

# Introduction to stochastic optimization for large-scale problems

Dr. Ir. Valentin Leplat

Skoltech  
Center for Artificial Intelligence and Technology

*Email: [v.leplat@skoltech.ru](mailto:v.leplat@skoltech.ru)*

*Website: <https://sites.google.com/view/valentinleplat>*

**ISP Seminar - Part 2**

24 th Jan. 2023

# Table of Contents

- 1 Proximal stochastic gradient
- 2 Noise Reduction Methods
- 3 Numerical tests - hard-coded NN

# Context

Here we tackle the following problems:

$$\min_w F(w) + g(w)$$

where  $F(w)$  can be the *expected risk* or the **empirical risk**,  $g(w)$  is a regularization function for promoting structure on the solution  $w$  or avoid over-fitting (ex: Tikhonov regularization).

**Strongly-convex assumption:** given  $F(w)$  convex, using a strongly-convex function  $g(w)$  makes the objective function  $F(w) + g(w)$  strongly convex, hence ensuring good rates for SGD, and a unique minimizer.

**Moreover:** for general functions, it improves the local smoothness, in particular we may get "locally Lipschitz" function, even if it is not globally, improving local convergence.

# Assumptions

**For simplicity**, we assume:

$g$  has a simple proximal operator (in the sense that it can be computed easily).

**Recall:** the proximal operator associated to a function  $g$  with parameter  $\mu > 0$ , is a mapping  $\text{prox}_{g; \mu} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , calculated as follows:

$$\text{prox}_{g; \mu}(v) = \arg \min_{w \in \mathbb{R}^d} \left\{ g(w) + \frac{1}{2\mu} \|w - v\|^2 \right\}$$

$\mathbb{R}^d$  an unconstrained opt. prob. on its own !!

Some proximal operators can be computed in closed form :) !  $\mathbb{R}^d$

**Example:**  $g(w) = \lambda \|w\|_1$ , see [ADMM lecture notes](#) (slides 39 -47) for an introduction and examples.

# Proximal Stochastic Gradient Descent Algorithm

## Initialization

Set  $w_1 \in W$ .

## Iteration ( $k \geq 1$ ):

Generate a realization of the random variable  $\xi_k$

Compute a stochastic vector  $r_k = \nabla f(w_k; \xi_k)$

Choose a step size  $\alpha_k > 0$

Set

$$v_k = w_k - \alpha_k r_k \quad (1)$$

Set

$$w_{k+1} = \text{prox}_{g; \frac{1}{2\alpha_k}}(v_k) = \arg \min_{w \in W} \left\{ g(w) + \frac{1}{2\alpha_k} \|w - v_k\|^2 \right\} \quad (2)$$

# Proximal Stochastic Gradient Descent Algorithm

## Interpretation:

$$w_{k+1} \in \text{prox}_{g; \frac{1}{2k}}(w_k - \eta \nabla f(w_k; \xi_k))$$

from the definition of the proximal operator:

$$w_{k+1} = \arg \min_{w \in \mathcal{W}} \left\{ g(w) + \frac{1}{2k} \|w - w_k + \eta \nabla f(w_k; \xi_k)\|^2 \right\}$$

$$\arg \min_{w \in \mathcal{W}} \left\{ g(w) + \frac{1}{2k} \|w - w_k + \eta \nabla f(w_k; \xi_k)\|^2 \right\}$$

$w_{k+1}$  minimizes  $g(w)$  plus a simple quadratic local model of  $f(w; \xi_k)$  around  $w_k$ .

# Proximal Stochastic Gradient Descent Algorithm

## Interpretation:

$$w_{k+1} \in \text{prox}_{g; \gamma} \{ w_k - \gamma \nabla f(w_k) \}$$

The term **proximal** refers to the presence of the third term in the minimization problem on the right-hand side,

encourages the new iterate to be close to  $w_k$ ,

Notice that if the regularization (i.e., last) term were not present, we would exactly recover the stochastic gradient method

# Prox. SGD for strongly cvx functions - decreasing $k$

**Here:** we give a brief convergence result for Proximal SGD, first we make several assumptions

## Assumptions 3

$g$  is a proper convex, closed and proper function,

$F_{pwq}$  is  $c$ -strongly convex and has a  $L$ -Lipschitz gradient,

there exists  $M > 0$  such that  $V_k(r) \leq M$ .

there exists  $w \in \arg \min F \cap g$ ,

the sequence  $\alpha_k$  is deterministic, satisfies  $\alpha_k \leq \frac{a}{c^k b^k}$  for given  $a > 1$

and  $b > 0$  such that  $\frac{a}{cb} > \frac{1}{2L}$



# Prox. SGD for strongly cvx functions - decreasing $k$

## Theorem 5 (Proof Part 3)

Under assumptions 3, the iterates of the proximal stochastic gradient algorithm satisfy the convergence rate (in terms of distance to the optimum):

$$\mathbb{E} \|w_k - w^*\|^2 \leq \frac{4a^2 M}{\rho a - 1} \frac{1}{k} + \frac{c^2}{b} \quad (3)$$

# Table of Contents

- 1 Proximal stochastic gradient
- 2 Noise Reduction Methods
- 3 Numerical tests - hard-coded NN

# Motivation

In previous sections:

The theoretical arguments developed have led many in the machine learning community to view SGD as the ideal optimization approach for large-scale applications.

however: that this is far from settled !

Indeed: SGD suffers from, among other things, the adverse effect of noisy gradient estimates:

prevents it from converging to the solution when fixed stepsizes are used,

leads to a "slow", sublinear rate of convergence when a diminishing stepsize sequence  $\alpha_k \propto \frac{1}{k}$  is employed.

# Motivation

In this section:

Solution: develop methods endowed with noise reduction capabilities (reduce the errors in the gradient estimates and/or iterate sequence)

Three main classes

**dynamic sampling methods**, achieve noise reduction by gradually increasing the minibatch size used in the gradient computation

• employing increasingly more accurate gradient estimates

**Gradient aggregation methods**: improve the quality of the search directions by storing gradient estimates corresponding to samples employed in previous iteration.

**iterate averaging methods**: accomplish noise reduction not by averaging gradient estimates, but by maintaining an average of iterates computed during the optimization process.

## high-level summary for convergence rates

**dynamic sampling methods:** *linear* rate of convergence to the optimal value using a  $\alpha \cdot k^{-\beta}$  stepsize.

**Gradient aggregation methods:** *linear* rate of convergence to the optimal value using a  $\alpha \cdot k^{-\beta}$  stepsize.

**iterate averaging methods:** uses a more aggressive stepsize sequence, of order  $O_p \frac{1}{k^{\beta}}$  instead of  $O_p \frac{1}{k}$   
 $\tilde{N}$  really appealing because:....it is this sequence of averaged iterates that converges to the solution !

# Organization of this section

formal motivation of noise reduction concept: a fundamental result that stipulates a rate of decrease in noise that allows an SGD-type method to converge at a linear rate.

three gradient aggregation methods — SVRG, SAGA, and SAG.

discussion of the practical and theoretical properties of iterate averaging methods.

# Reducing Noise at a Geometric Rate

Let us recall Lemma 1:

$$\mathbb{E}_k \langle r, Fp_{w_k} - q \rangle = \langle Fp_{w_k} - q, r \rangle \propto \mathbb{E}_k \langle r, Fp_{w_k} - q \rangle; \mathbb{E}_k \langle r, Fp_{w_k} - q \rangle \leq -\frac{2L}{2} \mathbb{E}_k \langle r \rangle \langle Fp_{w_k} - q \rangle^2 \quad (4)$$

**Intuitively:**

if  $\langle r, Fp_{w_k} - q \rangle$  is a descent direction in expectation (first term of RHS is negative),

and if we "manage" to decrease  $\mathbb{E}_k \langle r \rangle \langle Fp_{w_k} - q \rangle^2$  "fast enough",

**then:** the effect of having noisy directions will not impede a fast rate of convergence !

## Reducing Noise at a Geometric Rate

**Formally:** Consider SGD with fixed step sizes for **strongly convex functions**; if  $V_k(r) = \mathbb{E} \|w_k - w^*\|^2$  decreases *geometrically*, then sequence of expected optimality gaps vanishes at *linear* rate.

Theorem 6 (Proof on the board)

Suppose that assumptions 1, 2 and 3, but with the existence of constants  $M \geq 0$  and  $\rho \in (0, 1)$  such that for all  $k \in \mathbb{N}$ :

$$V_k(r) = \mathbb{E} \|w_k - w^*\|^2 \leq M \rho^k \quad (5)$$

Moreover, suppose SGD is ran with fixed step size  $\alpha \in \min\{\frac{L}{2}, \frac{1}{c}\} \mu$ .

Then the expected optimality gap satisfies:

$$\mathbb{E} F(w_k) - F^* \leq \mu \rho^k \quad (6)$$

with  $\rho = \max\{\frac{LM}{c}, F(w_1) - F^*\}$ , and  $\mu = \max\{1 - \frac{c}{2}, \mu\}$ .



# Gradient Aggregation

**Main idea:** achieve a lower variance by *reusing* and or/ *revising* previously computed information.

**Key:** say that the current iterate has not been displaced too far from previous iterates,  
 $\tilde{N}$  stochastic gradient information from previous iterates may still be useful !

**Moreover:** say one maintains indexed gradient estimates in storage, then one can revise specific estimates as new information is collected.

**Attention:** we focus on the min. of the empirical risk.

# Stochastic variance-reduced gradient (SVRG)

## Main principle

operates in cycles.

At the beginning of each cycle: compute a batch gradient:

$$\bar{r} F(w_k) = \frac{1}{n} \sum_i^n r_i f_i(w_k)$$

Initialize:  $\tilde{w}_1 \leftarrow w_k$

Perform a set of  $m$  inner iterations indexed by  $j$  as follows:

$$\tilde{w}_{j+1} \leftarrow \tilde{w}_j - \eta (f_{i_j}(\tilde{w}_j) - \bar{r} f_{i_j}(w_k) + \bar{r} F(w_k))$$

with  $i_j \in \{1, \dots, n\}$  chosen random.

## SVRG

**Initialization:** Set  $w_1 \in W$ , stepsize  $\eta > 0$  and positive integer  $m$ .

Iteration ( $k \geq 1$ ):

Compute the batch gradient  $r = \nabla F(w_k)$

Initialize  $\tilde{w}_1 \in W_k$

**for**  $j = 1 : m$

Choose  $i_j \in \{1, \dots, n\}$

$\tilde{w}_{j+1} \in W_j \quad r_{f_{i_j}}(\tilde{w}_j) - r_{f_{i_j}}(w_k) + r = \nabla F(w_k)$

**end for**

Option (a):  $w_{k+1} = \tilde{w}_{m+1}$

Option (b):  $w_{k+1} = \frac{1}{m} \sum_{j=1}^m \tilde{w}_{j+1}$

Option (c): Choose  $j \in \{1, \dots, m\}$ , and  $w_{k+1} = \tilde{w}_{j+1}$

## SVRG

**Insights on step 3**, for simplicity pose

$$\tilde{g}_j : \quad \mathbb{E}_{i_j \sim \mathcal{P}} \left[ \nabla f_{i_j}(\tilde{w}_j) - \nabla f_{i_j}(w_k) \right] = \nabla F(w_k)$$

$$\mathbb{E}_{i_j \sim \mathcal{P}_{1, \dots, n}} \left[ \nabla f_{i_j}(w_k) - \nabla F(w_k) \right] = 0$$

$\nabla f_{i_j}(w_k) - \nabla F(w_k)$  seen as the bias in the gradient estimate  $\nabla f_{i_j}(w_k)$  in every iteration: SVRG draws a stochastic gradient  $\nabla f_{i_j}(\tilde{w}_j)$  evaluated at current inner iterate  $\tilde{w}_j$  and correct it based on perceived bias

*Overall*: the stochastic vector  $\tilde{g}_j$  represents an unbiased estimator of  $\nabla F(\tilde{w}_j)$

*but* with a variance that one can expect to be smaller than if one were simply to chose  $\tilde{g}_j : \nabla f_{i_j}$  (as for SGD)

↯ "Variance reduced method"

## SVRG

## Comments on "Options"

For both Options (b) and (c) and for min. of strongly convex functions, SVRG achieves a linear rate of convergence.

**Formally:** if we choose  $m$  (number of inner loops) and  $\alpha$  such that:

$$\alpha = \frac{1}{1 + \frac{2L}{mc}} \quad (7)$$

then we obtain:

$$\mathbb{E} F(w_{k+1}) - F^* \leq \alpha \mathbb{E} F(w_k) - F^* \quad (8)$$

This result applies for the outer iterates  $w_k$  where step from  $w_k$  to  $w_{k+1}$  requires  $2m + n$  evaluations of the gradients.

**Therefore:** one iteration of SVRG is much more expensive than one of SGD, and in fact is comparable to a full gradient iteration.

# SAGA

## Main principle

iteration that is closer in form to SGD: does not operate in cycles, nor does it compute batch gradients (except possibly at the initial point).

**At each it.:** computes a stochastic vector  $g_k$  as the average of stochastic gradients evaluated at previous iterates.

**Specifically:** at it.  $k$ , method will have stored  $f_i(w_{r_{is}})$  for all  $i \in \{1, \dots, n\}$ , where  $w_{r_{is}}$  represents the latest iterate at which  $f_i$  was evaluated.

**Then:** choose  $j \in \{1, \dots, n\}$ , the stochastic vector is set to:

$$g_k = f_j(w_k) - f_j(w_{r_{js}}) + \frac{1}{n} \sum_{i=1}^n f_i(w_{r_{is}}) \quad (9)$$

# SAGA

Initialization : Set  $w_1 = P W$ , stepsize  $\mu > 0$ .

```

for i = 1 : n
    Computer  $f_i(p w_1)$ 
    Store  $r_i = f_i(p w_{r_i}) - f_i(p w_1)$ 
end for

```

Choose  $j = \text{Up}(1; \dots; n)$

Computer  $f_j(p w_k)$

Set  $g_k = f_j(p w_k) - r_j = f_j(p w_{r_j}) - \frac{1}{n} \sum_{i=1}^n r_i$

Store  $r_j = f_j(p w_{r_j}) - f_j(p w_k)$

Set  $w_{k+1} = w_k + \mu g_k$

# SAGA

Insights:

Taking the expectation of  $g_k$  with respect to all choices of  $j \in \{1, \dots, n\}$ :  $E g_k = -\nabla F(w_k)$

Hence: method employs unbiased gradient estimates, but with variances less than the stochastic gradients employed by SGD

Putting aside the "storage" phase: the per-iteration cost of SAGA is the **same** as in a basic SGD method.

However: SAGA enjoys linear rate of convergence for min. of strongly convex  $F(w)$  (= empirical risk).

Result: with  $\frac{1}{2\mu n L \eta}$ :

$$E \|w_k - w^*\|^2 \leq \frac{1}{2\mu n L \eta} \sum_{i=1}^k \sum_{j=1}^n \sigma_{ij}^2 \|w_j - w^*\|^2 \leq C \eta^k \quad (10)$$

for some constant  $C$  (not detailed here).



# SAGA

## Insights:

Tips: previous result requires knowledge of the strong convexity constant  $c$  and  $L$ . If  $c$  is not known, then the stepsize can instead be chosen to be  $\frac{1}{3L}$  and a similar convergence result can be established.

important drawback : need to store  $n$  stochastic gradient vectors, which would be prohibitive in many large-scale applications.

History :SAGA has its origins in the stochastic average gradient (SAG) algorithm, the main difference:

$$g_k \in \frac{1}{n} \sum_{j=1}^n \nabla f_j(w_k) - \frac{1}{n} \sum_{j=1}^n \nabla f_j(w_{r_j}) + \sum_{i=1}^n \nabla f_i(w_{r_i})$$

# Gradient Aggregation - commentary

Even if linear rates have been proved (for strongly convex case), Gradient Aggregation methods should not be considered as superior to SGD !

One can prove that the computing time:

for SGD:  $T \propto n; \sigma \propto \frac{1}{c} \sigma^2$ ,

for SVRG, SAGA and SAG:  $T \propto n; \sigma \propto n \frac{1}{c} \sigma \log p \{ \sigma$

$\tilde{N}$  increase with!

Hence: for large  $n$ , gradient aggregation methods are comparable to batch methods, and cannot beat SGD.

Moreover: for  $L \leq 1$ , SGD is optimal.

but for  $L \leq n$ , gradient aggregation methods may be superior.

Anyway: interesting methods due to clever use of past info.

## Iterate Averaging Methods

Starting point : SGD generates noisy iterate sequences that tend to oscillate around minimizers.

Natural idea : compute a corresponding sequence of iterate averages (automatically possess less noisy behavior).

Algorithmic principle : employ the iteration

$$w_{k+1} = w_k + \alpha_k r_k + \beta_k (w_k - w_{k-1})$$

$$w_{k+1} = \frac{1}{k+1} \sum_{j=1}^k w_j \quad (11)$$

Hope: the auxiliary sequence  $w_{k+1}$  possesses better convergence properties than SGD iterates.

But.. with diminishing step size sequence with rate  $\mathcal{O}(1/k)$ , not the case !!

# Iterate Averaging Methods

**Fundamental advancement** - the work of (Polyak, 1991): rate for diminishing step size should be  $O_p\left(\frac{1}{k^a}\right)$  with a  $P \geq \frac{1}{2}$ ;  $1 \leq a \leq 2$

See works of (Juditsky et al., 1992), (Ruppert, 1988) and (Nemirovskiy et al., 1978) for more details.

Results:

$$\mathbb{E} \|w_k - w\|^2 \leq O_p\left(\frac{1}{k^a}\right) \quad (12)$$

Fruitful consequences: this idea of sequence of averaged iterates has been fundamental in many central and even recent works.

Main goal: allow longer steps while maintaining desired rates of convergence.

## So many things we skipped :)

More stochastic methods:

Stochastic quasi-Newton methods: estimate the inverse of the Hessian  $B_k$  and use the step:

$$w_{k+1} = w_k - B_k^{-1} \nabla f(w_k); \quad (13)$$

Adaptive step-sizes: step-size sequence adapted to the problem at stake, most famous methods are Adagrad and Adam.

Coordinate descent methods: update a subset of variables at a time

Accelerated SGD; see, e.g., Ghadimi and Lan (2016), Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming*, 156(1-2), 59-99.

# Table of Contents

- 1 Proximal stochastic gradient
- 2 Noise Reduction Methods
- 3 Numerical tests - hard-coded NN**

# Numerical tests

## Goals:

implement and benchmark some methods presented in this course,  
on a small and hand-able NN: build the so-called "back-propagation"  
along with the selected stochastic method

# Test setup

Context: consider the set of points: labeled data|some points are in class 1, indicated by red circles, and the rest are in class 2, indicated by blue crosses.



# Test setup

Our job: construct a prediction function  $h; w; q$  that takes any point in  $x_i \in \mathbb{R}^3$  and returns either a circle or a cross.

Form of  $h$ : before training  $h$  to do this task, we need to fix its form: we consider the network presented in the figure below:

# Forward Propagation

If new for you: here-under a very simple and step by step explanation:

Layer 1: represented by three circles, since our input vector have three components.

Layer 2: four solid circles, indicating that four neurons are being employed.

First connections: arrows from layer 1 to layer 2 indicate that all the three components of the input data are made available to the four neurons in layer 2.

Weight and bias: input data has the form  $x_i \in \mathbb{R}^3$ : weights and biases for layer 2 may be represented by a matrix  $W \in \mathbb{R}^{4 \times 3}$  and a vector  $b \in \mathbb{R}^4$  resp.

Layer output: output from layer 2 has the form:

$$p = Wx + b \quad (14)$$

with  $\sigma$  the so-called activation function of the layer at hand, here we choose the sigmoid function  $\sigma(x) = \frac{1}{1 + e^{-x}}$  applied component-wise.  $(\sigma(p_1), \sigma(p_2), \sigma(p_3), \sigma(p_4))$



# Forward Propagation - loss function

In this numerical test, we have access to realizations  $x_i, y_i$  with  $i = 1, \dots, n$ .

Now: it is time to train the network, that is finding values for parameters/variables  $w$  such that  $h : \mathbb{R}^3 \rightarrow \mathbb{R}^2$  such that  $h(x_i; w) = y_i$  for most  $i$  ( $1 \leq i \leq n$ ).

Loss function: with the notations of Part 1, we consider the loss function  $L(w) = \frac{1}{2} \sum_{i=1}^n \|h(x_i; w) - y_i\|_2^2$

Optimization Problem : we want to min. the empirical risk function, that is we want to solve :

$$\min_w F(w) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|h(x_i; w) - y_i\|_2^2 \quad (17)$$

# Forward Propagation - generalization

General notation : for a general network

The network has  $L$  layers, with layers 1 and  $L$  being the input and output layers,

Suppose that layer  $l$ , for  $l = 1; 2; 3; \dots; L$  contains  $n_l$  neurons;  $n_1$  is the dimension of the input data (we used  $d_x$  in Introduction of Part 1).

Overall: the network maps from  $\mathbb{R}^{n_1}$  to  $\mathbb{R}^{n_L}$ .

Weight and biases: use  $W^{l,s} \in \mathbb{R}^{n_l \times n_{l-1}}$  (weight matrix at layer  $l$ ) and  $b^{l,s} \in \mathbb{R}^{n_l}$  the vector of biases for layer  $l$  (neuron  $j$  at layer  $l$  uses the bias  $b_j^{l,s}$ ).

Activation : given an input  $x_i \in \mathbb{R}^{n_1}$ , we summarize the action of the network by letting  $a_j^{l,s}$  denote the output, or activation, from neuron  $j$  at layer  $l$ :

$$a_j^{l,s} = \sum_i x_i W_{ij}^{l,s} + b_j^{l,s} \quad \text{for } l = 2; 3; \dots; L \quad (18)$$

# Back Propagation

To apply our stochastic methods to min. Problem (17) and then train the network, we need compute the partial derivatives (the gradients for each parameter) of the objective function w.r.t each  $W_{jk}^{r/l}$  and  $b_j^{r/l}$ .

For a selected and *xed* input point  $x_i$ , denote  $a^{r/Ls} : hp x_i; wq$ .

To ease the expressions of the partial derivatives, we introduce:

$$z^{r/l} = W^{r/l} a^{r/l} + b^{r/l} \in \mathbb{R}^{n_l} \quad \text{for } l = 2; 3; \dots; L \quad (19)$$

## Fundamental induced relations:

$$a^{r/l} = \Phi(z^{r/l}),$$

let  $\frac{\partial C}{\partial z_j^{r/l}}$  for  $1 \leq j \leq n_l$  and  $2 \leq l \leq L$ . (measures the sensitivity of the loss function to the weighted input for neuron  $j$  at layer  $l$ ), with  $C = f(p; w; r; s; q) = \frac{1}{2} \|y_i - a^{r/Ls}\|_2^2$

# Back Propagation - Partial derivatives

**Useful operator:** the Hadamard product of two vectors  $u; v \in \mathbb{R}^n$ , that is  $(u \odot v)_i = u_i v_i$ .

We obtain the following results that are a consequence of the *chain rule*:

## Lemma 3

We have:

$$\frac{\partial \Phi^1}{\partial p^{rLs}} = \frac{\partial \Phi^1}{\partial p^{rLs}} \odot p^{rLs} y_j$$

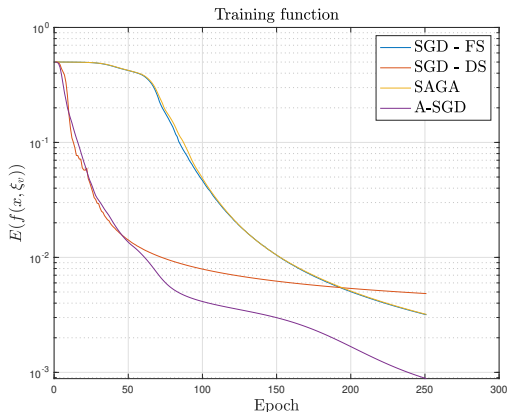
$$\frac{\partial \Phi^1}{\partial p^{r/s}} = \frac{\partial \Phi^1}{\partial p^{r/s}} \odot p^{r/s} W^{r/l} \odot q^{r/l} \odot q^{r/l} \text{ for } 2 \leq l \leq L$$

$$\frac{\partial \text{BC}}{\partial b_j^{r/s}} = \frac{\partial \text{BC}}{\partial b_j^{r/s}} \text{ for } 2 \leq l \leq L,$$

$$\frac{\partial \text{BC}}{\partial W_{jk}^{r/s}} = \frac{\partial \text{BC}}{\partial W_{jk}^{r/s}} a_k^{r/l} \text{ for } 2 \leq l \leq L.$$

# Demo - SGD (FS+DS) vs SAGA vs A-SGD

**Test case:**  $n = 1000$  (split into 80-20 % training-validation points),  
 $m = 10$  (size mini-batches), seed = 5000



Code [▶ Link](#), Code Ocean [▶ Link](#)



# Demo - SGD (FS+DS) vs SAGA vs A-SGD

Test case:  $n = 1000$  (split into 80-20 % training-validation points),  
 $m = 10$  (size mini-batches), seed = 5000

Code [▶ Link](#), Code Ocean [▶ Link](#)

# Demo - SGD (FS+DS) vs SAGA vs A-SGD

Test case:  $n = 1000$  (split into 80-20 % training-validation points),  
 $m = 10$  (size mini-batches), seed = 5000

Code [▶ Link](#), Code Ocean [▶ Link](#)

# References

Boyd, Parikh, Chu, Peleato and Eckstein (2010)

Distributed Optimization and Statistical Learning via the Alternating Method of Multipliers

Machine Learning Vol. 3, No. 1 (2010) 1{122

Y.Nesterov. (2018)

Lectures on Convex Optimization.

Springer.

G. E. Hinton. (2012)

A practical guide to training restricted Boltzmann machines.

Neural Networks: Tricks of the Trade, volume 7700 of Lecture Notes in Computer Science pages 599{619. Springer, 2012.

# References

N. Gillis (2021)

First-order methods for large-scale optimization

*Lectures UMon - Advanced Opt. course*

Bottou, Curtis, and Nocedal (2018)

Optimization Methods for Large-Scale Machine Learning.

*SIAM REVIEW* Vol. 60, No. 2, pp. 223{311, 2018.

G. Papamakarios (2014)

Comparison of Modern Stochastic Optimization Algorithms.

*Technical report, University of Edinburgh*

N. Keskar et al. (2017)

On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima.

*arXiv:1609.04836*

# References

A. Agarwal et al. (2012)

Information- theoretic lower bounds on the oracle complexity of stochastic convex optimization

*IEEE Trans. Inform. Theory*, 58, pp. 3235-3249

B.T. Polyak. (1991)

New method of stochastic approximation type

*Automat. Remote Control*, 51, pp. 937-946

B.T. Polyak and A.B. Juditsky (1992)

Acceleration of stochastic approximation by averaging

*SIAM J. Control Optim.*, 30, pp. 838-855

D. Ruppert (1988)

Efficient Estimations from a Slowly Convergent Robbins-Monro Process

*Tech report*, 781.

# References

A.S. Nemirovski and D.B. Yudin (1978)

On Cesaro's convergence of the steepest descent method for approximating saddle point of convex-concave functions

*Soviet Math. Dokl.*, 19

C.F. Higham and D.J. Higham (2019)

Deep Learning: An Introduction for Applied Mathematicians

*SIAM REVIEW* Vol. 61, No. 4, pp. 860{891, 2019.

