Introduction
 SGD
 General analysis
 Special cases
 Discussion
 Bibliography

 A Unified Theory of SGD: Variance Reduction,
 Sampling, Quantization and Coordinate Descent
 The talk is based on the work [2] and slides [6]

#### Eduard Gorbunov<sup>1</sup> Filip Hanzely<sup>2</sup> Peter Richtárik<sup>2 1</sup>

 $^1{\rm Moscow}$  Institute of Physics and Technology, Russia  $^2{\rm King}$  Abdullah University of Science and Technology, Saudi Arabia

INRIA, Paris, 18 October, 2019



Unified SGD

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography

## Table of Contents

#### 1 Introduction

SGD

- 3 General analysis
- **4** Special cases
- 6 Discussion
- 6 Bibliography

Introduction	SGD	<b>General analysis</b>	Special cases	Discussion	Bibliography
●0	oo	0000000		0	0

## The Problem



Eduard Gorbunov (MIPT)

Unified SGD

INRIA, 18.10.2019 3 / 27

<ロト <四ト < 臣ト < 臣ト

E

DQC

(1)

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
•0					

### The Problem

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

#### • f(x) - convex function with Lipschitz gradient.

4 円

э

(1)

590

Introduction ●○	SGD oo	General analysis	Special cases	Discussion O	Bibliography 0

#### The Problem

$$\min_{x \in \mathbb{R}^d} f(x) + R(x) \tag{1}$$

- f(x) convex function with Lipschitz gradient.
- $R: \mathbb{R}^d \to \mathbb{R} \cup \{+\infty\}$  proximable (proper closed convex) regularizer

3

# Structure of f in supervised learning theory and practice

Introduction SGD General analysis Special cases Discussion Bibliograph ⊙● 00 000000 00000000 0 0

#### Structure of f in supervised learning theory and practice

$$f(x) = \mathbf{E}_{\xi \sim \mathcal{D}} \left[ f_{\xi}(x) \right]$$
(2)

Introduction SGD General analysis Special cases Discussion Bibliography ⊙● 00 0000000 00000000 0 0

#### Structure of f in supervised learning theory and practice

$$f(x) = \mathbf{E}_{\xi \sim \mathcal{D}} \left[ f_{\xi}(x) \right]$$
(2)

$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x).$$
 (3)

 Introduction
 SGD
 General analysis
 Special cases
 Discussion
 Bibliography

 0●
 00
 0000000
 0
 0
 0

#### Structure of f in supervised learning theory and practice

We focus on the following situations

$$f(x) = \mathbf{E}_{\xi \sim \mathcal{D}} \left[ f_{\xi}(x) \right]$$
(2)

$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x).$$
 (3)

•  $f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x)$  and

 Introduction
 SGD
 General analysis
 Special cases
 Discussion
 Bibliography

 ○●
 00
 0000000
 0
 0
 0

#### Structure of f in supervised learning theory and practice

$$f(x) = \mathbf{E}_{\xi \sim \mathcal{D}} \left[ f_{\xi}(x) \right]$$
(2)

$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x).$$
 (3)

• 
$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x)$$
 and  
 $f_i(x) = \frac{1}{m} \sum_{j=1}^{m} f_{ij}(x),$ 
(4)

 Introduction
 SGD
 General analysis
 Special cases
 Discussion
 Bibliography

 ○●
 00
 0000000
 0
 0
 0

#### Structure of f in supervised learning theory and practice

We focus on the following situations

$$f(x) = \mathbf{E}_{\xi \sim \mathcal{D}} \left[ f_{\xi}(x) \right]$$
(2)

$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x).$$
 (3)

• 
$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x)$$
 and  
 $f_i(x) = \frac{1}{m} \sum_{j=1}^{m} f_{ij}(x),$ 
(4)

Typical case:

 Introduction
 SGD
 General analysis
 Special cases
 Discussion
 Bibliography

 o●
 00
 0000000
 0
 0
 0

#### Structure of f in supervised learning theory and practice

We focus on the following situations

$$f(x) = \mathbf{E}_{\xi \sim \mathcal{D}} \left[ f_{\xi}(x) \right]$$
(2)

$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x).$$
 (3)

• 
$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x)$$
 and  
 $f_i(x) = \frac{1}{m} \sum_{j=1}^{m} f_{ij}(x),$ 
(4)

Typical case: Using the exact gradient of  $\nabla f(x)$  is too expensive

 Introduction
 SGD
 General analysis
 Special cases
 Discussion
 Bibliography

 ○●
 00
 0000000
 0
 0
 0

#### Structure of f in supervised learning theory and practice

We focus on the following situations

$$f(x) = \mathbf{E}_{\xi \sim \mathcal{D}} \left[ f_{\xi}(x) \right]$$
(2)

$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x).$$
 (3)

•  $f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x)$  and

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x),$$
 (4)

Typical case: Using the exact gradient of  $\nabla f(x)$  is too expensive, but an unbiased estimator of  $\nabla f(x)$  can be computed efficiently.

Eduard Gorbunov (MIPT)

Unified SGD

INRIA, 18.10.2019

4/27

Introduction

SGD

•0

General analysis

Special cases

Discussion 0 Bibliography 0

### Popular approach to solve problem (1)

$$x^{k+1} = \operatorname{prox}_{\gamma R}(x^k - \gamma g^k).$$
(5)

4 円

DQC

 Introduction
 SGD
 General analysis
 Special cases
 Discussion
 Bibliograp

 00
 ●0
 0000000
 0
 0
 0
 0

### Popular approach to solve problem (1)

$$x^{k+1} = \operatorname{prox}_{\gamma R}(x^k - \gamma g^k).$$
(5)

•  $g^k$  is unbiased gradient estimator:  $\mathbf{E}\left[g^k \mid x^k\right] = \nabla f(x^k)$ 

4 円

#### Popular approach to solve problem (1)

$$x^{k+1} = \operatorname{prox}_{\gamma R}(x^k - \gamma g^k).$$
(5)

- $g^k$  is unbiased gradient estimator:  $\mathbf{E}\left[g^k \mid x^k\right] = \nabla f(x^k)$
- $\operatorname{prox}_{R}(x) \stackrel{\text{def}}{=} \operatorname{argmin}_{u} \left\{ \frac{1}{2} \|u x\|^{2} + R(x) \right\}$

Introduction SGD General analysis S

Special cases

Discussion 0 Bibliography 0

### Popular approach to solve problem (1)

$$x^{k+1} = \operatorname{prox}_{\gamma R}(x^k - \gamma g^k).$$
(5)

- $g^k$  is unbiased gradient estimator:  $\mathbf{E}\left[g^k \mid x^k\right] = \nabla f(x^k)$
- $\operatorname{prox}_{R}(x) \stackrel{\text{def}}{=} \operatorname{argmin}_{u} \left\{ \frac{1}{2} \|u x\|^{2} + R(x) \right\}$
- $\|\cdot\|$  standard Euclidean norm

・ 何 ト ・ ヨ ト ・ ヨ ト

 Introduction
 SGD
 General analysis
 Special cases
 Discussion

 00
 00
 0000000
 000000000
 0

Bibliography 0

## Popular approach to solve problem (1)

$$x^{k+1} = \operatorname{prox}_{\gamma R}(x^k - \gamma g^k).$$
(5)

- $g^k$  is unbiased gradient estimator:  $\mathbf{E}\left[g^k \mid x^k\right] = \nabla f(x^k)$
- $\operatorname{prox}_{R}(x) \stackrel{\text{def}}{=} \operatorname{argmin}_{u} \left\{ \frac{1}{2} \|u x\|^{2} + R(x) \right\}$
- $\|\cdot\|$  standard Euclidean norm

#### The prox operator

•  $x \to \operatorname{prox}_{R}(x)$  is a function

< 回 ト く ヨ ト く ヨ ト

Introduction SGD General analysis Special cases Dis 00 0000000 00000000 0

## Popular approach to solve problem (1)

$$x^{k+1} = \operatorname{prox}_{\gamma R}(x^k - \gamma g^k).$$
(5)

- $g^k$  is unbiased gradient estimator:  $\mathbf{E}\left[g^k \mid x^k\right] = \nabla f(x^k)$
- $\operatorname{prox}_{R}(x) \stackrel{\text{def}}{=} \operatorname{argmin}_{u} \left\{ \frac{1}{2} \|u x\|^{2} + R(x) \right\}$
- $\|\cdot\|$  standard Euclidean norm

#### The prox operator

•  $x \to \operatorname{prox}_{R}(x)$  is a function

• 
$$\|\operatorname{prox}_R(x) - \operatorname{prox}_R(y)\| \le \|x - y\|$$

- 4 回 ト 4 ヨ ト

Introduction	SGD	<b>General analysis</b>	Special cases	Discussion	Bibliography
00	○●	0000000		O	0
Stochastic	gradien	t			

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
	00				

There are infinitely many ways of getting unbiased estimator with "good" properties.

• Flexibility to construct stochastic gradients in order to target desirable properties:

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
	00				

- Flexibility to construct stochastic gradients in order to target desirable properties:
  - convergence speed

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
	00				

- Flexibility to construct stochastic gradients in order to target desirable properties:
  - convergence speed
  - iteration cost

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
	00				

- Flexibility to construct stochastic gradients in order to target desirable properties:
  - convergence speed
  - iteration cost
  - overall complexity

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
	00				

- Flexibility to construct stochastic gradients in order to target desirable properties:
  - convergence speed
  - iteration cost
  - overall complexity
  - parallelizability

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
	00				

- Flexibility to construct stochastic gradients in order to target desirable properties:
  - convergence speed
  - iteration cost
  - overall complexity
  - parallelizability
  - communication cost and etc.

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
	00				

- Flexibility to construct stochastic gradients in order to target desirable properties:
  - convergence speed
  - iteration cost
  - overall complexity
  - parallelizability
  - communication cost and etc.
- Too many methods

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
	00				

- Flexibility to construct stochastic gradients in order to target desirable properties:
  - convergence speed
  - iteration cost
  - overall complexity
  - parallelizability
  - communication cost and etc.
- Too many methods
  - Hard to keep up with new results

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
	00				

- Flexibility to construct stochastic gradients in order to target desirable properties:
  - convergence speed
  - iteration cost
  - overall complexity
  - parallelizability
  - communication cost and etc.
- Too many methods
  - Hard to keep up with new results
  - Challenges in terms of analysis

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
	00				

- Flexibility to construct stochastic gradients in order to target desirable properties:
  - convergence speed
  - iteration cost
  - overall complexity
  - parallelizability
  - communication cost and etc.
- Too many methods
  - Hard to keep up with new results
  - Challenges in terms of analysis
  - Some problems with fair comparison: different assumptions are used

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
		000000			

## Bregman divergence

By  $D_f(x, y)$  we denote the Bregman divergence associated with f:

$$D_f(x,y) \stackrel{\text{def}}{=} f(x) - f(y) - \langle \nabla f(y), x - y \rangle.$$

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
		000000			

#### Assumption 1

Let  $\{x^k\}$  be the random iterates produced by proximal SGD.

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
		000000			

#### Assumption 1

Let  $\{x^k\}$  be the random iterates produced by proximal SGD.

1 The stochastic gradients  $g^k$  are unbiased

$$\mathsf{E}\left[g^{k} \mid x^{k}\right] = \nabla f(x^{k}) \quad \forall k \geq 0.$$

(6)

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
		000000			

#### Assumption 1

Let  $\{x^k\}$  be the random iterates produced by proximal SGD.

1 The stochastic gradients  $g^k$  are unbiased

$$\mathsf{E}\left[g^{k} \mid x^{k}\right] = \nabla f(x^{k}) \quad \forall k \ge 0.$$
(6)

2 There exist non-negative constants A, B, C, D<sub>1</sub>, D<sub>2</sub>, ρ and a (possibly) random sequence {σ<sup>2</sup><sub>k</sub>}<sub>k≥0</sub> such that the following two relations hold

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
		000000			

#### Assumption 1

Let  $\{x^k\}$  be the random iterates produced by proximal SGD.

1 The stochastic gradients  $g^k$  are unbiased

$$\mathsf{E}\left[g^{k} \mid x^{k}\right] = \nabla f(x^{k}) \quad \forall k \ge 0.$$
(6)

2 There exist non-negative constants A, B, C, D<sub>1</sub>, D<sub>2</sub>, ρ and a (possibly) random sequence {σ<sup>2</sup><sub>k</sub>}<sub>k≥0</sub> such that the following two relations hold

$$\mathsf{E}\left[\left\|g^{k}-\nabla f(x^{*})\right\|^{2}\mid x^{k}\right] \leq 2AD_{f}(x^{k},x^{*})+B\sigma_{k}^{2}+D_{1},\qquad(7)$$

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
		000000			

#### Assumption 1

Let  $\{x^k\}$  be the random iterates produced by proximal SGD.

1 The stochastic gradients  $g^k$  are unbiased

$$\mathsf{E}\left[g^{k} \mid x^{k}\right] = \nabla f(x^{k}) \quad \forall k \ge 0.$$
(6)

2 There exist non-negative constants  $A, B, C, D_1, D_2, \rho$  and a (possibly) random sequence  $\{\sigma_k^2\}_{k\geq 0}$  such that the following two relations hold

$$\mathsf{E}\left[\left\|g^{k}-\nabla f(x^{*})\right\|^{2}\mid x^{k}\right] \leq 2AD_{f}(x^{k},x^{*})+B\sigma_{k}^{2}+D_{1},\qquad(7)$$

$$\mathsf{E}\left[\sigma_{k+1}^{2} \mid \sigma_{k}^{2}\right] \leq (1-\rho)\sigma_{k}^{2} + 2CD_{f}(x^{k}, x^{*}) + D_{2}, \quad \forall k \geq 0$$
 (8)
General analysis

Special cases

Discussion 0 Bibliography 0

### Gradient Descent satisfies Assumption 1

#### Assumption 2

Assume that f is convex, i.e.

$$D_f(x,y) \ge 0 \quad \forall x,y \in \mathbb{R}^d,$$

and L-smooth, i.e.

$$\|
abla f(x) - 
abla f(y)\| \leq L \|x - y\| \quad \forall x, y \in \mathbb{R}^d.$$

Eduard Gorbunov (MIPT)

DQR

General analysis

Special cases

Discussion 0 Bibliography 0

## Gradient Descent satisfies Assumption 1

#### Assumption 2

Assume that f is convex, i.e.

$$D_f(x,y) \geq 0 \quad \forall x,y \in \mathbb{R}^d,$$

and L-smooth, i.e.

$$\|
abla f(x) - 
abla f(y)\| \le L \|x - y\| \quad \forall x, y \in \mathbb{R}^d.$$

Assumption 2 implies that (see Nesterov's book [4])  $\|\nabla f(x) - \nabla f(y)\|^2 \le 2LD_f(x, y) \quad \forall x, y \in \mathbb{R}^d.$ 

< 回 ト く ヨ ト く ヨ ト

General analysis

Special cases

Discussion 0 Bibliography 0

## Gradient Descent satisfies Assumption 1

#### Assumption 2

Assume that f is convex, i.e.

$$D_f(x,y) \ge 0 \quad \forall x,y \in \mathbb{R}^d,$$

and L-smooth, i.e.

$$\|
abla f(x) - 
abla f(y)\| \leq L \|x - y\| \quad \forall x, y \in \mathbb{R}^d.$$

Assumption 2 implies that (see Nesterov's book [4])

$$\|
abla f(x) - 
abla f(y)\|^2 \leq 2LD_f(x,y) \quad \forall x, y \in \mathbb{R}^d.$$

Therefore, if f satisfies Assumption 2, then gradient descent satisfies Assumption 1 with

$$A=L,B=0,D_1=0,\sigma_k=0,\rho=1, \ C_{\rm eff}, 0, \ D_2, \ C_{\rm eff}, 0, \ C_{\rm eff},$$

Eduard Gorbunov (MIPT)

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
		000000			

### Other assumptions

#### Assumption 3 (Unique solution)

The problem (1) has an unique minimizer  $x^*$ .

Introduction	SGD oo	General analysis 000●000	Special cases	Discussion 0	Bibliography 0

### Other assumptions

#### Assumption 3 (Unique solution)

The problem (1) has an unique minimizer  $x^*$ .

#### Assumption 4 ( $\mu$ -strong quasi-convexity)

There exists  $\mu > 0$  such that  $f : \mathbb{R}^d \to \mathbb{R}$  is  $\mu$ -strongly quasi-convex, i.e.

$$f(x^*) \ge f(x) + \langle \nabla f(x), x^* - x \rangle + \frac{\mu}{2} \|x^* - x\|^2, \qquad \forall x \in \mathbb{R}^d.$$
(9)

Introduction	SGD oo	General analysis 0000●00	Special cases	Discussion 0	Bibliography 0			
Main result								
Theorem 1								
Let Assumptions 1, 3 and 4 be satisfied.								

-1404

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
		0000000			

### Main result

#### Theorem 1

Let Assumptions 1, 3 and 4 be satisfied. Choose constant *M* such that  $M > \frac{B}{a}$ . Choose a stepsize satisfying

$$\mathsf{O} < \gamma \leq \min\left\{rac{1}{\mu}, rac{1}{\mathcal{A} + \mathcal{CM}}
ight\}.$$

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
		0000000			

### Main result

#### Theorem 1

Let Assumptions 1, 3 and 4 be satisfied. Choose constant M such that  $M > \frac{B}{a}$ . Choose a stepsize satisfying

$$0 < \gamma \le \min\left\{\frac{1}{\mu}, \frac{1}{A + CM}\right\}.$$
(10)

Then the iterates  $\{x^k\}_{k\geq 0}$  of proximal SGD satisfy

$$\mathbf{E}\left[V^{k}\right] \leq \max\left\{(1-\gamma\mu)^{k}, \left(1+\frac{B}{M}-\rho\right)^{k}\right\}V^{0}+\frac{(D_{1}+MD_{2})\gamma^{2}}{\min\left\{\gamma\mu, \rho-\frac{B}{M}\right\}},\tag{11}$$

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
		0000000			

### Main result

#### Theorem 1

Let Assumptions 1, 3 and 4 be satisfied. Choose constant *M* such that  $M > \frac{B}{a}$ . Choose a stepsize satisfying

$$0 < \gamma \le \min\left\{\frac{1}{\mu}, \frac{1}{A + CM}\right\}.$$
(10)

Then the iterates  $\{x^k\}_{k\geq 0}$  of proximal SGD satisfy

$$\mathbf{E}\left[V^{k}\right] \leq \max\left\{\left(1-\gamma\mu\right)^{k}, \left(1+\frac{B}{M}-\rho\right)^{k}\right\}V^{0} + \frac{\left(D_{1}+MD_{2}\right)\gamma^{2}}{\min\left\{\gamma\mu,\rho-\frac{B}{M}\right\}},$$
(11) where the Lyapunov function  $V^{k}$  is defined by  $V^{k} \stackrel{\text{def}}{=} \left\|x^{k}-x^{*}\right\|^{2} + M\gamma^{2}\sigma_{k}^{2}.$ 

w

Int oc	roduction	SGD oo	Genera 00000	l analysis ●●○		Special o	: <b>ases</b> 0000	Discus 0	sion	Bibliograp 0	hy
N	lethods										
Γ	Problem	Method	Alg #	Citation	VR?	AS?	Quant?	RCD?	Section	Result	1
l	(1)+(2)	SGD	Alg 1	[26]	X	X	X	×	A.1	Cor A.1	
	(1)+(3)	SGD-SR	Alg 2	[6]	×	1	X	×	A.2	Cor A.2	
	(1)+(3)	SGD-MB	Alg 3	NEW	×	×	X	×	A.3	Cor A.3	
	(1)+(3)	SGD-star	Alg 4	NEW	1	1	×	×	A.4	Cor A.4	
	(1)+(3)	SAGA	Alg 5	[5]	1	X	×	×	A.5	Cor A.5	
	(1)+(3)	N-SAGA	Alg 6	NEW	×	X	×	×	A.6	Cor A.6	
	(1)	SEGA	Alg 7	[11]	1	X	X	1	A.7	Cor A.7	
	(1)	N-SEGA	Alg 8	NEW	×	×	×	1	A.8	Cor A.8	
	(1)+(3)	$\mathtt{SVRG}^a$	Alg 9	[15]	1	X	×	×	A.9	Cor A.9	
	$(1)_{\pm}(3)$	I_SVRC	$\Delta \log 10$	[13 18]	1	X	×	×	A 10	Cor A 10	

	H-Dvitta	mgio	[15, 10]	•				11.10	00171.10
(1)+(3)	DIANA	Alg 11	[20, 14]	×	×	1	×	A.11	Cor A.11
(1)+(3)	$DIANA^b$	Alg 12	[20, 14]	1	×	1	×	A.11	Cor A.12
(1)+(3)	Q-SGD-SR	Alg 13	NEW	×	1	1	×	A.12	Cor A.13
(1)+(3)+(4)	VR-DIANA	Alg 14	[14]	1	×	1	×	A.13	Cor A.15
(1)+(3)	JacSketch	Alg 15	[9]	1	✓X	×	×	A.14	Cor A.16
The last of an arite stricting (in some some properties d) and now with de which fit are some and									

Table 1: List of specific existing (in some cases generalized) and new methods which fit our general analysis framework. VR = variance reduced method, AS = arbitrary sampling, Quant = supports gradient quantization, RCD = randomized coordinate descent type method. <sup>a</sup> Special case of SVRG with 1 outer loop only; <sup>b</sup> Special case of DIANA with 1 node and quantization of exact gradient.

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
		000000			

### Parameters

Method	A	В	ρ	C	$D_1$	$D_2$
SGD	2L	0	1	0	$2\sigma^2$	0
SGD-SR	$2\mathcal{L}$	0	1	0	$2\sigma^2$	0
SGD-MB	$\frac{A'+L(\tau-1)}{\tau}$	0	1	0	$\frac{D'}{\tau}$	0
SGD-star	$2\mathcal{L}$	0	1	0	Ó	0
SAGA	2L	2	1/n	L/n	0	0
N-SAGA	2L	2	1/n	L/n	$2\sigma^2$	$\frac{\sigma^2}{n}$
SEGA	2dL	2d	1/d	L/d	0	Ő
N-SEGA	2dL	2d	1/d	L/d	$2d\sigma^2$	$\frac{\sigma^2}{d}$
$\mathtt{SVRG}^a$	2L	2	0	0	0	Õ
L-SVRG	2L	2	p	Lp	0	0
DIANA	$\left(1+\frac{2\omega}{n}\right)L$	$\frac{2\omega}{n}$	$\alpha$	L lpha	$\frac{(1+\omega)\sigma^2}{n}$	$\alpha\sigma^2$
$\mathtt{DIANA}^b$	$(1+2\omega)L$	$2\omega$	$\alpha$	Llpha	õ	0
Q-SGD-SR	$2(1+\omega)\mathcal{L}$	0	1	0	$2(1+\omega)\sigma^2$	0
VR-DIANA	$\left(1+\frac{4\omega+2}{n}\right)L$	$\frac{2(\omega+1)}{n}$	$\alpha$	$\left(\frac{1}{m}+4\alpha\right)L$	0	0
JacSketch	$2\mathcal{L}_1$	$\frac{2\lambda_{\max}}{n}$	$\lambda_{ m min}$	$\frac{\mathcal{L}_2}{n}$	0	0

E

900

A B + 
 A B +
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Introduction	SGD	<b>General analysis</b>	Special cases	Discussion	Bibliography
	oo	0000000	●000000000	0	0

We checked that GD satisfies Assumption 1 when f is convex and  $L\mbox{-smooth}$  with

$$A = L, B = 0, D_1 = 0, \sigma_k = 0, \rho = 1, C = 0, D_2 = 0.$$

Introduction	<b>SGD</b>	<b>General analysis</b>	Special cases	Discussion	Bibliography
	00	0000000	●000000000	0	0

We checked that GD satisfies Assumption 1 when f is convex and L-smooth with

$$A = L, B = 0, D_1 = 0, \sigma_k = 0, \rho = 1, C = 0, D_2 = 0.$$

We can choose M = 1 in Theorem 1 and get

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
			•00000000		

We checked that GD satisfies Assumption 1 when f is convex and L-smooth with

$$A = L, B = 0, D_1 = 0, \sigma_k = 0, \rho = 1, C = 0, D_2 = 0.$$

We can choose M = 1 in Theorem 1 and get

• 
$$0 < \gamma \leq \frac{1}{L}$$
 (since  $\mu \leq L$  and  $C = 0$ )

Introduction	SGD	<b>General analysis</b>	Special cases	Discussion	Bibliography
	oo	0000000	●000000000	0	0

We checked that GD satisfies Assumption 1 when f is convex and L-smooth with

$$A = L, B = 0, D_1 = 0, \sigma_k = 0, \rho = 1, C = 0, D_2 = 0.$$

We can choose M = 1 in Theorem 1 and get

•  $0 < \gamma \leq \frac{1}{L}$  (since  $\mu \leq L$  and C = 0)

• 
$$V^{k} = ||x^{k} - x^{*}||^{2}$$
 (since  $\sigma_{k} = 0$ )

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
			000000000		

We checked that GD satisfies Assumption 1 when f is convex and L-smooth with

$$A = L, B = 0, D_1 = 0, \sigma_k = 0, \rho = 1, C = 0, D_2 = 0.$$

We can choose M = 1 in Theorem 1 and get

- $0 < \gamma \leq \frac{1}{L}$  (since  $\mu \leq L$  and C = 0)
- $V^{k} = ||x^{k} x^{*}||^{2}$  (since  $\sigma_{k} = 0$ )
- The rate:  $\|x^k x^*\|^2 \le (1 \gamma \mu)^k \|x^0 x^*\|^2$

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
			•000000000		

We checked that GD satisfies Assumption 1 when f is convex and L-smooth with

$$A = L, B = 0, D_1 = 0, \sigma_k = 0, \rho = 1, C = 0, D_2 = 0.$$

We can choose M = 1 in Theorem 1 and get

- $0 < \gamma \leq \frac{1}{L}$  (since  $\mu \leq L$  and C = 0)
- $V^{k} = ||x^{k} x^{*}||^{2}$  (since  $\sigma_{k} = 0$ )
- The rate:  $\|x^k x^*\|^2 \le (1 \gamma \mu)^k \|x^0 x^*\|^2$
- In particular, for  $\gamma = \frac{1}{L}$  we recover the standard rate for proximal gradient descent:

$$k \geq rac{L}{\mu}\lograc{1}{arepsilon} \implies \|x^k - x^*\|^2 \leq arepsilon\|x^0 - x^*\|^2.$$

General analysis

Special cases 000000000

# SGD-SR (see Gower et al. (2019), [3])

### Stochastic reformulation

$$\min_{x \in \mathbb{R}^d} f(x) + R(x), \quad f(x) = \mathbf{E}_{\mathcal{D}} \left[ f_{\xi}(x) \right], \quad f_{\xi}(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \xi_i f_i(x) \quad (12)$$

э

Image: A marked black

General analysis

Special cases 000000000

----

Image: A marked black

# SGD-SR (see Gower et al. (2019), [3])

### Stochastic reformulation

$$\min_{x \in \mathbb{R}^d} f(x) + R(x), \quad f(x) = \mathbf{E}_{\mathcal{D}} \left[ f_{\xi}(x) \right], \quad f_{\xi}(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \xi_i f_i(x) \quad (12)$$

**1** 
$$\xi \sim \mathcal{D}$$
:  $\mathbf{E}_{\mathcal{D}}[\xi_i] = 1$  for all *i*

э

General analysis

Special cases

Discussion 0

n

Bibliography 0

## SGD-SR (see Gower et al. (2019), [3])

#### Stochastic reformulation

$$\min_{x \in \mathbb{R}^d} f(x) + R(x), \quad f(x) = \mathbf{E}_{\mathcal{D}} \left[ f_{\xi}(x) \right], \quad f_{\xi}(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \xi_i f_i(x) \quad (12)$$

$$\mathbf{0} \ \xi \sim \mathcal{D}: \ \mathbf{E}_{\mathcal{D}}[\xi_i] = 1 \text{ for all } i$$

2  $f_i$  (for all i) is smooth, possibly non-convex function

э

∃ ► < ∃ ►</p>

Image: A marked black

General analysis

Special cases

Discussion 0

n

Bibliography 0

## SGD-SR (see Gower et al. (2019), [3])

#### Stochastic reformulation

$$\min_{x\in\mathbb{R}^d} f(x) + R(x), \quad f(x) = \mathsf{E}_{\mathcal{D}}\left[f_{\xi}(x)\right], \quad f_{\xi}(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \xi_i f_i(x) \quad (12)$$

**1** 
$$\xi \sim \mathcal{D}$$
:  $\mathbf{E}_{\mathcal{D}}[\xi_i] = 1$  for all  $i$ 

**2**  $f_i$  (for all *i*) is smooth, possibly non-convex function

#### Algorithm 4 SGD-SR

Input: learning rate  $\gamma > 0$ , starting point  $x^0 \in \mathbb{R}^d$ , distribution  $\mathcal{D}$  over  $\xi \in \mathbb{R}^n$ such that  $\mathbf{E}_{\mathcal{D}}[\xi]$  is vector of ones 1: for k = 0, 1, 2, ... do 2: Sample  $\xi \sim \mathcal{D}$ 3:  $g^k = \nabla f_{\xi}(x^k)$ 4:  $x^{k+1} = \operatorname{prox}_{\gamma R}(x^k - \gamma g^k)$ 5: end for

## SGD-SR (see Gower et al. (2019), [3])

### Expected smoothness

We say that f is  $\mathcal{L}$ -smooth in expectation with respect to distribution  $\mathcal{D}$  if there exists  $\mathcal{L} = \mathcal{L}(f, \mathcal{D}) > 0$  such that

$$\mathbf{E}_{\mathcal{D}}\left[\left\|\nabla f_{\xi}(x) - \nabla f_{\xi}(x^{*})\right\|^{2}\right] \leq 2\mathcal{L}D_{f}(x, x^{*}), \tag{13}$$

for all  $x \in \mathbb{R}^d$  and write  $(f, \mathcal{D}) \sim ES(\mathcal{L})$ .

General analysis

Special cases

Discussion 0 Bibliography 0

## SGD-SR (see Gower et al. (2019), [3])

Lemma 1 (Generalization of Lemma 2.4, [3])

If  $(f, \mathcal{D}) \sim \textit{ES}(\mathcal{L})$ , then

$$\mathbf{E}_{\mathcal{D}}\left[\|\nabla f_{\xi}(x) - \nabla f(x^*)\|^2\right] \leq 4\mathcal{L}D_f(x, x^*) + 2\sigma^2. \tag{14}$$

where 
$$\sigma^2 \stackrel{\text{def}}{=} \mathbf{E}_{\mathcal{D}} \left[ \| \nabla f_{\xi}(x^*) - \nabla f(x^*) \|^2 \right].$$

э

(B)

Image: A marked black

General analysis

Special cases

Discussion 0 Bibliography 0

## SGD-SR (see Gower et al. (2019), [3])

Lemma 1 (Generalization of Lemma 2.4, [3])

If  $(f, \mathcal{D}) \sim \textit{ES}(\mathcal{L})$ , then

$$\Xi_{\mathcal{D}}\left[\left\|\nabla f_{\xi}(x) - \nabla f(x^{*})\right\|^{2}\right] \leq 4\mathcal{L}D_{f}(x, x^{*}) + 2\sigma^{2}.$$
 (14)

where 
$$\sigma^2 \stackrel{\text{def}}{=} \mathbf{E}_{\mathcal{D}} \left[ \| \nabla f_{\xi}(x^*) - \nabla f(x^*) \|^2 \right].$$

That is, SGD-SR satisfies Assumption 1 with

$$A = 2\mathcal{L}, B = 0, D_1 = 2\sigma^2, \sigma_k = 0, \rho = 1, C = 0, D_2 = 0.$$

3

(B)

Image: A marked black

General analysis

Special cases

Discussion 0 Bibliography 0

## SGD-SR (see Gower et al. (2019), [3])

Lemma 1 (Generalization of Lemma 2.4, [3])

If  $(f,\mathcal{D})\sim \textit{ES}(\mathcal{L})$ , then

$$\mathsf{E}_{\mathcal{D}}\left[\left\|\nabla f_{\xi}(x) - \nabla f(x^{*})\right\|^{2}\right] \leq 4\mathcal{L}D_{f}(x, x^{*}) + 2\sigma^{2}. \tag{14}$$

where 
$$\sigma^2 \stackrel{\text{def}}{=} \mathbf{E}_{\mathcal{D}} \left[ \| \nabla f_{\xi}(x^*) - \nabla f(x^*) \|^2 \right].$$

That is, SGD-SR satisfies Assumption 1 with

$$A = 2\mathcal{L}, B = 0, D_1 = 2\sigma^2, \sigma_k = 0, \rho = 1, C = 0, D_2 = 0.$$

If  $\gamma^k \equiv \gamma \leq \frac{1}{2\mathcal{L}}$  then Theorem 1 implies that

$$\mathsf{E}\|x^{k} - x^{*}\|^{2} \leq (1 - \gamma\mu)^{k}\|x^{0} - x^{*}\|^{2} + \frac{2\gamma\sigma^{2}}{\mu}.$$

IntroductionSGDGeneral analysisSpecial casesDiscussionBibliograph○0○0○000000○0000000○0

### SAGA (see Defazio, Bach & Lacoste-Julien (2014), [1])

Consider the finite-sum minimization problem

$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x) + R(x), \qquad (15)$$

4 円

## IntroductionSGDGeneral analysisSpecial casesDiscussionBibliography0000000000000000000000

## SAGA (see Defazio, Bach & Lacoste-Julien (2014), [1])

Consider the finite-sum minimization problem

$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x) + R(x), \qquad (15)$$

where  $f_i$  is convex, *L*-smooth for each *i* and *f* is  $\mu$ -strongly convex.

IntroductionSGDGeneral analysisSpecial casesDiscussionBibliography0000000000000000000000

## SAGA (see Defazio, Bach & Lacoste-Julien (2014), [1])

Consider the finite-sum minimization problem

$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x) + R(x), \qquad (15)$$

where  $f_i$  is convex, L-smooth for each i and f is  $\mu$ -strongly convex.

#### Algorithm 7 SAGA [1]

Input: learning rate  $\gamma > 0$ , starting point  $x^0 \in \mathbb{R}^d$ 1: Set  $\psi_j^0 = x^0$  for each  $j \in [n]$ 2: for k = 0, 1, 2, ... do 3: Sample  $j \in [n]$  uniformly at random 4: Set  $\phi_j^{k+1} = x^k$  and  $\phi_i^{k+1} = \phi_i^k$  for  $i \neq j$ 5:  $g^k = \nabla f_j(\phi_j^{k+1}) - \nabla f_j(\phi_j^k) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\phi_i^k)$ 6:  $x^{k+1} = \operatorname{prox}_{\gamma R} (x^k - \gamma g^k)$ 7: end for

## SAGA (see Defazio, Bach & Lacoste-Julien (2014), [1])

#### Lemma 2

We have

$$\mathsf{E}\left[\left\|g^{k}-\nabla f(x^{*})\right\|^{2} \mid x^{k}\right] \leq 4LD_{f}(x^{k},x^{*})+2\sigma_{k}^{2}$$

-

< 47 ▶

Special cases

Discussion 0 Bibliography 0

## SAGA (see Defazio, Bach & Lacoste-Julien (2014), [1])

#### Lemma 2

We have

$$\mathsf{E}\left[\left\|g^{k}-\nabla f(x^{*})\right\|^{2}\mid x^{k}\right] \leq 4LD_{f}(x^{k},x^{*})+2\sigma_{k}^{2} \tag{16}$$

and

$$\mathsf{E}\left[\sigma_{k+1}^2 \mid x^k\right] \le \left(1 - \frac{1}{n}\right)\sigma_k^2 + \frac{2L}{n}D_f(x^k, x^*),\tag{17}$$

Image: A marked black

F

э

- ∢ ≣ →

## SAGA (see Defazio, Bach & Lacoste-Julien (2014), [1])

#### Lemma 2

We have

$$\mathsf{E}\left[\left\|g^{k}-\nabla f(x^{*})\right\|^{2}\mid x^{k}\right] \leq 4LD_{f}(x^{k},x^{*})+2\sigma_{k}^{2} \tag{16}$$

and

where

$$\mathbf{E}\left[\sigma_{k+1}^{2} \mid x^{k}\right] \leq \left(1 - \frac{1}{n}\right)\sigma_{k}^{2} + \frac{2L}{n}D_{f}(x^{k}, x^{*}),$$
(17)  
$$\sigma_{k}^{2} = \frac{1}{n}\sum_{i=1}^{n}\left\|\nabla f_{i}(\phi_{i}^{k}) - \nabla f_{i}(x^{*})\right\|^{2}.$$

э

- ∢ ≣ →

Image: A marked black

## SAGA (see Defazio, Bach & Lacoste-Julien (2014), [1])

That is, SAGA satisfies Assumption 1 with

$$A = 2L, B = 1, D_1 = 0, \rho = \frac{1}{n}, C = \frac{L}{n}, D_2 = 0.$$

roduction SC

## SAGA (see Defazio, Bach & Lacoste-Julien (2014), [1])

That is, SAGA satisfies Assumption 1 with

$$A = 2L, B = 1, D_1 = 0, \rho = \frac{1}{n}, C = \frac{L}{n}, D_2 = 0.$$

Theorem 1 with M = 4n implies that for  $\gamma = \frac{1}{6L}$ 

$$\mathbf{E}V^k \leq \left(1 - \min\left\{\frac{\mu}{6L}, \frac{1}{2n}\right\}\right)^2 V^0.$$

Introduction	<b>SGD</b>	<b>General analysis</b>	Special cases	Discussion	Bibliography
	00	0000000	0000000●00	0	0

### SGD-star

Consider the same situation as for SGD-SR. Recall that for SGD-SR we got

$$\mathsf{E}\|x^k - x^*\|^2 \leq (1 - \gamma \mu)^k \|x^0 - x^*\|^2 + \frac{2\gamma \sigma^2}{\mu}, \quad \gamma \in \left(0, \frac{1}{2\mathcal{L}}\right]$$

- 4 E

3

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
	oo	0000000	0000000●00	0	0

#### SGD-star

Consider the same situation as for SGD-SR. Recall that for SGD-SR we got

$$\mathbf{\mathsf{E}} \|x^k - x^*\|^2 \leq (1 - \gamma \mu)^k \|x^0 - x^*\|^2 + \frac{2\gamma \sigma^2}{\mu}, \quad \gamma \in \left(0, \frac{1}{2\mathcal{L}}\right]$$

#### Algorithm 9 SGD-star

**Input:** learning rate  $\gamma > 0$ , starting point  $x^0 \in \mathbb{R}^d$ , distribution  $\mathcal{D}$  over  $\xi \in \mathbb{R}^n$  such that  $\mathbf{E}_{\mathcal{D}}[\xi]$  is vector of ones

1: for k = 0, 1, 2, ... do

2: Sample 
$$\xi \sim \mathcal{D}$$

3: 
$$g_{k+1}^{k} = \nabla f_{\xi}(x^{k}) - \nabla f_{\xi}(x^{*}) + \nabla f(x^{*})$$

- 4:  $x^{k+1} = \operatorname{prox}_{\gamma R}(x^k \gamma g^k)$
- 5: end for

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
			0000000000		

### SGD-star

### Lemma 3 (Lemma 2.4 from [3])

If  $(f,\mathcal{D})\sim \textit{ES}(\mathcal{L})$ , then

$$\mathbf{E}_{\mathcal{D}}\left[\left\|g^{k}-\nabla f(x^{*})\right\|^{2}\right] \leq 4\mathcal{L}D_{f}(x^{k},x^{*}).$$
(18)

Thus, SGD-star satisfies Assumption 1 with

$$A = 2\mathcal{L}, B = 0, D_1 = 0, \sigma_k = 0, \rho = 1, C = 0, D_2 = 0.$$

3

(B)

< 47 ▶
Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
			0000000000		

#### SGD-star

### Lemma 3 (Lemma 2.4 from [3])

If  $(f,\mathcal{D})\sim \textit{ES}(\mathcal{L})$ , then

$$\mathbf{E}_{\mathcal{D}}\left[\left\|g^{k}-\nabla f(x^{*})\right\|^{2}\right] \leq 4\mathcal{L}D_{f}(x^{k},x^{*}).$$
(18)

Thus, SGD-star satisfies Assumption 1 with

$$A = 2\mathcal{L}, B = 0, D_1 = 0, \sigma_k = 0, \rho = 1, C = 0, D_2 = 0.$$

If  $\gamma^k \equiv \gamma \leq \frac{1}{2L}$  then Theorem 1 implies that

$$\mathsf{E} \|x^{k} - x^{*}\|^{2} \leq (1 - \gamma \mu)^{k} \|x^{0} - x^{*}\|^{2}.$$

1

# Quantitative definition of variance reduction

#### Variance-reduced methods

We say that SGD with stochastic gradients satisfying Assumption 1 is variance-reduced if  $D_1 = D_2 = 0$ .

# Quantitative definition of variance reduction

#### Variance-reduced methods

We say that SGD with stochastic gradients satisfying Assumption 1 is variance-reduced if  $D_1 = D_2 = 0$ .

• Variance-reduced methods converge to the exact solution

# Quantitative definition of variance reduction

#### Variance-reduced methods

We say that SGD with stochastic gradients satisfying Assumption 1 is variance-reduced if  $D_1 = D_2 = 0$ .

- Variance-reduced methods converge to the exact solution
- The sequence  $\{\sigma_k^2\}_{k\geq 0}$  reflects the progress of the variance reduction process.

 Introduction
 SGD
 General analysis
 Special cases
 Discussion
 Bibliography

 00
 0000000
 000000000
 ●
 0

### Limitations and Extensions

**1** We consider only *L*-smooth  $\mu$ -strongly quasi-convex case.

< □ > < 凸

- **1** We consider only *L*-smooth  $\mu$ -strongly quasi-convex case.
- It would interesting to unify the theory for biased gradients estimator (e.g. SAG [7], SARAH [5], zero-order optimization and etc.)

- **1** We consider only *L*-smooth  $\mu$ -strongly quasi-convex case.
- It would interesting to unify the theory for biased gradients estimator (e.g. SAG [7], SARAH [5], zero-order optimization and etc.)
- Our analysis doesn't recover the best known rates for RCD type methods with importance sampling.

- **1** We consider only *L*-smooth  $\mu$ -strongly quasi-convex case.
- It would interesting to unify the theory for biased gradients estimator (e.g. SAG [7], SARAH [5], zero-order optimization and etc.)
- Our analysis doesn't recover the best known rates for RCD type methods with importance sampling.
- **4** An extension of iteration dependent parameters  $A, B, C, D_1, D_2, \rho$  would cover a new methods, such as SGD with decreasing stepsizes.

- **1** We consider only *L*-smooth  $\mu$ -strongly quasi-convex case.
- It would interesting to unify the theory for biased gradients estimator (e.g. SAG [7], SARAH [5], zero-order optimization and etc.)
- Our analysis doesn't recover the best known rates for RCD type methods with importance sampling.
- **4** An extension of iteration dependent parameters  $A, B, C, D_1, D_2, \rho$  would cover a new methods, such as SGD with decreasing stepsizes.
- We consider only non-accelerated methods. It would be interesting to provide an unified analysis of stochastic methods with acceleration and momentum.

IntroductionSGDGeneral analysis00000000000

Special cases

# Bibliography I

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*, pages

1646–1654, 2014.

Eduard Gorbunov, Filip Hanzely, and Peter Richtárik. A unified theory of sgd: Variance reduction, sampling, quantization and coordinate descent. arXiv preprint arXiv:1905.11261, 2019.

Introduction	SGD	<b>General analysis</b>	Special cases	Discussion	Bibliography
00	oo	0000000	0000000000	0	●
Bibliogra	phy II				

 Robert M Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik.
 SGD: General analysis and improved rates. arXiv preprint arXiv:1901.09401, 2019.

Yurii Nesterov.

Introductory lectures on convex optimization: a basic course.

2004.

Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. SARAH: A novel method for machine learning problems using stochastic recursive gradient.

In Proceedings of the 34th International Conference on Machine Learning, volume 70 of Proceedings of Machine Learning Research, pages 2613–2621. PMLR, 2017.

Introduction	SGD	General analysis	Special cases	Discussion	Bibliography
					•

# Bibliography III



#### Peter Richtárik.

A guided walk through the zoo of stochastic gradient descent methods.

#### ICCOPT 2019 Summer School.

Nicolas Le Roux, Mark Schmidt, and Francis Bach.

A stochastic gradient method with an exponential convergence rate for finite training sets.

In Advances in Neural Information Processing Systems, pages 2663–2671, 2012.