



SDNA

Stochastic Dual Newton Ascent for Empirical Risk Minimization



Peter Richtárik

Coauthors



Zheng Qu
(Edinburgh)



Martin Takáč
(Lehigh)



Olivier Fercoq
(Telecom ParisTech)



Zheng Qu, P.R., Martin Takáč and Olivier Fercoq
SDNA: Stochastic Dual Newton Ascent for empirical risk minimization
arXiv:1502.02268, 2015

Randomized Methods with Arbitrary Sampling



P.R. and Martin Takáč

On optimal probabilities in stochastic coordinate descent methods

In NIPS Workshop on Optimization for Machine Learning, 2013 (arXiv:1310.3438)



Zheng Qu, P.R. and Tong Zhang

Randomized dual coordinate ascent with arbitrary sampling

arXiv:1411.5873, 2014

I will
mention this



Zheng Qu and P.R.

Coordinate descent with arbitrary sampling I: algorithms and complexity

arXiv:1412.8060, 2014

Zheng Qu:
tomorrow

Zheng Qu and P.R.

Coordinate descent with arbitrary sampling II: ESO

arXiv:1412.8063, 2014



Zheng Qu, P.R., Martin Takáč and Olivier Fercoq

SDNA: Stochastic Dual Newton Ascent for empirical risk minimization

arXiv:1502.02268, 2015

Robert M. Gower and P.R.

Randomized iterative methods for linear systems

Robert M. Gower: today

Part A

Minimization of a Smooth & Strongly Convex Function

The Problem & Assumptions

$$\min_{x \in \mathbb{R}^n} f(x)$$

Strong convexity

Large dimension

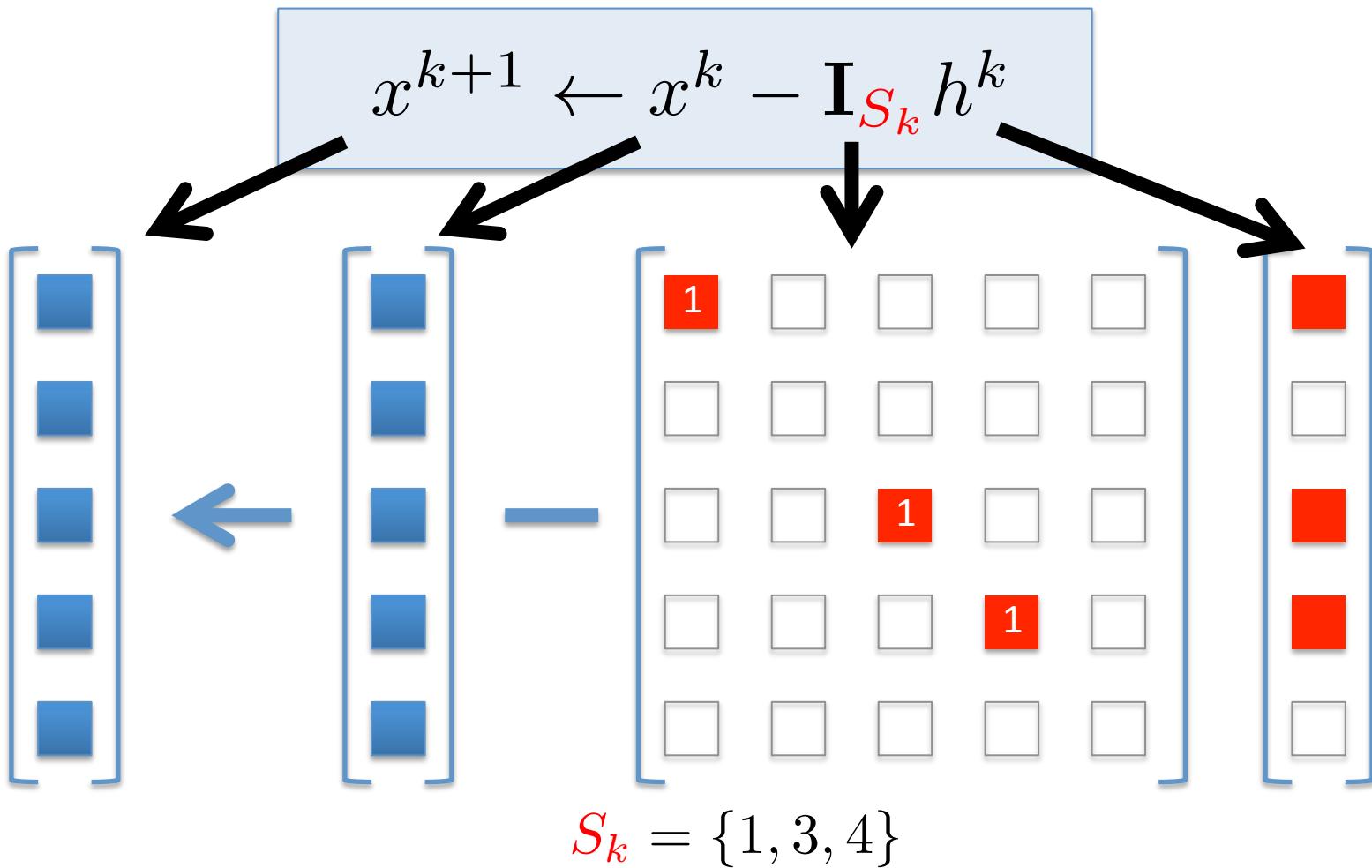
$$f(x) + (\nabla f(x))^\top h + \frac{1}{2} h^\top \mathbf{G} h \leq f(x + h)$$

Smoothness

Positive definite matrices

$$f(x + h) \leq f(x) + (\nabla f(x))^\top h + \frac{1}{2} h^\top \mathbf{M} h$$

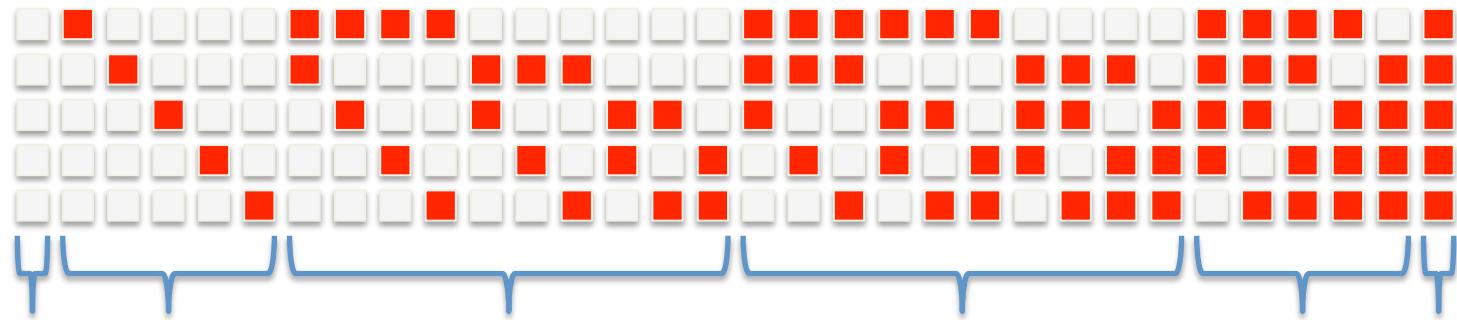
Randomized Update



Arbitrary Sampling

$$\mathbb{P}(S_k = S) = ? \quad \forall S \subseteq \{1, 2, \dots, n\}$$

Example: $n = 5$



$$|S| = 0 \quad |S| = 1 \quad |S| = 2 \quad |S| = 3 \quad |S| = 4 \quad |S| = 5$$

Method 3



P.R. and Martin Takáč

On optimal probabilities in stochastic coordinate descent methods

In NIPS Workshop on Optimization for Machine Learning, 2013

(arXiv:1310.3438)

Key Inequality

$$f(x + h) \leq f(x) + (\nabla f(x))^\top h + \frac{1}{2} h^\top \mathbf{M} h$$

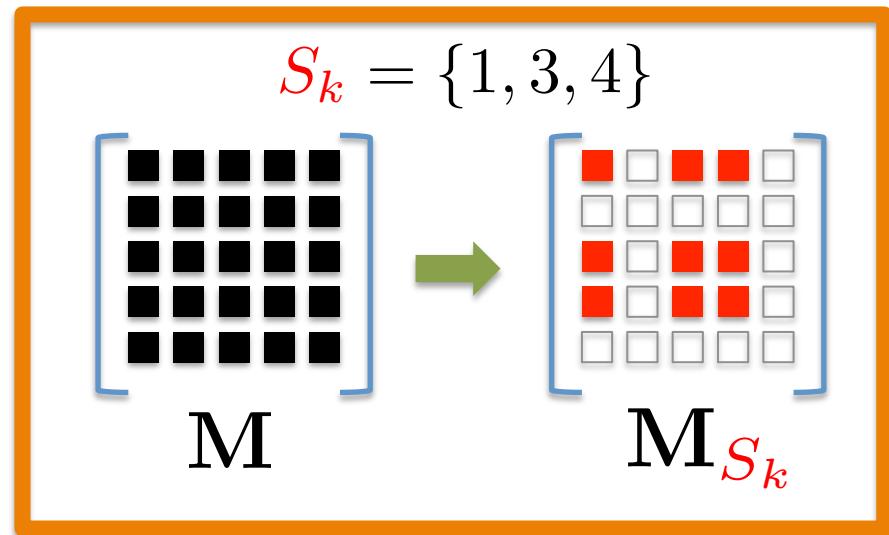


$$x \leftarrow x^k$$

$$h \leftarrow \mathbf{I}_{S_k} h = \sum_{i \in S_k} h_i e_i$$



$$f(x^k + \mathbf{I}_{S_k} h) \leq f(x^k) + (\nabla f(x^k))^\top (\mathbf{I}_{S_k} h) + \frac{1}{2} (\mathbf{I}_{S_k} h)^\top \mathbf{M} (\mathbf{I}_{S_k} h)$$



$$h^\top \mathbf{M}_{S_k} h$$

$$\frac{1}{2} (\mathbf{I}_{S_k} h)^\top \mathbf{M} (\mathbf{I}_{S_k} h)$$



Method 3

$$f(x^k + \mathbf{I}_{S_k} h) \leq f(x^k) + (\mathbf{I}_{S_k} \nabla f(x^k))^{\top} h + \frac{1}{2} h^{\top} \mathbf{M}_{S_k} h$$

1. take expectations on both sides



$$p_i = \mathbb{P}(i \in S_k)$$

$$\mathbb{E}[f(x^k + \mathbf{I}_{S_k} h)] \leq f(x^k) + (\text{Diag}(p) \nabla f(x^k))^{\top} h + \frac{1}{2} h^{\top} \mathbb{E}[\mathbf{M}_{S_k}] h$$

2. diagonalize



$$\mathbb{E}[\mathbf{M}_{S_k}] \preceq \text{Diag}(p \circ v)$$

$$\mathbb{E}[f(x^k + \mathbf{I}_{S_k} h)] \leq f(x^k) + (\text{Diag}(p) \nabla f(x^k))^{\top} h + \frac{1}{2} h^{\top} \text{Diag}(p \circ v) h$$

3. minimize the RHS in h



$$x^{k+1} \leftarrow x^k - \mathbf{I}_{S_k} (\text{Diag}(v))^{-1} \nabla f(x^k)$$

Method 3

i.i.d. with arbitrary distribution

Choose a random set S_k of coordinates

For $i \in S_k$ do

$$x_i^{k+1} \leftarrow x_i^k - \frac{1}{v_i} (\nabla f(x^k))^{\top} e_i$$

For $i \notin S_k$ do

$$x_i^{k+1} \leftarrow x_i^k$$

Convergence

Theorem (RT'13)

$$\mathbb{E}[f(x^k) - f(x^*)] \leq (1 - \sigma_3)^k (f(x^0) - f(x^*))$$



$$\sigma_3 = \lambda_{\min} \left(\mathbf{G}^{1/2} \mathbf{Diag}(p \circ v^{-1}) \mathbf{G}^{1/2} \right)$$

Alternative formulation:

$$k \geq \frac{1}{\sigma_3} \log \left(\frac{f(x^0) - f(x^*)}{\epsilon} \right) \Rightarrow \mathbb{E}[f(x^k) - f(x^*)] \leq \epsilon$$

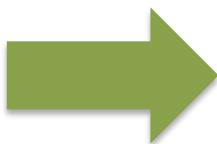
Uniform vs Optimal Sampling

Special case:

$$\mathbf{G} = \lambda \mathbf{I} \quad \Rightarrow \quad \frac{1}{\sigma_3} = \max_i \frac{v_i}{\lambda p_i}$$

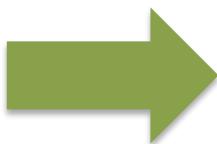
$$\mathbb{P}(|S_k| = 1) = 1 \quad \Rightarrow \quad v_i = \mathbf{M}_{ii}$$

$$p_i = \frac{1}{n}$$



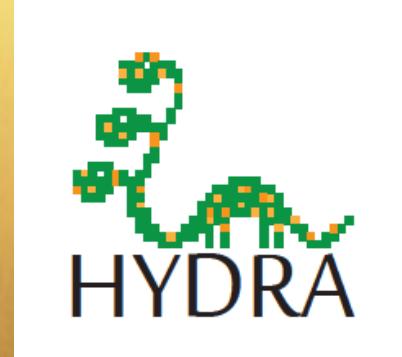
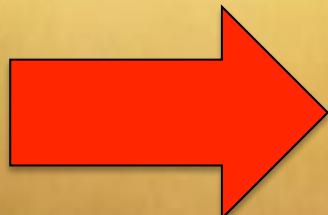
$$\frac{1}{\sigma_3} = \frac{n \max_i \mathbf{M}_{ii}}{\lambda}$$

$$p_i = \frac{\mathbf{M}_{ii}}{\sum_i \mathbf{M}_{ii}}$$



$$\frac{1}{\sigma_3} = \frac{\sum_{i=1}^n \mathbf{M}_{ii}}{\lambda}$$

Big Data Detour

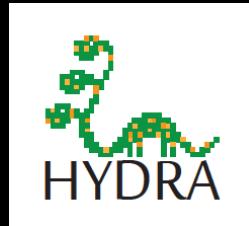


Experiment

Machine: 128 nodes of Hector Supercomputer (4096 cores)

Problem: LASSO, $n = 1$ billion, $d = 0.5$ billion, 3 TB

Algorithm:



with $c = 512$

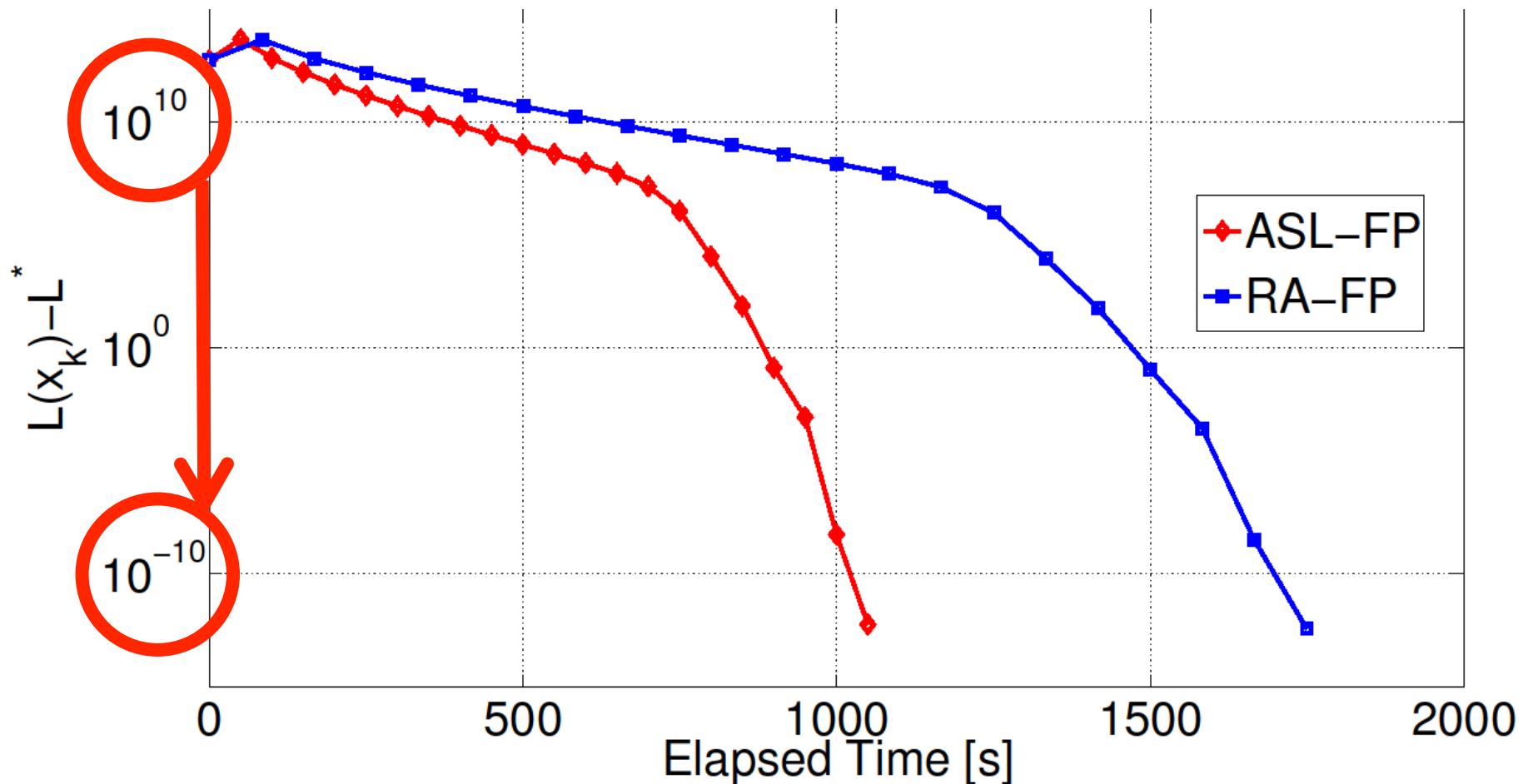


P.R. and Martin Takáč

Distributed coordinate descent method for learning with big data

arXiv:1310.2059, 2013

LASSO: 3TB data + 128 nodes



Method 2

Method 2

$$f(x^k + \mathbf{I}_{S_k} h) \leq f(x^k) + (\mathbf{I}_{S_k} \nabla f(x^k))^{\top} h + \frac{1}{2} h^{\top} \mathbf{M}_{S_k} h$$



1. take expectations on both sides

$$\mathbb{E}[f(x^k + \mathbf{I}_{S_k} h)] \leq f(x^k) + (\mathbf{Diag}(p) \nabla f(x^k))^{\top} h + \frac{1}{2} h^{\top} \mathbb{E}[\mathbf{M}_{S_k}] h$$



$$p_i = \mathbb{P}(i \in S_k)$$

2. minimize the RHS in h

$$x^{k+1} \leftarrow x^k - \mathbf{I}_{S_k} (\mathbb{E}[\mathbf{M}_{S_k}])^{-1} \mathbf{Diag}(p) \nabla f(x^k)$$

Convergence of Method 2

Theorem (QRTF'15)

$$\mathbb{E}[f(x^k) - f(x^*)] \leq (1 - \sigma_2)^k (f(x^0) - f(x^*))$$



$$\sigma_2 = \lambda_{\min} \left(\mathbf{G}^{1/2} \mathbf{Diag}(p) (\mathbb{E} [\mathbf{M}_{S_k}])^{-1} \mathbf{Diag}(p) \mathbf{G}^{1/2} \right)$$

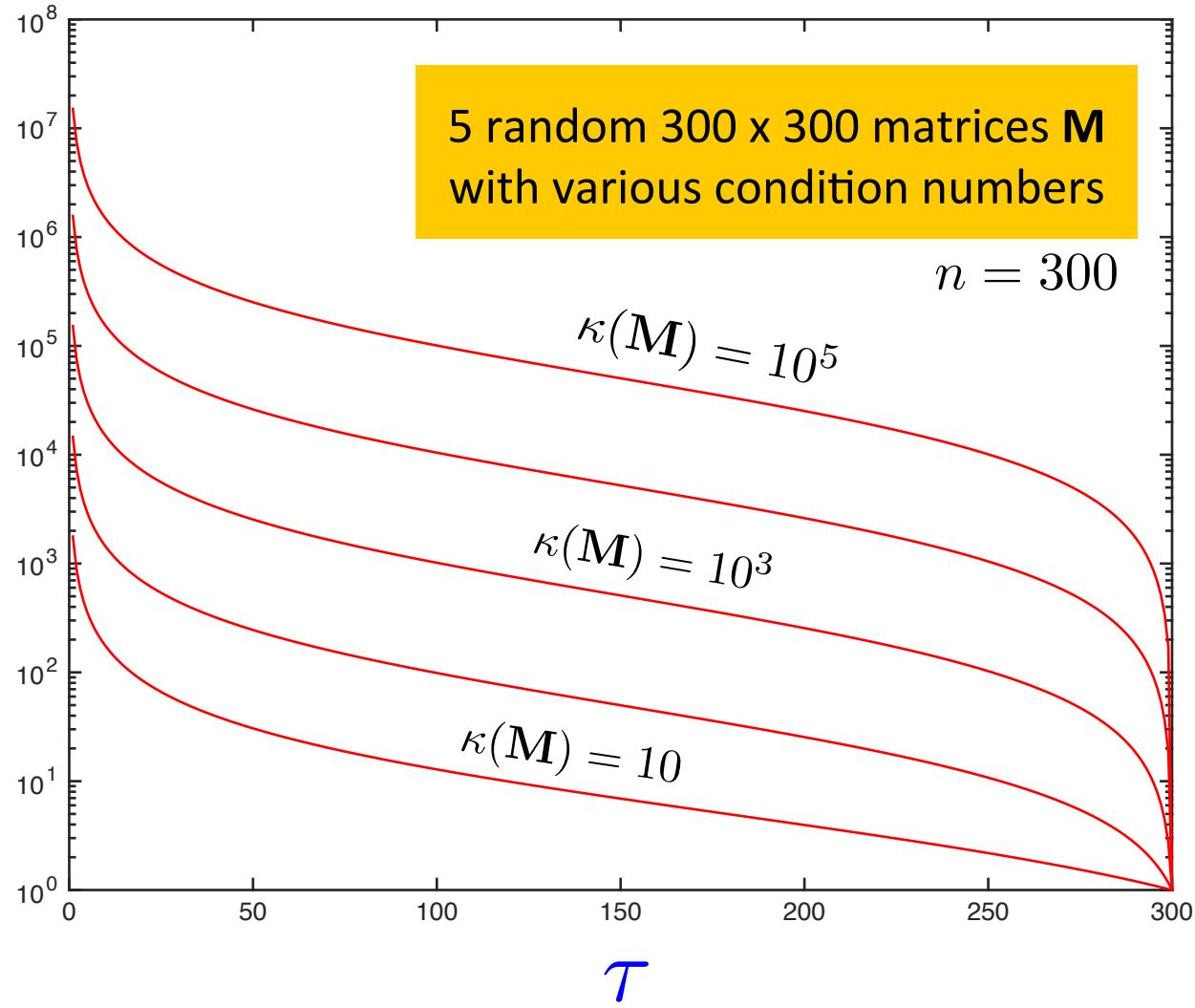
Alternative formulation:

$$k \geq \frac{1}{\sigma_2} \log \left(\frac{f(x^0) - f(x^*)}{\epsilon} \right) \Rightarrow \mathbb{E}[f(x^k) - f(x^*)] \leq \epsilon$$

Leading term in the complexity of Method 2 as a function of $\tau = \mathbb{E}[|S_k|]$

$$\frac{1}{\sigma_2(\tau)} = \frac{n}{n-1} \lambda_{\max} \left(G^{-1/2} \left[\left(\frac{n}{\tau} - 1 \right) \text{Diag}(M) + \left(1 - \frac{1}{\tau} \right) M \right] G^{-1/2} \right)$$

$$\frac{1}{\sigma_2(\tau)}$$



Method 1

Randomized Newton

Method

Method 1: Randomized Newton

$$f(x^k + \mathbf{I}_{S_k} h) \leq f(x^k) + (\mathbf{I}_{S_k} \nabla f(x^k))^{\top} h + \frac{1}{2} h^{\top} \mathbf{M}_{S_k} h$$



minimize the RHS in h

$$x^{k+1} \leftarrow x^k - (\mathbf{M}_{S_k})^{-1} \nabla f(x^k)$$

$$S_k = \{1, 3, 4\}$$

$$\mathbf{M}_{S_k} = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$$

$$(\mathbf{M}_{S_k})^{-1} = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$$

$$\mathbf{I}_{S_k} = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 1 \end{bmatrix}$$

$$\mathbf{M}_{S_k}$$

$$(\mathbf{M}_{S_k})^{-1}$$

$$\mathbf{I}_{S_k}$$

Convergence of Method 1 (Randomized Newton Method)

Theorem (QRTF'15)

$$\mathbb{E}[f(x^k) - f(x^*)] \leq (1 - \sigma_1)^k (f(x^0) - f(x^*))$$



$$\sigma_1 = \lambda_{\min} \left(\mathbf{G}^{1/2} \mathbb{E} \left[(\mathbf{M}_{S_k})^{-1} \right] \mathbf{G}^{1/2} \right)$$

Alternative formulation:

$$k \geq \frac{1}{\sigma_1} \log \left(\frac{f(x^0) - f(x^*)}{\epsilon} \right) \Rightarrow \mathbb{E}[f(x^k) - f(x^*)] \leq \epsilon$$

Three Convergence Rates

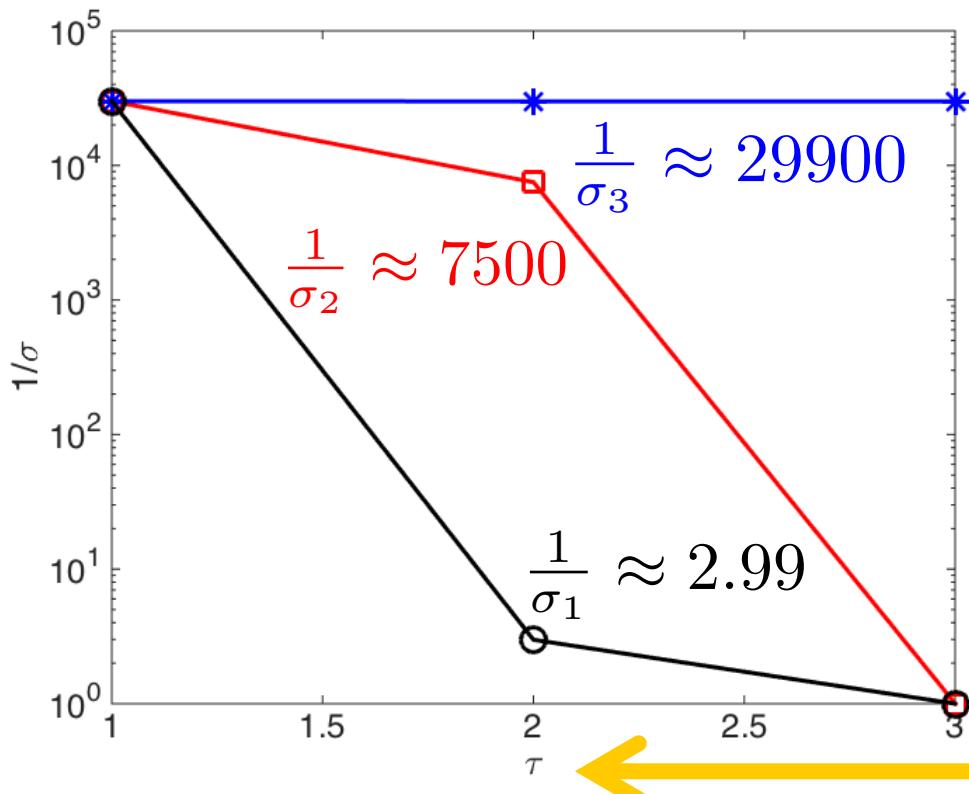
3 Convergence Rates

Theorem (QRTF'15)

$$0 < \sigma_3 \leq \sigma_2 \leq \sigma_1 \leq 1$$

Example: A Quadratic in 3D ($n = 3$)

$$\min_{x \in \mathbb{R}^3} \left[f(x) = \frac{1}{2} x^T \mathbf{M} x + b^T x + c \right]$$



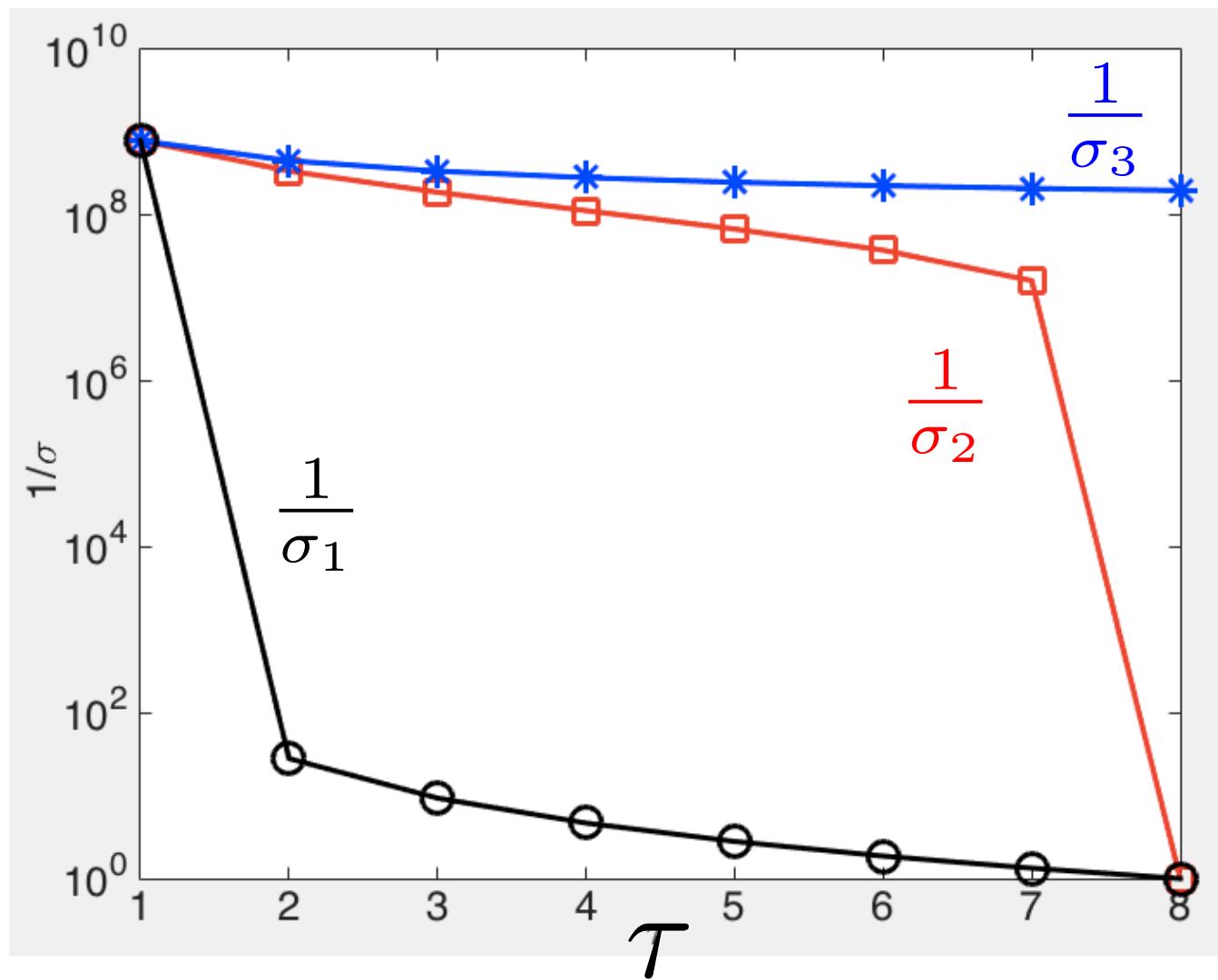
$$\mathbf{M} = \begin{pmatrix} 1.0000 & 0.9900 & 0.9999 \\ 0.9900 & 1.0000 & 0.9900 \\ 0.9999 & 0.9900 & 1.0000 \end{pmatrix}$$

condition number $\approx 3 \times 10^4$

$$\tau = \mathbb{E} [|S_k|]$$

Another Quadratic

$n = 8$
 $f = \text{quadratic}$
 $(\mathbf{G} = \mathbf{M})$



Part B

Empirical Risk Minimization

Primal Problem

$$|\phi'_i(a) - \phi'_i(b)| \leq \frac{1}{\gamma} |a - b| \quad \forall a, b \in \mathbb{R}$$

$P = \text{Regularized Empirical Risk}$

$1/\gamma$ - smooth & convex
functions (“risk”)

positive
regularization
parameter

$$\min_{w \in \mathbb{R}^d} P(w) \equiv \frac{1}{n} \sum_{i=1}^n \phi_i(A_i^\top w) + \lambda g(w)$$

$w = \text{linear predictor}$

n data vectors
 (“examples”)

$d = \# \text{ features}$
(parameters)

1 - strongly convex function (“regularizer”)

$$g(w) \geq g(w') + \langle \nabla g(w'), w - w' \rangle + \frac{1}{2} \|w - w'\|^2, \quad w, w' \in \mathbb{R}^d$$

Dual Problem

n dual variables: as many as
examples in the primal

$\in \mathbb{R}^d$

$$\max_{\alpha \in \mathbb{R}^n} \left[D(\alpha) \equiv -\lambda g^* \left(\frac{1}{\lambda n} \sum_{i=1}^n A_i \alpha_i \right) - \frac{1}{n} \sum_{i=1}^n \phi_i^*(-\alpha_i) \right]$$

1 – smooth & convex

γ - strongly convex

$$g^*(w') = \max_{w \in \mathbb{R}^d} \{(w')^\top w - g(w)\}$$

$$\phi_i^*(a') = \max_{a \in \mathbb{R}^m} \{(a')^\top a - \phi_i(a)\}$$

SDNA

Initialization: $\alpha^0 \in \mathbb{R}^n \quad \bar{\alpha}^0 = \frac{1}{\lambda n} \mathbf{A} \alpha^0$

Iterate:

Primal update: $w^k = \nabla g^*(\bar{\alpha}^k)$

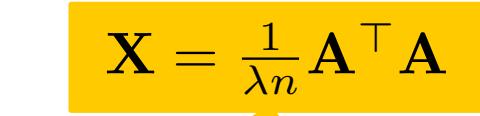
Generate a random set S_k

Compute:

$$h^k = \arg \min_{h \in \mathbb{R}^n} ((\mathbf{A}^\top w^k)_{S_k})^\top h + \frac{1}{2} h^\top \mathbf{X}_{S_k} h + \sum_{i \in S_k} \phi_i^*(-\alpha_i^k - h_i)$$

Dual update: $\alpha^{k+1} \leftarrow \alpha^k + \sum_{i \in S_k} h_i^k e_i$

Maintain average: $\bar{\alpha}^{k+1} = \bar{\alpha}^k + \frac{1}{\lambda n} \sum_{i \in S_k} h_i^k A_i$



Convergence of SDNA

Theorem (QRTF'15)

Assume that S_k is uniform

$$\mathbb{E}[P(w^k) - D(\alpha^k)] \leq (1 - \sigma_1^{prox})^k \frac{D(\alpha^*) - D(\alpha^0)}{\theta(S_k)}$$

Expected duality gap
after k iterations

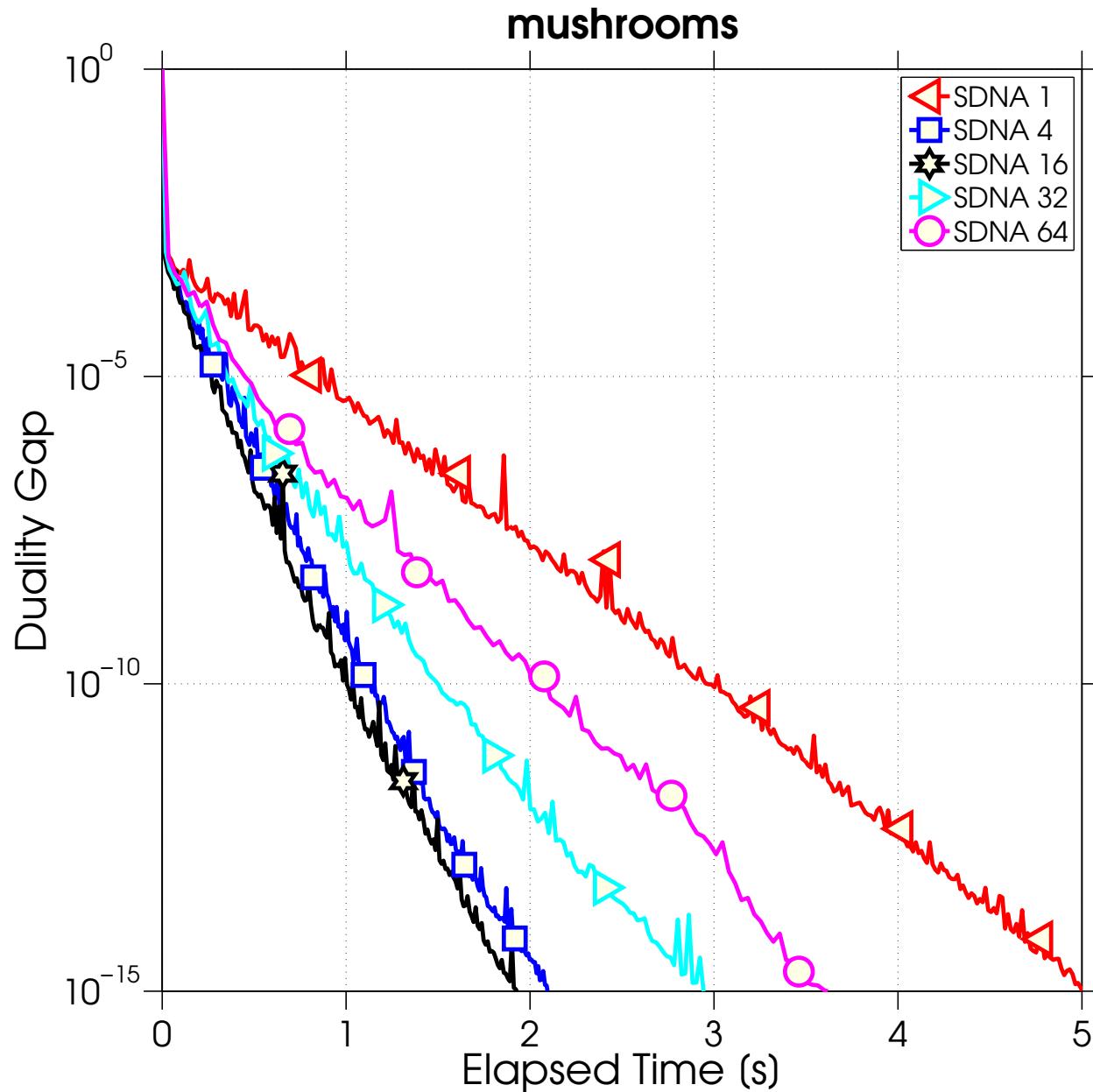
$$\sigma_1^{prox} = \frac{\tau}{n} \min\{1, s_1\}$$

$$\tau = \mathbb{E}[|S_k|] \quad s_1 = \lambda_{\min} \left[\left(\frac{1}{\tau \gamma \lambda} \mathbb{E}[(\mathbf{A}^\top \mathbf{A})_{S_k}] + \mathbf{I} \right)^{-1} \right]$$

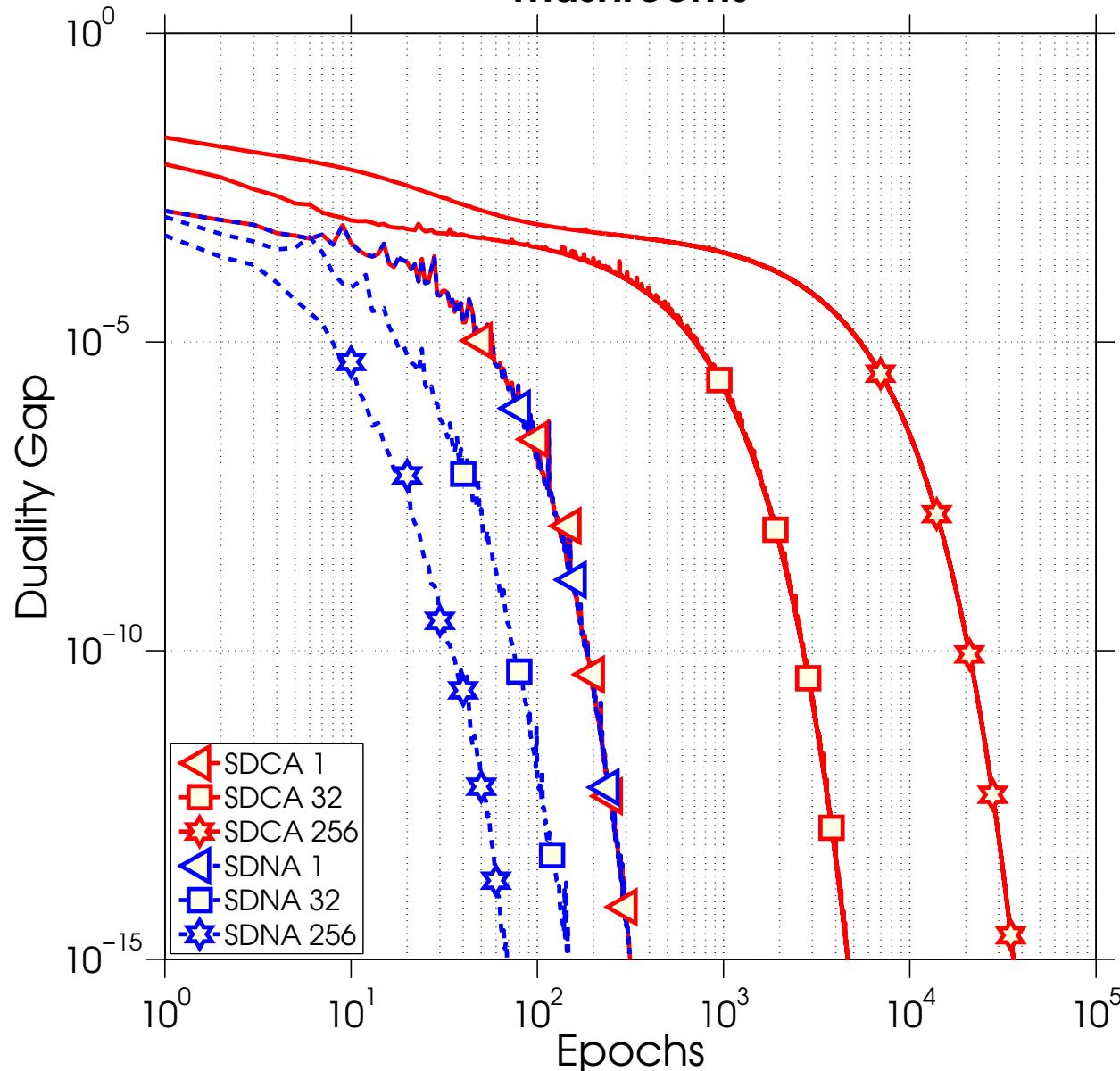
We prove this is a **better rate than that of SDCA** (which we needed to generalize to uniform samplings for comparison)

Real Dataset:
mushrooms

$d = 112$ $n = 8,124$



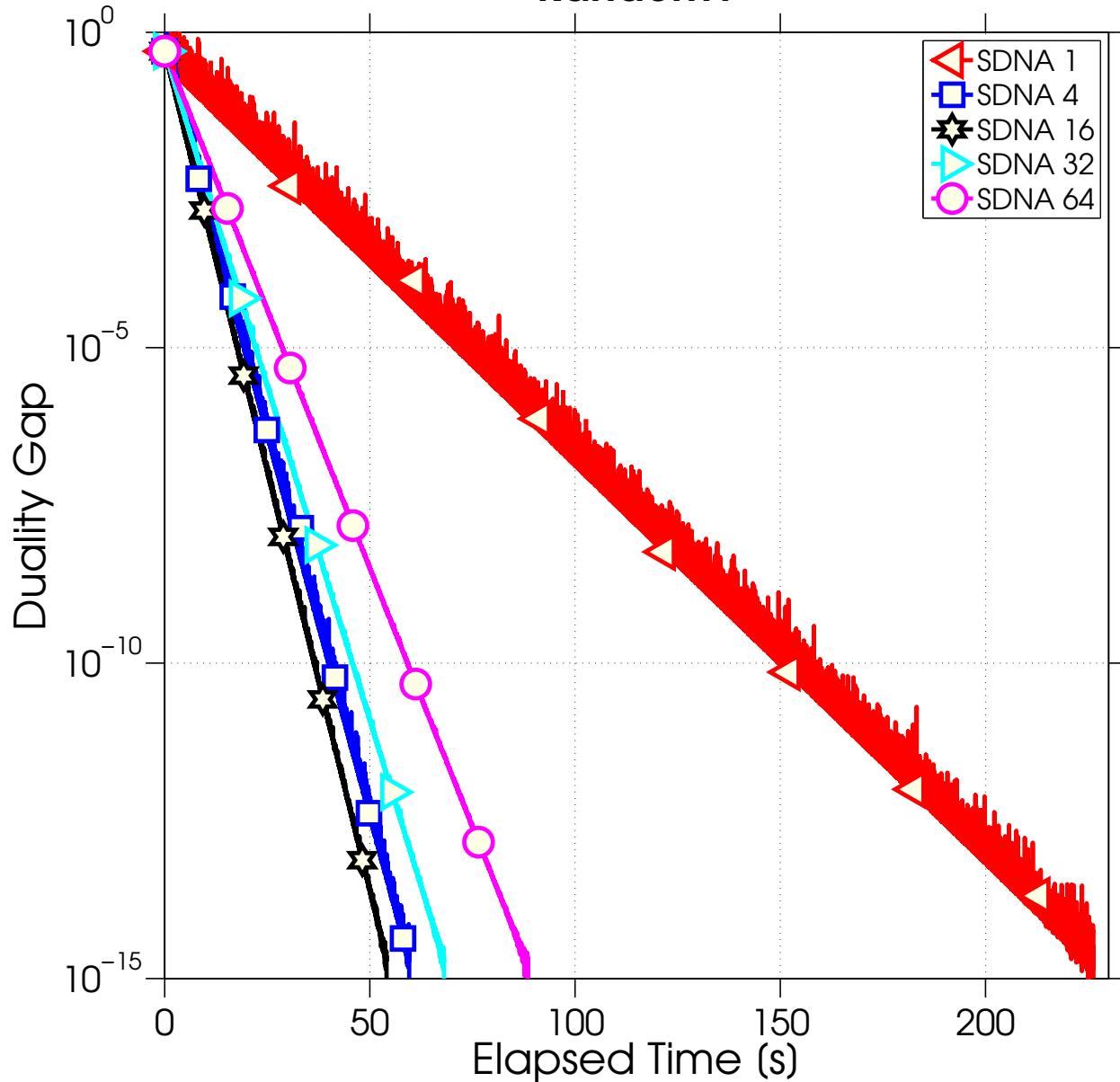
mushrooms

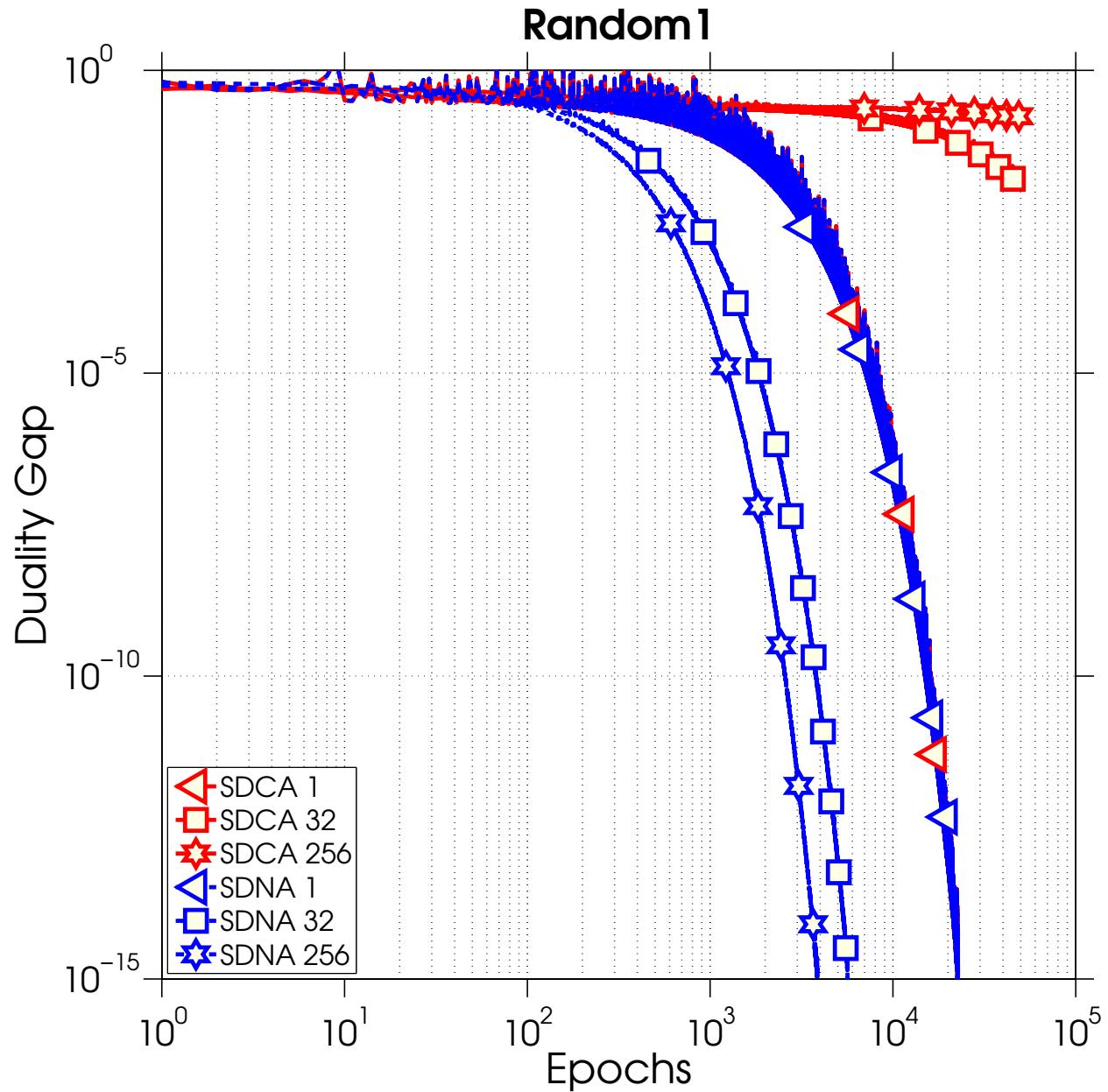


Synthetic Dataset

$d = 1,024 \quad n = 2,048$

Random 1

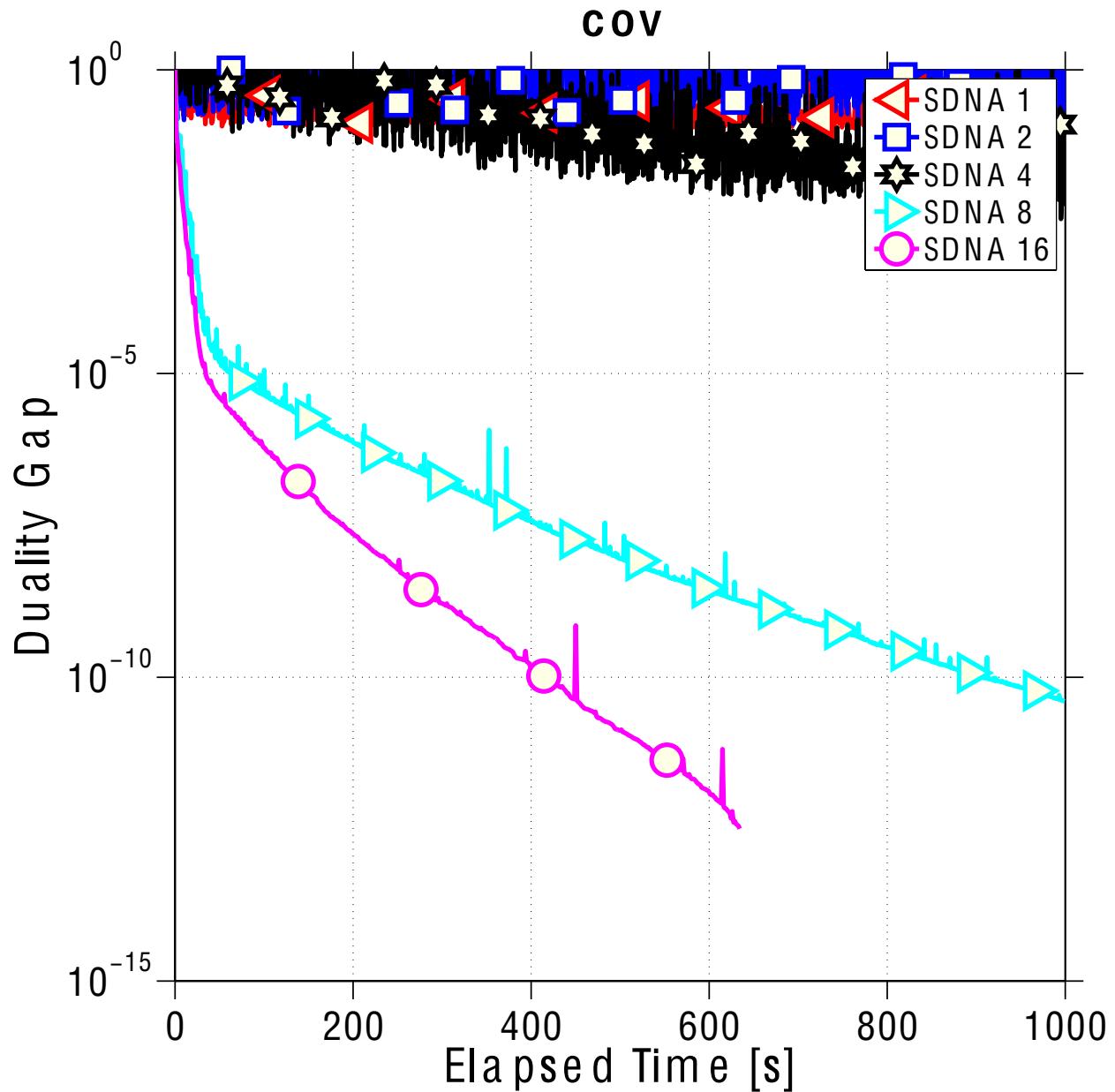




Real Dataset:

COV

$d = 54$ $n = 581,012$



END