Learning Research*, 14(1):567–599, 2013.
- Zheng Qu, Peter Richtárik, and Tong Zhang. Quartz: Randomized dual coordinate ascent with arbitrary sampling. In *Advances in Neural Information Processing Systems 28*, pages 865–873. Curran Associates, Inc., 2015.
- R. M. Gower, P. Richtárik, and F. Bach. Stochastic quasi-gradient methods: Variance reduction via Jacobian sketching. *arXiv Preprint arXiv: 1805.02632*, 2018.
- Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems 27*, 2014.



**The End**