# Stochastic and randomized convex optimization

(Incomplete version without transitions)

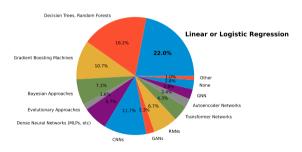
Adrien Taylor

INRIA & École Normale Supérieure

# **Program** for today

- ♦ Convex stochastic optimization,
- batch gradient methods,
- stochastic gradient descent,
- ♦ finite-sum algorithms,
- (randomized) coordinate methods,
- ... on a few running examples.

#### Which of the following ML algorithms do you use on a regular basis?



See Kaggle survey 2022.

#### Table of contents

- 1. Stochastic optimization problems
- 2. Plain gradient methods
- 3. Stochastic gradient methods
- 4. Finite sums
- 5. Popular stochastic algorithms
- 6. Randomized coordinate descent
- 7. Conclusion

# Stochastic optimization problems

#### Table of contents

- 1. Stochastic optimization problems
- 2. Plain gradient methods
- 3. Stochastic gradient methods
- 4. Finite sums
- 5. Popular stochastic algorithms
- 6. Randomized coordinate descent
- 7. Conclusion

- $\diamond$  Input measurement  $x \in \mathcal{X}$ ,
- $\diamond$  output measurement  $y \in \mathcal{Y}$ ,
- $\diamond$   $(x, y) \sim \mathcal{D}$  with  $\mathcal{D}$  unknown,
- $\diamond$  training data:  $\mathcal{D}_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$  (i.i.d.  $\sim \mathcal{D}$ ).

#### Often:

- $x \in \mathbb{R}^d$  and  $y \in \{-1,1\}$  (classification),
- or  $x \in \mathbb{R}^d$  and  $y \in \mathbb{R}$  (regression).

We search a predictor function  $p: \mathcal{X} \to \mathcal{Y}$ .



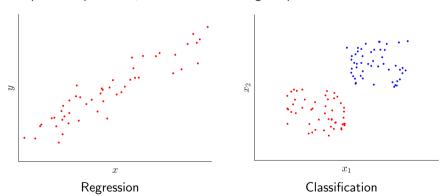
Target: find  $p: \mathcal{X} \to \mathcal{Y}$ 

$$hoigg(egin{pmatrix} 
hoigg) & 
ightarrow & 1 \ 
hoigg(igg) & 
ightarrow & -1 \ \end{pmatrix}$$

#### Often:

- $x \in \mathbb{R}^d$  and  $y \in \{-1, 1\}$  (classification),
- or  $x \in \mathbb{R}^d$  and  $y \in \mathbb{R}$  (regression).

We search a predictor  $p: \mathcal{X} \to \mathcal{Y}$ . How to construct good predictors?



How to construct a good predictor?

- $\diamond$  Pick a **loss function**:  $\ell(p(x), y)$  to measure quality of  $p(x) \approx y$ .
- ♦ Examples:
  - $0 1 \text{ loss: } \ell(p(x), y) = \mathbf{1}_{y \neq f(x)},$
  - quadratic loss:  $\ell(p(x), y) = |p(x) y|^2$ .

#### Risk function

 $\diamond$  Risk measures the average loss over  ${\cal D}$ 

$$\mathcal{R}(p) = \mathbb{E}_{(x,y)\sim\mathcal{D}} \left[\ell(p(x),y)\right].$$

- ♦ Examples:
  - -0-1 risk:  $\mathcal{R}(p) = \mathbb{P}(y \neq p(x))$ .
  - Quadratic risk:  $\mathcal{R}(p) = \mathbb{E}\left[|y p(x)|^2\right]$ .

Learning a predictor via decision variable  $\theta$ 

$$\underset{\theta \in \mathbb{R}^d}{\text{minimize}} \ \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(p_{\theta}(x),y)].$$

Here:  $\mathcal{D}$  is distribution of datapoints  $\xi = (x, y) \in \mathbb{R}^{d+1}$ , and linear  $p_{\theta}(x) = \langle \theta, x \rangle$ . Examples:

- $\diamond$  linear regression:  $\ell(p_{\theta}(x), y) = (\langle \theta, x \rangle y)^2$ ,
- $\diamond$  support vector machines:  $\ell(p_{\theta}(x), y) = \max\{0, 1 y\langle \theta, x \rangle\}.$

For all of those beyond pure linear models: see kernel versions.

Name	$\ell(y_p,y)$	Graph $\ell(y_p,1)$
0-1 loss	$\ell(y_p, y) = \left\{ egin{array}{ll} 0 &  ext{if } y_p = y \ 1 &  ext{if } y_p  eq y \end{array}  ight.$	<u> </u>
quadratic loss	$\ell(y_p,y)=(y_p-y)^2$	
logistic loss	$\ell(y_p,y) = \log\left(1 + \exp(-y_p y)\right)$	<del></del>
hinge loss	$\ell(y_\rho,y)=max\{0,1-y_\rho y\}$	

# Stochastic optimization framework

Learning a predictor via decision variable  $\theta$ 

$$\underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\text{minimize}} \ \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(p_{\boldsymbol{\theta}}(x),y)] \triangleq \mathbb{E}_{\xi \sim \mathcal{D}}\left[f(\underline{\boldsymbol{\theta}};\xi)\right].$$

with  $\mathcal{D}$ : distribution of datapoints  $\xi = (x, y) \in \mathbb{R}^{d+1}$ .

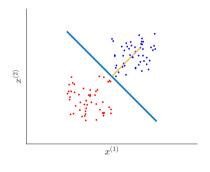
Often approached via empirical risk minimization:

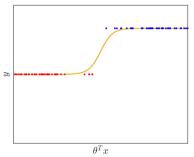
$$\underset{\theta \in \mathbb{R}^d}{\mathsf{minimize}} \ \frac{1}{n} \sum_{i=1}^n \ell(\langle \theta, x_i \rangle, y_i) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(\theta) \triangleq F(\theta).$$

# Classification via logistic regression

We have  $\mathcal{D}_n = \{(x_1, y_i), i = 1, \dots, n\}$ , with  $y_i \in \{-1, 1\}$ .

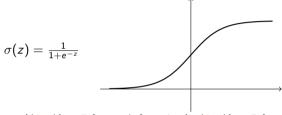
Objective: find  $\theta$  such that  $y_i\langle \theta, x_i \rangle \geqslant 0$  for all  $i = 1, \ldots, n$ .





# Classification via logistic regression

♦ Pick sigmoid function



- $\diamond$  interpret:  $\sigma(\langle \theta, x \rangle) = \mathbb{P}\{y = 1 | x\}$  and  $\sigma(-\langle \theta, x \rangle) = \mathbb{P}\{y = -1 | x\}$
- $\diamond$  with maximum likelihood / cross-entropy loss, yields logistic regression

$$\underset{\theta \in \mathbb{R}^d}{\mathsf{minimize}} \, \frac{1}{n} \sum_{i=1}^n \log \left( 1 + \exp \left( -y_i \langle \theta, x_i \rangle \right) \right).$$

Convex! (How to show that?)

#### Table of contents

- 1. Stochastic optimization problems
- 2. Plain gradient methods
- 3. Stochastic gradient methods
- 4. Finite sums
- 5. Popular stochastic algorithms
- 6. Randomized coordinate descent
- 7. Conclusion

# Plain gradient methods

# **Stochastic optimization**

Empirical risk minimization as

$$\underset{\theta \in \mathbb{R}^d}{\mathsf{minimize}} \left\{ F(\theta) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(\theta) \right\}.$$

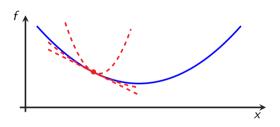
Starting assumptions (we will use variations around this):

- $\diamond$  each  $f_i(\cdot)$  has a Lipschitz gradient (constant L),
- $\diamond$  each  $f_i(\cdot)$  is strongly convex (constant  $\mu$ ).

When  $f_i$  twice continuously differentiable:  $\mu I_d \preccurlyeq \nabla^2 f_i(\theta) \preccurlyeq L I_d$  for all  $\theta \in \text{dom } f_i$ .

# **About the assumptions**

A differentiable function  $f: \mathbb{R}^d \to \mathbb{R}$  is  $\mu$ -strongly convex and L-smooth iff  $\forall x, y \in \mathbb{R}^d$ :



(1) (Convexity) 
$$f(x) \ge f(y) + \langle \nabla f(y), x - y \rangle$$
.

(1b) (
$$\mu$$
-strong convexity)  $f(x) \ge f(y) + \langle \nabla f(y), x - y \rangle + \frac{\mu}{2} ||x - y||_2^2$ ,

(2) (L-smoothness) 
$$f(x) \leqslant f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} ||x - y||_2^2$$
,

$$(1\&2) \ f(x) \geqslant f(y) + \langle \nabla f(y), x - y \rangle + \frac{1}{2L} \|\nabla f(x) - \nabla f(y)\|_2^2 + \frac{\mu}{2(1 - \mu/L)} \|x - y - \frac{1}{L} (\nabla f(x) - \nabla f(y))\|_2^2,$$

(1&2b) 
$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \ge \frac{1}{L+\mu} \|\nabla f(x) - \nabla f(y)\|_2^2 + \frac{\mu L}{L+\mu} \|x - y\|_2^2$$
.

# **About the assumptions**

First-order optimization: condition number  $\kappa=\frac{L}{\mu}\geqslant 1$  discriminates "easy" vs. "hard".

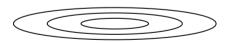
- ♦ Smoothness *L* given by curvature in direction with fastest variation,
- $\diamond$  Strong convexity  $\mu$  given by curvature in direction with slowest variation.

Insights from level curves:

very well conditioned problem (  $\kappa\approx 1)$  :



more poorly conditioned one (  $\kappa\gg 1$  ):



## **Examples**

- ♦ Regularized least squares (Ridge regression):  $f_i(\theta) = (\langle \theta, x_i \rangle y_i)^2 + \frac{\lambda}{2} \|\theta\|_2^2$ . Hessian:  $\nabla^2 f_i(\theta) = 2x_i x_i^T + \lambda I_d$ . Hence:  $L = 2 \max_{1 \leqslant i \leqslant n} \|x_i\|_2^2 + \lambda$  and  $\mu = \lambda$ .
- $\diamond$  Regularized logistic regression:  $f_i(\theta) = \log(1 + \exp(-y_i\langle \theta, x_i \rangle)) + \frac{\lambda}{2} \|\theta\|_2^2$ , we have:

$$\nabla f_i(\theta) = \frac{-y_i x_i}{1 + \exp(y_i \langle \theta, x_i \rangle)} + \lambda \theta, \quad \nabla^2 f_i(\theta) = \frac{\exp(y_i \langle \theta, x_i \rangle)}{(1 + \exp(y_i \langle \theta, x_i \rangle))^2} x_i x_i^T + \lambda I_d.$$

Therefore, for any  $z \in \mathbb{R}^d$ : (hint: use  $\frac{e^u}{(1+e^u)^2} \leqslant \frac{1}{4}$ )

$$z^{T} \nabla^{2} f_{i}(\theta) z = z^{T} x_{i} x_{i}^{T} z \frac{\exp(y_{i}\langle\theta,x_{i}\rangle)}{(1+\exp(y_{i}\langle\theta,x_{i}\rangle))^{2}} + \lambda I_{d} ||z||_{2}^{2} \leqslant z^{T} \left(\frac{1}{4} x_{i} x_{i}^{T} + \lambda I_{d}\right) z$$

Hence 
$$L = \frac{1}{4} \max_{1 \leq i \leq n} ||x_i||_2^2 + \lambda$$
 and  $\mu = \lambda$ .

# Plain gradient descent

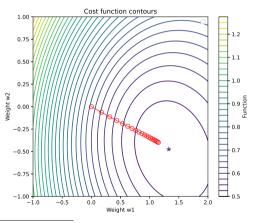
Algorithm: Plain gradient descent Set 
$$\theta^0 \in \mathbb{R}^d$$
,  $\alpha > 0$ . for  $t = 0, 1, \ldots$  do  $\theta^{t+1} = \theta^t - \frac{\alpha}{n} \sum_{i=1}^n \nabla f_i(\theta^t)$  end

In this context, for  $\alpha = \frac{1}{L}$ :

$$\begin{split} F(\theta^t) - F(\theta^\star) &\leqslant \min\left\{\tfrac{1}{t}, \left(1 - \tfrac{1}{\kappa}\right)^t\right\} \frac{L\|\theta^0 - \theta^\star\|_2^2}{2}. \\ \|\theta^t - \theta^\star\|_2^2 &\leqslant \left(1 - \tfrac{1}{\kappa}\right)^t\|\theta^0 - \theta^\star\|_2^2. \end{split}$$

### Plain GD

Gradient descent  $(\alpha = \frac{1}{L})$ :<sup>1</sup>



<sup>&</sup>lt;sup>1</sup>Logistic regression problem: "fourclass" dataset from LIBSVM (n, d) = (862, 2).

# Classical GD convergence analysis

General idea: studying a single iteration is simpler. Need recursable bounds.

One can prove  $V^{t+1} \leqslant V^t$  for all  $\theta^t$ ,  $\theta^{t+1} = \theta^t - \frac{1}{L} \nabla F(\theta^t)$  and L-smooth convex function, with

$$V^t \triangleq V(A_t, \theta^t) \triangleq A_t(F(\theta^t) - F(\theta^\star)) + \frac{L}{2} \|\theta^t - \theta^\star\|_2^2$$

as soon as  $A_{t+1} \leqslant A_t + 1$  (with  $A_t \geqslant 0$ ). Valid choice:  $A_t = t$ .

Why is this nice?

$$A_t (F(\theta^t) - F(\theta^*)) \leqslant V^t \leqslant V^{t-1} \leqslant \ldots \leqslant V^0,$$

so  $F(\theta^t) - F(\theta^\star) \leqslant \frac{V^0}{A_t} = \frac{L\|\theta^0 - \theta^\star\|_2^2}{2A_t}$  when choosing  $A_0 = 0$ . Choice  $A_t = t$  gives:

$$F(\theta^t) - F(\theta^\star) \leqslant \frac{L\|\theta^0 - \theta^\star\|_2^2}{2t}.$$

# GD: recall convergence analysis — a simple case

For GD, proof example of a simple bound:

$$\begin{split} \|\theta^{t+1} - \theta^\star\|_2^2 &= \|\theta^t - \theta^\star\|_2^2 - 2\alpha \langle \nabla F(\theta^t), \theta^t - \theta^\star \rangle + \alpha^2 \|\nabla F(\theta^t)\|_2^2 \\ & \qquad \qquad \Big| \text{ Inequality (1\&2b)} \\ & \leqslant \left(1 - \frac{2\alpha L\mu}{L+\mu}\right) \|\theta^t - \theta^\star\|_2^2 + \alpha \left(\alpha - \frac{2}{L+\mu}\right) \|\nabla F(\theta^t)\|_2^2 \\ & \qquad \qquad \Big| \text{ if } 0 \leqslant \alpha \leqslant \frac{2}{L+\mu} \\ & \leqslant (1 - \alpha\mu)^2 \|\theta^t - \theta^\star\|_2^2. \end{split}$$

# Regularized linear approximations and momentum

Recall gradient descent  $\theta^{t+1} = \theta^t - \alpha \nabla F(\theta^t)$ . Equivalently:

$$\theta^{t+1} = \operatorname*{argmin}_{\theta \in \mathbb{R}^d} \left\{ \left| F(\theta^t) + \langle \nabla F(\theta^t), \theta - \theta^t \rangle \right| + \frac{2}{\alpha} \|\theta - \theta^t\|_2^2 \right\}$$

essentially: regularized linear approximation.

Similar template for accelerated gradient descent (iterates  $\{(\theta^t, \phi^t, \lambda^t)\}_{t=0,1,...}$ 

$$\phi^{t} = (1 - \tau_{t}) \theta^{t} + \tau_{t} \lambda^{t}$$

$$\lambda^{t+1} = \operatorname*{argmin}_{\lambda \in \mathbb{R}^{d}} \left\{ \sum_{i=0}^{t} \left[ (A_{i+1} - A_{i}) \left( F(\phi^{i}) + \langle \nabla F(\phi^{i}), \theta - \theta^{i} \rangle \right) \right] + \frac{2}{\alpha} \|\lambda - \phi^{t}\|_{2}^{2} \right\}$$

$$\theta^{t+1} = (1 - \tilde{\tau}_{t}) \theta^{t} + \tilde{\tau}_{t} \lambda^{t+1}$$

... similarly: based on regularized (weighted) linear approximations of  $F(\cdot)$ 

(with growing sequence  $\{A_t\}_{t=0,1,...}$  and some  $\{(\tau_k, \tilde{\tau}_k)\}_{t=0,1,...}$  for convex combinations).

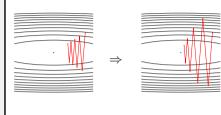
# Plain accelerated gradient descent

#### Algorithm: Plain acceleration for ERM

Set 
$$\theta^0 = \tilde{\theta}^0 \in \mathbb{R}^d$$
,  $\alpha, \{\beta_t\} > 0$ .  
for  $t = 0, 1, \dots, T - 1$  do

$$\theta^{t+1} = \tilde{\theta}^t - \frac{\alpha}{n} \sum_{i=1}^n \nabla f_i(\tilde{\theta}^t)$$
  
$$\tilde{\theta}^{t+1} = \theta^{t+1} + \beta_t(\theta^{t+1} - \theta^t)$$

end



In this context, for appropriate choices of  $(\alpha, \beta)$ : (for some C > 0)

$$\begin{split} F(\theta^t) - F(\theta^\star) &\leqslant \min\left\{\frac{2}{t^2}, \left(1 - \sqrt{\frac{1}{\kappa}}\right)^t\right\} L \|\theta^0 - \theta^\star\|_2^2. \\ \|\theta^t - \theta^\star\|_2^2 &\leqslant C\left(1 - \sqrt{\frac{1}{\kappa}}\right)^t \|\theta^0 - \theta^\star\|_2^2, \end{split}$$

using similar proof patterns.<sup>2</sup>

<sup>&</sup>lt;sup>2</sup>See, e.g., d'Aspremont, Scieur, T (2021). "Acceleration methods."

# Classical AGD convergence analysis

General idea: studying a single iteration is simpler. Need recursable bounds.

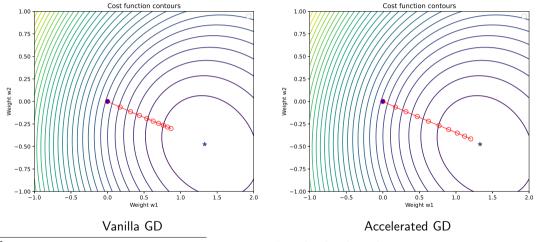
One can prove  $V^{t+1} \leqslant V^t$  with

$$V^t \triangleq V(A_t, \theta^t, \lambda^t) \triangleq A_t(F(\theta^t) - F(\theta^*)) + \frac{L}{2} ||\lambda^t - \theta^*||_2^2$$

and  $A_t \approx t^2$  when  $A_0 = 0$ .

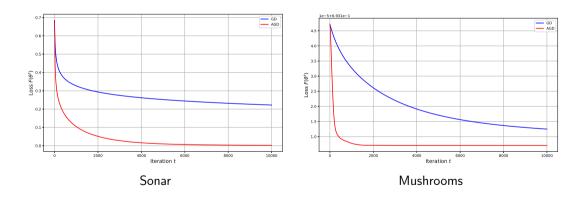
In short: all coefficient choices made for greedily making  $A_t$  large.

### GD vs. AGD<sup>3</sup>



 $<sup>^3</sup>$ Logistic regression problem: "fourclass" dataset from LIBSVM (n, d) = (862, 2).

#### GD vs. AGD4



 $<sup>^4</sup>$ Logistic regression: "sonar" (n,d)=(208,60) and "mushrooms" (n,d)=(8124,112) datasets from LIBSVM.

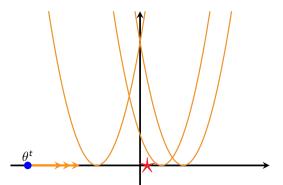
# Plain gradients for ERM - takeaways

Were we exploiting what we can?

- $\diamond$  Momentum?  $\rightarrow$  accelerated convergence rates.
- $\diamond$  Adaptative step-size selection?  $\to$  backtracking line-search, online estimation of L,...

But when far away from solution: single  $\nabla f_i(\theta^t)$  is already informative!

ightarrow useful to evaluate the full batch?



#### Table of contents

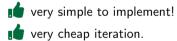
- 1. Stochastic optimization problems
- 2. Plain gradient methods
- 3. Stochastic gradient methods
- 4. Finite sums
- 5. Popular stochastic algorithms
- 6. Randomized coordinate descent
- 7. Conclusion

Stochastic gradient methods

# Stochastic gradient descent (SGD)

$$\underset{\theta \in \mathbb{R}^d}{\mathsf{minimize}} \left\{ F(\theta) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(\theta) \right\}.$$

 $\textbf{Algorithm:} \ \mathsf{SGD}, \ \mathsf{constant} \ \mathsf{step}\text{-}\mathsf{size}$  $\begin{aligned} &\mathsf{Set}\; \theta^0 \in \mathbb{R}^d,\; \alpha > 0 \\ &\mathsf{for}\; t = 0,1,\dots,T-1 \; \mathsf{do} \\ &\mid \mathsf{sample}\; i_t \sim \mathcal{U}[[1,n]] \\ &\mid \theta^{t+1} = \theta^t - \alpha \nabla f_{i_t}(\theta^t) \\ &\mathsf{end} \end{aligned}$ 





# Stochastic gradient descent (SGD)

#### Observations:

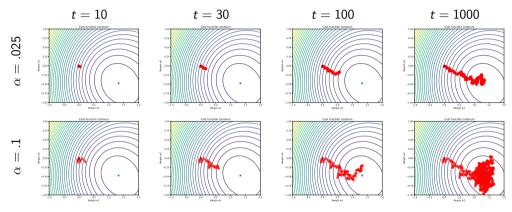
 $\diamond \ \nabla f_i(\theta)$  (with  $i \sim \mathcal{U}[[1,n]]$ ) is unbiased estimate of gradient (more later):

$$\mathbb{E}_i[\nabla f_i(\theta)] = \nabla F(\theta).$$

- $\diamond$  What if  $\nabla f_i(\theta)$ 's  $(i=1,\ldots,n)$  are very different? If very similar?
  - ightarrow variance of gradient estimation drives behavior of SGD!

# Stochastic gradient descent: empirical behavior

Short step sizes<sup>5</sup>

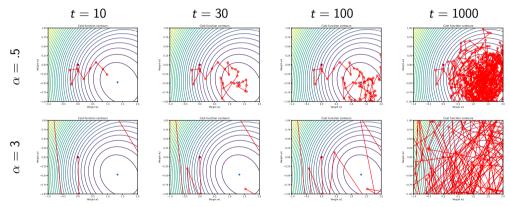


 $\rightarrow\,$  very slow to converge & relatively accurate.

<sup>&</sup>lt;sup>5</sup>Logistic regression problem: "fourclass" dataset from LIBSVM (n, d) = (862, 2).

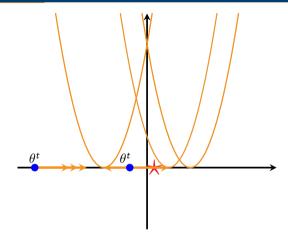
# Stochastic gradient descent: empirical behavior

## Larger step sizes



- ightarrow faster to reach "stationnary behavior" (forget about initial conditions) & not accurate.
- $\rightarrow$  we want: initially large  $\alpha,$  then short  $\alpha$  on the long run.

# Stochastic gradient descent: empirical behavior



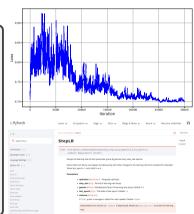
Morally, two extreme regimes:

- $\diamond$  "error due to initial conditions" dominates  $\rightarrow$  stochastic gradients are very informative
- $\diamond$  "error due to noise" dominates  $\rightarrow$  need to accomodate noise.

# Mitigating noise via step-size schedulers

Naive scheduler:6

```
Algorithm: SGD, naive step-size scheduler
Set \theta^0 \in \mathbb{R}^d, \alpha^0 > 0, c \in (0,1), K \in \mathbb{N}
for t = 0, 1, ..., T - 1 do
     sample i_t \sim \mathcal{U}[[1, n]]
     \theta^{t+1} = \theta^t - \alpha \nabla f_{i,t}(\theta^t)
     if mod(t+1, K) = 0 then
       \alpha = c\alpha
     end
end
```



<sup>&</sup>lt;sup>6</sup>Experiment with the "mushroom" dataset from LIBSVM (n, d) = (8124, 112).

# Mitigating noise via minibatches

$$\begin{aligned} &\textbf{Algorithm:} \text{ minibatch-SGD, constant step-size} \\ &\text{Set } \theta^0 \in \mathbb{R}^d, \ \alpha > 0, \ b \in \mathbb{N}. \\ &\textbf{for } t = 0, 1, \dots, T-1 \ \textbf{do} \\ & \left| \begin{array}{l} \text{sample } i_t^{(1)}, \dots, i_t^{(b)} \sim \mathcal{U}[[1, n]], \ \mathcal{I}_t = \{i_t^{(1)}, \dots, i_t^{(b)}\} \\ & \theta^{t+1} = \theta^t - \frac{\alpha}{b} \sum_{i \in \mathcal{I}_t} \nabla f_i(\theta^t) \end{array} \right. \end{aligned}$$
 end

Ь	name	gradient estimate	computational cost
1	(pure) SGD	$ abla f_{i_t}( heta^t)$ with $i_t \in \mathcal{U}[[1,n]]$	O(d)
1 < b < n	minibatch SGD	$\sum_{i \in \mathcal{I}_t}  abla f_i( heta^t), \  \mathcal{I}_t  = b$	O(bd)
b = n	full batch/plain GD	$\sum_{i=1}^{n} \nabla f_i(\theta^t)$	O(nd)

- $\diamond$  Commonly: pick  $b=2^a$  (a=5,6,...) to benefit from parallelization on GPU/CPU.
- $\diamond$  For theory, focus on b = 1.

# Stochastic gradient descent – unbiasedness

If batch chosen uniformly at random & independently from past  $\Rightarrow$  unbiased gradient estimate.

 $\diamond$  pure SGD: pick  $i_t \in \mathcal{U}[[1, n]]$  independently from past iterates then

$$\mathbb{E}\left[\nabla f_{i_t}(\theta^t)|\theta^t\right] = \frac{1}{n}\sum_{i=1}^n \nabla f_i(\theta^t) \triangleq \nabla F(\theta^t).$$

 $\diamond$  Minibatch SGD: pick  $\mathcal{I}_t$  uniformly at random in  $\{1,\ldots,n\}$  (with or without resampling) & independently from past iterates then

$$\mathbb{E}\left[\frac{1}{b}\sum_{i\in\mathcal{I}_t}\nabla f_i(\theta^t)\bigg|\theta^t\right] = \frac{1}{bn}\sum_{i=1}^b\sum_{j=1}^n\nabla f_j(\theta^t) = \nabla F(\theta^t).$$

unbiased but noisy estimations. Effect of b on variance?

# Stochastic gradient descent – bounds

$$\underset{\theta \in \mathbb{R}^d}{\mathsf{minimize}} \left\{ F(\theta) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(\theta) \right\}.$$

Classical assumptions (variations on this theme in what follows)

- $\diamond$  each  $f_i(\cdot)$  is L-smooth and  $\mu$ -strongly convex,
- $\diamond$  bounded variance at  $\theta^*$ :  $\mathbb{E}_i \left[ \|\nabla F(\theta^*) \nabla f_i(\theta^*)\|_2^2 \right] = \mathbb{E}_i \left[ \|\nabla f_i(\theta^*)\|_2^2 \right] \leqslant \sigma^2$ .

One can show: (with  $\alpha = 1/L$  for simplicity)

$$\mathbb{E}_{i}\left[\left\|\theta^{t+1}-\theta^{\star}\right\|_{2}^{2}\left|\theta^{t}\right]\leqslant\left(1-\frac{\mu}{L}\right)^{2}\left\|\theta^{t}-\theta^{\star}\right\|_{2}^{2}+\frac{2\sigma^{2}}{L^{2}}.$$

### Stochastic gradient descent – bounds

**Proof.** reformulate the inequality (due to smoothness and strong convexity), namely (1&2b):

$$\begin{split} 0 \geqslant \mathbb{E}_{i_t} \left[ -\langle \nabla f_{i_t}(\theta^t) - \nabla f_{i_t}(\theta^\star), \theta^t - \theta^\star \rangle + \frac{1}{L} \|\nabla f_{i_t}(\theta^t) - \nabla f_{i_t}(\theta^\star)\|_2^2 \right. \\ & + \frac{\mu}{1 - \mu/L} \|\theta^t - \theta^\star - \frac{1}{L} (\nabla f_{i_t}(\theta^t) - \nabla f_{i_t}(\theta^\star))\|_2^2 \right] \end{split}$$

multiplied by  $2\alpha(1-\alpha\mu)\geqslant 0$  (with  $0\leqslant \alpha\leqslant 1/L$ ) as

$$\begin{split} \mathbb{E}_{i_{t}}[\|\theta^{t+1} - \theta^{\star}\|_{2}^{2}] \leqslant & (1 - \alpha\mu)^{2} \|\theta^{t} - \theta^{\star}\|_{2}^{2} + \frac{2\alpha^{2}(1 - \alpha\mu)}{2 - \alpha(L + \mu)} \mathbb{E}_{i_{t}}[\|\nabla f_{i_{t}}(\theta^{\star})\|_{2}^{2}] \\ & - \frac{\alpha(2 - \alpha(L + \mu))}{L - \mu} \mathbb{E}_{i_{t}} \|\mu(\theta^{\star} - \theta^{t}) + \nabla f_{i_{t}}(\theta^{t}) + 2\frac{1 - \alpha\mu}{\alpha(L + \mu) - 2} \nabla f_{i_{t}}(\theta^{\star})\|_{2}^{2} \\ \leqslant & (1 - \alpha\mu)^{2} \|\theta^{t} - \theta^{\star}\|_{2}^{2} + \frac{2\alpha^{2}(1 - \alpha\mu)}{2 - \alpha(L + \mu)} \sigma^{2} \end{split}$$

(using unbiasedness:  $\mathbb{E}_{i_t}[\langle \nabla f_{i_t}(\theta^\star), \theta^t \rangle] = \mathbb{E}_{i_t}[\langle \nabla f_{i_t}(\theta^\star), \theta^\star \rangle] = 0$ ).

Desired result by evaluating  $\alpha \leftarrow \frac{1}{L}$ .

# Stochastic gradient descent – bounds

By chaining inequalities, we arrive to  $(t \ge 0)$ 

$$\mathbb{E}_{i} \left[ \| \theta^{t} - \theta^{\star} \|_{2}^{2} | \theta^{0} \right] \leq \left( 1 - \frac{\mu}{L} \right)^{2t} \| \theta^{0} - \theta^{\star} \|_{2}^{2} + \frac{2\sigma^{2}}{L^{2}} \sum_{i=0}^{t-1} \left( 1 - \frac{\mu}{L} \right)^{2i}$$

$$\leq \left( 1 - \frac{\mu}{L} \right)^{2t} \| \theta^{0} - \theta^{\star} \|_{2}^{2} + \frac{2\sigma^{2}}{L^{2}} \left( \frac{L}{\mu} - \frac{L}{\mu} \left( 1 - \frac{\mu}{L} \right)^{t} \right)$$

$$\leq \left( 1 - \frac{\mu}{L} \right)^{2t} \| \theta^{0} - \theta^{\star} \|_{2}^{2} + \frac{2\sigma^{2}}{L\mu} .$$

Hence, for SGD with constant  $\alpha = \frac{1}{L}$  reaches

$$\mathbb{E}_{i}\left[\|\theta^{t} - \theta^{\star}\|_{2}^{2} |\theta^{0}\right] \leqslant \left(1 - \frac{\mu}{L}\right)^{2t} \|\theta^{0} - \theta^{\star}\|_{2}^{2} + \frac{2\sigma^{2}}{L\mu}.$$

 $\rightarrow$  convergence to a ball around  $\theta^{\star}$ .

# Stochastic gradient descent – bounds & takeaways

#### Theory and experience agree on:

- small step-size: slowly forget initial condition; convergence to a small ball around solution.
- ♦ Large step-size: better forget initial conditions; convergence to a larger ball.

#### Can we do better?

- averaging,
- decreasing step-sizes (step-size schedules),
- decrease variance (minibatching).

Here: let's simplify the assumptions for this study.

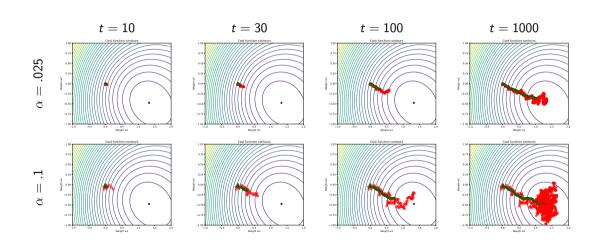
$$\underset{\theta \in \mathbb{R}^d}{\mathsf{minimize}} \left\{ F(\theta) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(\theta) \right\}.$$

Simplifying assumptions here:

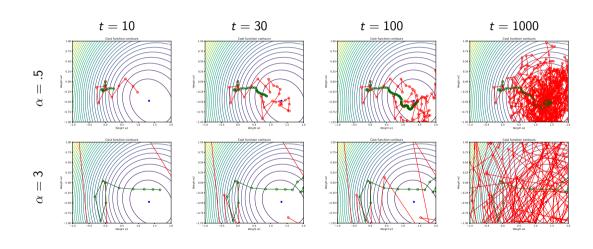
- $\diamond$  each  $f_i(\cdot)$  is convex
- $\diamond$  bounded stochastic gradients  $\mathbb{E}\left[\|\nabla f_i(\theta)\|_2^2\right] \leqslant M^2$ .

(one can get refined analyses using smoothness/strong convexity).

# **SGD**: averaging



# **SGD**: averaging



Suppose 
$$\|\theta^0 - \theta^\star\|_2 \leqslant R$$
 for some  $\theta^0 \in \mathbb{R}^d$ , and  $\mathbb{E}_i[\|\nabla f_i(\theta^t)\|_2^2] \leqslant M^2$ , then

$$\mathbb{E}[F(\bar{\theta}^T)] - F(\theta^*) \leqslant \frac{R^2 + M^2 \sum_{t=0}^T \alpha_t^2}{2 \sum_{t=0}^T \alpha_t},$$

$$\mathbb{E}[F(\bar{\theta}^T)] - F(\theta)$$
 with  $\bar{\theta}^T = \frac{1}{\sum_{t=0}^T \alpha_t} \sum_{t=0}^T \alpha_t \theta^t$  (averaging).

- $\diamond\,$  Proof essentially similar to that of the subgradient method.
- Rates are similar (but in expectation).

**Proof.** Define  $r_t = \|\theta^t - \theta^*\|_2$ , we have:

$$r_{t+1}^2 = r_t^2 - 2\alpha_t \langle \nabla f_{i_t}(\theta^t), \theta^t - \theta^\star \rangle + \alpha_t^2 \|\nabla f_{i_t}(\theta^t)\|_2^2.$$

Taking expectations and using convexity and indepence of  $i_t$  and  $\theta^t$ 

$$\begin{split} \mathbb{E}[r_{t+1}^2] &\leqslant \mathbb{E}[r_t^2] - 2\alpha_t \mathbb{E}\left[ \langle \nabla f_{i_t}(\theta^t), \theta^t - \theta^* \rangle \right] + \alpha_t^2 \mathbb{E}[\|\nabla f_{i_t}(\theta^t)\|_2^2] \\ &\leqslant \mathbb{E}[r_t^2] - 2\alpha_t \mathbb{E}\left[ \langle \mathbb{E}\left[ \nabla f_{i_t}(\theta^t) \mid \theta^t \right], \theta^t - \theta^* \rangle \right] + \alpha_t^2 M^2 \\ &\leqslant \mathbb{E}[r_t^2] - 2\alpha_t (\mathbb{E}[F(\theta^t)] - F(\theta^*)) + \alpha_t^2 M^2. \end{split}$$

(with abusive drops of conditional expectations, and using  $\alpha_t \geqslant 0$ ).

Summing up and using convexity of  $F(\cdot)$ , we reach the desired

$$r_0^2 + M^2 \sum_{t=0}^T \alpha_t^2 \geqslant 2 \sum_{t=0}^T \alpha_t (\mathbb{E}[F(\theta^t)] - F(\theta^\star)) \geqslant 2 \left( \sum_{t=0}^T \alpha_t \right) (\mathbb{E}[F(\bar{\theta}^T)] - F(\theta^\star)).$$

Suppose 
$$\|\theta^0 - \theta^\star\|_2 \leqslant R$$
 for some  $\theta^0 \in \mathbb{R}^d$ , and  $\mathbb{E}_i[\|\nabla f_i(\theta^t)\|_2^2] \leqslant M^2$ , then 
$$\mathbb{E}[F(\bar{\theta}^T)] - F(\theta^\star) \leqslant \frac{R^2 + M^2 \sum_{t=0}^T \alpha_t^2}{2 \sum_{t=0}^T \alpha_t},$$
 with  $\bar{\theta}^T = \frac{1}{\sum_{t=0}^T \alpha_t} \sum_{t=0}^T \alpha_t \theta^t$  (averaging).

with 
$$\bar{\theta}^T = \frac{1}{\sum_{t=0}^T \alpha_t} \sum_{t=0}^T \alpha_t \theta^t$$
 (averaging).

#### Examples:

$$\diamond \ \mathsf{Pick} \ \alpha_t = \tfrac{\alpha}{M} \colon \ F(\bar{\theta}^T) - F(\theta^\star) \leqslant \tfrac{M \|\theta^0 - \theta^\star\|_2^2 + (T+1)\alpha^2 M}{2(T+1)\alpha} = \tfrac{M \|\theta^0 - \theta^\star\|_2^2}{2(T+1)\alpha} + \tfrac{\alpha M}{2}$$

 $\diamond$  Pick  $lpha_t = \frac{\| heta^0 - heta^\star\|_2}{M \cdot / T \perp 1}$  (constant step-size depending on horizon T) then

$$F(\bar{\theta}^T) - F(\theta^*) \leqslant \frac{M \|\theta^0 - \theta^*\|_2}{\sqrt{T+1}}.$$

Suppose  $\|\theta^0 - \theta^\star\|_2 \leqslant R$  for some  $\theta^0 \in \mathbb{R}^d$ , and  $\mathbb{E}_i[\|\nabla f_i(\theta^t)\|_2^2] \leqslant M^2$ , then

$$\mathbb{E}[F(\bar{\theta}^T)] - F(\theta^*) \leqslant \frac{R^2 + M^2 \sum_{t=0}^T \alpha_t^2}{2 \sum_{t=0}^T \alpha_t},$$

with  $ar{ heta}^T = rac{1}{\sum_{t=0}^T lpha_t} \sum_{t=0}^T lpha_t heta^t$  (averaging).

 $\diamond\,$  Square summable but not summable, e.g.:  $\alpha_t = \frac{\alpha}{\mathit{M}(t+1)}$ 

$$\mathbb{E}[F(\bar{\theta}^T)] - F(\theta^*) \leqslant M \frac{\|\theta^0 - \theta^*\|_2^2 + \alpha(1+\alpha)}{2\alpha \log(T+2)},$$

 $\diamond\,$  Non-summable diminishing, example  $\alpha_t = \frac{\alpha}{M\sqrt{t+1}}$  then

$$\mathbb{E}[F(\bar{\theta}^T)] - F(\theta^*) \leqslant M \frac{\|\theta^0 - \theta^*\|_2^2 + \alpha^2(1 + \log(T+2))}{4\alpha\sqrt{T+2}}.$$

# Summing up: rough computational estimates for smooth convex minimization

Computational cost to reach  $F(\theta) - F(\theta^*) \leq \epsilon$ ?

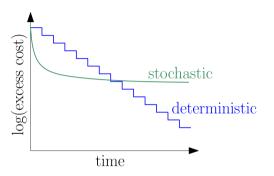
Method	Cost per iteration	# iterations	Computational cost
GD	O(nd)	$O\left(\frac{1}{\epsilon}\right)$	$O\left(\frac{nd}{\epsilon}\right)$
AGD	O(nd)	$O\left(\frac{1}{\sqrt{\epsilon}}\right)$	$O\left(rac{nd}{\sqrt{\epsilon}} ight)$
SGD	O(d)	$O\left(\frac{1}{\epsilon^2}\right)$	$O\left(\frac{d}{\epsilon^2}\right)$

- $\rightarrow$  SGD: total complexity does not depend on *n*.
- $\rightarrow$  For any  $\epsilon > 0$ , total complexity of SGD better than that of (A)GD if *n* large enough.

What target accuracy? Total computational cost:

$\epsilon$	GD	AGD	SGD	$\diamond$ Low/moderate accuracy wrt. $n$ : SGD better.
$\frac{1/\sqrt{n}}{1/n}$ $1/n^2$	$O(n^{3/2}d)$ $O(n^2d)$ $O(n^3d)$	$O(n^{5/4}d)$ $O(n^{3/2}d)$ $O(n^2d)$	$ \begin{array}{c c} O(nd) \\ O(n^2d) \\ O(n^4d) \end{array} $	<ul> <li>♦ Moderate/high accuracy wrt. n: (A)GD better.</li> <li>♦ ML: typically low/moderate accuracy.</li> </ul>

### GD vs. SGD



Example: smooth convex optimization:

- $\diamond\,$  from low to moderate accuracy requirements: SGD better.
- $\diamond$  from moderate to high accuracy requirements: (A)GD better.

# Momentum & stochasticity

```
 \begin{aligned} &\textbf{Algorithm:} \text{ Stochastic accelerated gradient} \\ &\textbf{Set } \theta^0 = \tilde{\theta}^0 \in \mathbb{R}^d, \ \{\alpha_t\}, \{\beta_t\} > 0. \\ &\textbf{for } t = 0, 1, \dots, T-1 \textbf{ do} \\ & | \text{ sample } i_t \sim \mathcal{U}[[1, n]] \\ & | \theta^{t+1} = \tilde{\theta}^t - \alpha_t \nabla f_{i_t}(\tilde{\theta}^t) \\ & | \tilde{\theta}^{t+1} = \theta^{t+1} + \beta_t(\theta^{t+1} - \theta^t) \end{aligned}   \end{aligned}
```

- $\diamond$  Classical choices: momentum  $\rightarrow$  critical noise accumulation!
- ♦ Can be mitigated via appropriate scheduling (but essentially no rate improvement).<sup>7,8</sup>

<sup>&</sup>lt;sup>7</sup>Devolder (2011). "Stochastic first order methods in smooth convex optimization."

 $<sup>^8</sup>$ Aybat, Fallah, Gurbuzbalaban, Ozdaglar (2019). "A universally optimal multistage accelerated stochastic gradient method."

### Table of contents

- 1. Stochastic optimization problems
- 2. Plain gradient methods
- 3. Stochastic gradient methods
- 4. Finite sums
- 5. Popular stochastic algorithms
- 6. Randomized coordinate descent
- 7. Conclusion

# Finite sums

### Finite sum optimization

$$\underset{\theta \in \mathbb{R}^d}{\mathsf{minimize}} \left\{ F(\theta) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(\theta) \right\}.$$

So far:

- $\diamond$  full batch (A)GD: accurate (but expensive) estimate of  $\nabla F(\theta^t)$  useless accuracy when far from solution, convergence to a solution.
- $\diamond$  SGD: cheap (but noisy) estimate of  $\nabla F(\theta^t)$  when far from solution:  $\nabla f_i(\theta^t)$  essentially points the right direction when close to solution: direction is not good.

Can we get best of both world?  $\rightarrow$  "variance reduction" techniques!

# Finite sum optimization

$$\underset{\theta \in \mathbb{R}^d}{\mathsf{minimize}} \left\{ F(\theta) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(\theta) \right\}.$$

#### Running assumptions:

- $\diamond$  each  $f_i(\cdot)$  has a Lipschitz gradient (constant L),
- $\diamond$  each  $f_i(\cdot)$  is strongly convex (constant  $\mu$ ).

Most methods below apply more generally to (but not discussed further):

- $\diamond$  each  $f_i(\cdot)$  has a Lipschitz gradient (constant L),
- $\diamond$   $F(\cdot)$  is strongly convex (constant  $\mu$ ), but all  $f_i(\cdot)$  not necessarily strongly convex.

# **Exploiting finite sums**

Instead of using  $\nabla f_{i_t}(\theta^t) \approx \nabla F(\theta^t)$ :

- $\diamond$  build running estimate  $g^t \approx \nabla F(\theta^t)$ ,
- $\diamond$  update estimate using new information  $\nabla f_{i_t}(\theta^t)$ .

Target/hopes: unbiased  $g^t \approx \nabla F(\theta^t)$  with  $\|g^t\|_2^2 \to 0$  (as  $\theta^t \to \theta^*$ ).

# **Exploiting finite sums**

Recall gradient descent  $\theta^{t+1} = \theta^t - \alpha \nabla F(\theta^t)$ . Equivalently:

$$\theta^{t+1} = \operatorname*{argmin}_{\theta \in \mathbb{R}^d} \left\{ \left| F(\theta^t) + \langle \nabla F(\theta^t), \theta - \theta^t \rangle \right| + \frac{2}{\alpha} \|\theta - \theta^t\|_2^2 \right\}$$

essentially: regularized linear approximation. What about the stochastic setting? Proposal:

$$\theta^{t+1} = \operatorname*{argmin}_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{n} \left( \sum_{i=1}^n \left| f_i(\phi_i^t) + \langle \nabla f_i(\phi_i^t), \theta - \phi_i^t \rangle \right| \right) + \frac{2}{\alpha} \|\theta - \theta^t\|_2^2 \right\}.$$

Difference with classical SGD? How to update  $\phi_i^t$ 's?

# SAG: Stochastic Average Gradient<sup>9</sup>

```
Algorithm: SAG
Set \theta^0 \in \mathbb{R}^d, \alpha > 0, \phi_i^0 = \theta^0 and g_i = \nabla f_i(\theta^0) for all i \in [[1, n]].
for t = 0, 1, ..., T - 1 do
      sample i_t \sim \mathcal{U}[[1, n]]
      \phi_i^t = \phi_i^{t-1} for all i \neq i_t
                        (save evaluated point for \nabla f_{i}.)
   egin{aligned} g_{i_t} &= 
abla f_{i_t}(	heta^t) & 	ext{(upgrade gradient of } f_{i_t}) \ g^t &= rac{1}{n} \sum_{i=1}^n g_i & 	ext{(estimate of } 
abla F(	heta^t)) \end{aligned}
      \theta^{t+1} = \overline{\theta^t} - \alpha \sigma^t
end
```

- simple to implement.
- ? more efficient computation of  $g^t$ ? ? do we really need to store matrix?
- stores  $d \times n$  matrix  $[g_1, g_2, \ldots, g_n]$ .

<sup>&</sup>lt;sup>9</sup>Schmidt. Le Roux, Bach, (2013). "Minimizing finite sums with the stochastic average gradient."

### **SAG:** observations

#### Observations:

⋄ Gradient estimate?

$$abla F(\theta^t) pprox g^t = rac{1}{n} \sum_{i=1}^n g_i.$$

♦ Unbiased?

$$\mathbb{E}_{i_t}[g^t \mid \theta^t, g^{t-1}] =$$

 $\rightarrow$ 

 $\diamond\,$  Can we do something about storage?  $\to$  for linear models, yes (later).

# Stochastic average gradient (SAG)

Let 
$$f_i$$
  $(i=1,\ldots,n)$  be  $L$ -smooth and  $\mu$ -strongly convex, and let  $\alpha=\frac{1}{16L}$ , we have 
$$\mathbb{E}[F(\theta^t)]-F(\theta^\star)\leqslant \left(1-\min\left\{\frac{\mu}{16L},\frac{1}{8n}\right\}\right)^tC_0\leqslant \exp\left(-\min\left\{\frac{\mu t}{16L},\frac{t}{8n}\right\}\right)C_0,$$
 with  $C_0=\frac{3}{2}\left(F(\theta^0)-F(\theta^\star)+\frac{4L}{n}\|\theta^0-\theta^\star\|_2^2+\frac{\sigma^2}{16L}\right)$ .

with 
$$C_0 = \frac{3}{2} \left( F(\theta^0) - F(\theta^*) + \frac{4L}{n} \|\theta^0 - \theta^*\|_2^2 + \frac{\sigma^2}{16L} \right)$$

Complexity?  $\mathbb{E}[F(\theta)] - F(\theta^*) \leq \epsilon$  in t at most

$$\exp\left(-\min\left\{\frac{\mu t}{16L}, \frac{t}{8n}\right\}\right) C_0 \leqslant \epsilon \iff t \geqslant \max\left\{16\frac{L}{\mu}, 8n\right\} \log\left(\frac{C_0}{\epsilon}\right)$$

Result actually not easy to prove. Proof relies on computer-aided verification steps.

# SAGA: Stochastic Average Gradient "Amélioré" 10

```
Algorithm: SAGA
Set \theta^0 \in \mathbb{R}^d, \alpha > 0, \phi_i^0 = \theta^0 and g_i = \nabla f_i(\theta^0) for all i \in [[1, n]].
for t = 0, 1, ..., T - 1 do
     sample i_t \sim \mathcal{U}[[1, n]]
     \phi_i^t = \phi_i^{t-1} for all i \neq i_t
     \begin{aligned} \phi_{i_t}^t &= \theta^t & \text{(save evaluated point} \\ g^t &= \nabla f_{i_t}(\theta^t) - g_{i_t} + \frac{1}{n} \sum_{i=1}^n g_i & \text{(estimate of } \nabla F(\theta^t)) \end{aligned}
                                                                     (save evaluated point for \nabla f_{i_t})
  g_{i_t} = \nabla f_{i_t}(\theta^t)
                                                                    (upgrade gradient of f_{i})
      \theta^{t+1} = \theta^t - \alpha \sigma^t
end
```



<sup>&</sup>lt;sup>10</sup>Defazio, Bach, Lacoste-Julien (2014). "SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives."

#### **SAGA:** observations

#### Observations:

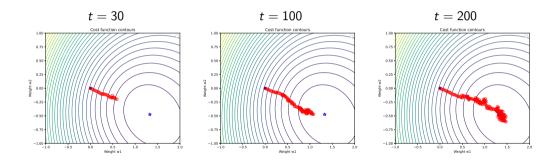
⋄ Gradient estimate?

$$abla F( heta^t) pprox oldsymbol{g}^t = 
abla f_{i_t}( heta^t) - oldsymbol{g}_{i_t} + rac{1}{n} \sum_{i=1}^n oldsymbol{g}_i.$$

♦ Unbiased?

$$\mathbb{E}_{i_t}[g^t \mid \theta^t, g^{t-1}] =$$

### SAGA: observations



# SAGA: Stochastic Average Gradient "Amélioré"

Let  $f_i$   $(i=1,\ldots,n)$  be L-smooth and  $\mu$ -strongly convex, and let  $\alpha=\frac{1}{3L}$ , we have

$$\mathbb{E}\left[\|\theta^t - \theta^\star\|_2^2\right] \leqslant \left(1 - \min\left\{\frac{1}{4n}, \frac{\mu}{3L}\right\}\right)^t C_0$$

with 
$$C_0 = \left[\|\theta^0 - \theta^\star\|_2^2 + \frac{2n}{3L}\left[F\left(\theta^0\right) - \left\langle \nabla F\left(\theta^*\right), \theta^0 - \theta^\star \right\rangle - F\left(\theta^\star\right)\right]\right].$$

Similar conclusions as for SAG: we reach  $\|\theta - \theta^{\star}\|_2^2 \leqslant \epsilon$  in at most

$$O\left(\max\left\{\kappa,n\right\}\log\left(\frac{1}{\epsilon}\right)\right)$$
.

Analysis of SAGA is considerably simpler than that of SAG.

# SAGA: Stochastic Average Gradient "Amélioré"

**Proof overview.** Show that (Lyapunov analysis): We have

$$\mathbb{E}\left[V^{t+1}\right] \leqslant \left(1 - \min\left\{\frac{1}{4n}, \frac{\mu}{3L}\right\}\right) V^{t}$$

with

$$V^{t} \triangleq V\left(\theta^{t}, \left\{\phi_{i}^{t}\right\}_{i=1}^{n}\right) \triangleq \frac{1}{n} \sum_{i=1}^{n} \left[f_{i}\left(\phi_{i}^{t}\right) - f_{i}\left(\theta^{*}\right) - \left\langle\nabla f_{i}\left(\theta^{*}\right), \phi_{i}^{t} - \theta^{*}\right\rangle\right] + c \left\|\theta^{t} - \theta^{*}\right\|_{2}^{2}$$

and 
$$c = \frac{1}{2\alpha(1-\alpha\mu)n}$$
.

Details: see arXiv.

#### **SAGA** for linear models

If learning problem can be written as

$$\underset{\theta \in \mathbb{R}^d}{\mathsf{minimize}} \frac{1}{n} \sum_{i=1}^n h_i(\langle \theta, \mathsf{x}_i \rangle) + \frac{\lambda}{2} \|\theta\|_2^2,$$

we have:  $\nabla f_i(\theta) = h_i'(\langle \theta, x_i \rangle) x_i + \lambda \theta$ . Hence, for each data point store only  $\beta_i = h_i'(\langle \theta^t, x_i \rangle)$ .

# SAGA for $\ell_2$ -regularized linear models

### Algorithm: SAGA for linear models

```
Set \theta^0 \in \mathbb{R}^d, \lambda \geqslant 0, \alpha > 0, \beta_i = h_i'(\langle \theta^0, x_i \rangle) for all i \in [[1, n]]. for t = 0, 1, \ldots, T - 1 do \begin{vmatrix} \text{sample } i_t \sim \mathcal{U}[[1, n]] \\ g^t = \\ \beta_{i_t} = \\ \theta^{t+1} = \end{vmatrix} end
```

- **?** Storage
- ? Stochastic gradient
- Gradient estimate?

No storage issue!

# Stochastic variance reduced method gradient (SVRG)<sup>11</sup>

```
Algorithm: SVRG
Set \tilde{\theta}^0 \in \mathbb{R}^d. \alpha > 0. m \in \mathbb{N}.
for s = 0, 1, ..., T do
     G^{s} = \frac{1}{n} \sum_{i=1}^{n} \nabla f_{i}(\tilde{\theta}^{s})
\theta^{0} = \tilde{\theta}^{s}
       for t = 0, 1, ..., m-1 do
             sample i_t \sim \mathcal{U}[[1, n]]
             g^t = \nabla f_i(\theta^t) - \nabla f_i(\tilde{\theta}^s) + G^s
             \theta^{t+1} = \theta^t - \alpha g^t
       end
       sample j_s \sim \mathcal{U}[[1, m]]
       \tilde{\rho}^{s+1} - \rho^{j_s}
end
```

? differences

no need to store  $d \times n$  matrix  $[g_1, g_2, \dots, g_n]$ .

 $\blacksquare$  need to tune m (inner loop).

 $<sup>^{11}</sup>$ Johnson, Zhang (2013). "Accelerating stochastic gradient descent using predictive variance reduction."

### **SVRG**: observations

⋄ Gradient estimate?

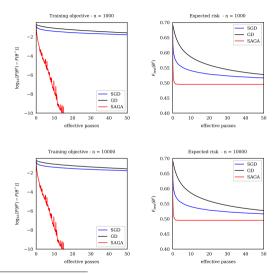
$$abla F( heta^t) pprox \mathbf{g}^t = 
abla f_{i_t}( heta^t) - 
abla f_{i_t}( ilde{ heta}^s) + rac{1}{n} \sum_{i=1}^n 
abla f_i( ilde{ heta}^s).$$

♦ Unbiasedness?

$$\mathbb{E}_{\textit{i}_t}\left[\textit{g}^t \mid \theta^t, \tilde{\theta}^\textit{s}\right] =$$

 $\rightarrow$ 

# Stochastic vs. variance reduction vs. full batch methods<sup>12</sup>



<sup>&</sup>lt;sup>12</sup>Bach (2024). "Learning theory from first principles."

# **Exploiting finite sums – momentum**

Recall template for accelerated gradient descent (iterates  $\{(\theta^t,\phi^t,\lambda^t)\}_{t=0,1,...}$ 

$$\phi^{t} = (1 - \tau_{t}) \theta^{t} + \tau_{t} \lambda^{t}$$

$$\lambda^{t+1} = \operatorname*{argmin}_{\lambda \in \mathbb{R}^{d}} \left\{ \sum_{i=0}^{t} \left[ (A_{i+1} - A_{i}) \left( F(\phi^{i}) + \langle \nabla F(\phi^{i}), \theta - \theta^{i} \rangle \right) \right] + \frac{2}{\alpha} \|\lambda - \phi^{t}\|_{2}^{2} \right\}$$

$$\theta^{t+1} = (1 - \tilde{\tau}_{t}) \theta^{t} + \tilde{\tau}_{t} \lambda^{t+1}$$

... similarly: based on regularized (weighted) linear approximations of  $F(\cdot)$ 

(with growing sequence  $\{A_t\}_{t=0,1,...}$  and some  $\{(\tau_k, \tilde{\tau}_k)\}_{t=0,1,...}$  for convex combinations).

#### Momentum versions

A few momentum variations exist. Among the simplest ones:<sup>13</sup>

```
Algorithm: SAGA with Sampled Negative Momentum
Set \theta^0 \in \mathbb{R}^d, \alpha, \tau > 0, \phi_i^0 = \nabla f_i(\theta^0) for all i \in [[1, n]].
for t = 0, 1, ..., T - 1 do
      sample i_t \sim \mathcal{U}[[1, n]]
      \tilde{\theta}^t = \tau \theta^t + (1 - \tau) \phi_i^t
      g^t = \nabla f_{i_t}(\tilde{\theta}^t) - \nabla f_{i_t}(\phi_{i_t}^t) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\phi_{i_t}^t)
      \theta^{t+1} = \theta^t - \alpha \sigma^t
      sample i_t \sim \mathcal{U}[[1, n]]
      \phi_{i}^{t+1} = \tau \theta^{t+1} + (1-\tau)\phi_{i}^{t}
end
```

differences

# gradient evaluations

 $<sup>^{13}</sup>$ Zhou et al. (2019). "Direct acceleration of SAGA using sampled negative momentum."

# Takeaways from variance reduction

- Finite-sums methods use only one stochastic gradient per iteration and converge linearly on strongly convex functions.
- Choice of fixed (nondecreasing) step-size possible.
- ♦ SAGA only needs to know the smoothness parameter, but requires storing *n* past stochastic gradients in general (but not for linear classifier).
- $\diamond$  SVRG only has O(d) storage in general, but requires full gradient computations every so often. Has an extra "number of inner iterations" parameter.



# Summing up: rough computational cost estimates

Method	# iterations	# gradient queries
GD	$O\left(\kappa\log\left(\frac{1}{\epsilon}\right)\right)$	$O\left(n\kappa\log\left(\frac{1}{\epsilon}\right)\right)$
AGD	$O\left(\sqrt{\kappa}\log\left(rac{1}{\epsilon} ight) ight)$	$O\left(n\sqrt{\kappa}\log\left(\frac{1}{\epsilon}\right)\right)$
SAG/SAGA/SVRG	$O\left( max\{n,\kappa\}\log\left(rac{1}{\epsilon} ight) ight)$	$O\left(\max\{n,\kappa\}\log\left(\frac{1}{\epsilon}\right)\right)$
$Katyushia^{14/MiG^{15}/SSNM^{16}/Pt\text{-}SAGA^{17}}$	$O\left(\max\{n,\sqrt{n\kappa}\}\log\left(\frac{1}{\epsilon}\right)\right)$	$O\left(\max\{n,\sqrt{n\kappa}\}\log\left(\frac{1}{\epsilon}\right)\right)$

So: finite-sum methods benefit from momentum when  $n \ll \kappa$ . That is:

- $\diamond \max\{n,\kappa\} = \kappa \to \text{computational complexities of SAG/SAGA/SVRG is } O\left(\kappa \log\left(\frac{1}{\epsilon}\right)\right).$
- $\diamond \ \max\{n, \sqrt{n\kappa}\} = \sqrt{n\kappa} \to \text{computational complexities of momentum variants is}$

$$O\left(\sqrt{n\kappa}\log\left(\frac{1}{\epsilon}\right)\right) \ll O\left(\kappa\log\left(\frac{1}{\epsilon}\right)\right)$$
 .

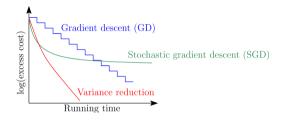
<sup>&</sup>lt;sup>14</sup>Allen-Zhu (2017). "Katyusha: The first direct acceleration of stochastic gradient methods."

 $<sup>^{15}</sup>$ Zhou, Shang, Cheng (2018). "A simple stochastic variance reduced algorithm with fast convergence rates."

 $<sup>^{16}</sup>$ Zhou et al. (2019). "Direct acceleration of SAGA using sampled negative momentum."

<sup>&</sup>lt;sup>17</sup>Defazio (2016). "A simple practical accelerated method for finite sums."

#### Stochastic vs. variance reduction vs. full batch methods



To experiment with those:



## Table of contents

- 1. Stochastic optimization problems
- 2. Plain gradient methods
- 3. Stochastic gradient methods
- 4. Finite sums
- 5. Popular stochastic algorithms
- 6. Randomized coordinate descent
- 7. Conclusion

# Popular stochastic algorithms

# **Practical improvements**

#### Practical improvements:

- adapt to observations,
- adapt componentwise,
- ⋄ momentum,
- different step-size schedules.

#### Generally, either

- o no existing analysis,
- or extremely technical.

Optimizers in Pytorch  $\rightarrow$ 



# Adagrad<sup>18,19</sup>

Adagrad<sup>16</sup> (update all components j):

$$g^{t} = \nabla f_{i_{t}} (\theta^{t-1})$$

$$v_{(j)}^{t} = v_{(j)}^{t-1} + (g_{(j)}^{t})^{2}$$

$$\theta_{(j)}^{t} = \theta_{(j)}^{t-1} - \frac{\alpha}{\sqrt{\epsilon + v_{(j)}^{t}}} g_{(j)}^{t}$$

For certain parameter choices:<sup>19</sup>

$$\mathbb{E}\left[\|\nabla F(\theta^t)\|_2^2\right] = O\left(\frac{1}{\sqrt{t}}\right)$$

for smooth objectives.

Adagrad in Pytorch  $\rightarrow$ 

#### Adagrad

CLASS torchoptim.Adagrad(params, Ir=0.01, Ir\_decay=0, meight\_decay=0, initial\_accumulator\_value=0, eps=1e-10, foreach=None, \*, maximize=False, differentiable=False, fused=None) [SOURCE]

Implements Adagrad algorithm.

input :  $\gamma$  (lr),  $\theta_0$  (params),  $f(\theta)$  (objective),  $\lambda$  (weight decay),  $\tau$  (initial accumulator value),  $\eta$  (lr decay) initialize : state, sum  $\leftarrow \tau$ 

 $\begin{aligned} & \text{for } t = 1 \text{ to } \dots \text{ do} \\ & g_t \leftarrow \nabla g_t(\theta_{t-1}) \\ & \tilde{\gamma} \leftarrow \gamma / (1 + (t-1)\eta) \\ & \text{ if } \lambda \neq 0 \\ & g_t \leftarrow g_t + \lambda \theta_{t-1} \\ & state\_sum_t \leftarrow state\_sum_{t-1} + g_t^2 \\ & \theta_t \leftarrow \theta_{t-1} - \tilde{\gamma} - \frac{g_t}{\sqrt{state\_sum_{t-1}}} + \varepsilon \end{aligned}$ 

 ${f return}\ heta_{f t}$ 

 $<sup>^{18}</sup>$ Duchi, Hazan, Singer (2011). "Adaptive subgradient methods for online learning and stochastic optimization."

<sup>&</sup>lt;sup>19</sup>Défossez, Bottou, Bach, Usunier (2020). "A simple convergence proof of Adam and Adagrad."

Adam<sup>20</sup> (update all components j):

$$egin{aligned} egin{aligned} egi$$

For certain parameter choices:<sup>21</sup>

$$\mathbb{E}\left[\|
abla F( heta^t)\|_2^2
ight] = O\left(rac{\log t}{\sqrt{t}}
ight)$$

for smooth objectives.

### Adam in Pytorch $\rightarrow$

#### Adam

CLASS torchoptim.Adam(params, 1r=0.001, betas=(0.9, 0.999), eps=1e-08, weight\_decay=0, amsgraoFalse, \*, foreach=None, maximize=False, capturable=False, differentiable=False, fusea-None) [SOURCE]

Implements Adam algorithm.

input :  $\gamma$  (lr),  $\beta_1$ ,  $\beta_2$  (betas),  $\theta_0$  (params),  $f(\theta)$  (objective) λ (weight decay), amsorad, maximize initialize:  $m_0 \leftarrow 0$  (first moment),  $n_0 \leftarrow 0$  (second moment),  $\overline{m}^{max} \leftarrow 0$ for t = 1 to ... do if maximize:  $a_i \leftarrow -\nabla_a f_i(\theta_{i-1})$  $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ if  $\lambda \neq 0$  $q_t \leftarrow q_t + \lambda \theta_{t-1}$  $m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) q_t$  $v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) a_t^2$  $\widehat{m_i} \leftarrow m_i/(1 - \beta_i^t)$  $\mathfrak{V}_i \leftarrow v_i/(1-\beta_i^t)$ if amsgrad  $\overline{v}_i^{max} \leftarrow \max(\overline{v}_i^{max}, \overline{v}_i)$  $\theta_{i} \leftarrow \theta_{i-1} - \gamma \widehat{m}_{i} / (\sqrt{\widehat{v}_{i}^{max}} + \epsilon)$ 

 $\mathbf{return}\ \theta_{\mathrm{t}}$ 

 $\theta_i \leftarrow \theta_{i-1} - \sqrt{m_i}/(\sqrt{v_i} + \epsilon)$ 

<sup>&</sup>lt;sup>20</sup>Kingma, Ba (2014). "Adam: A method for stochastic optimization."

<sup>&</sup>lt;sup>21</sup>Défossez, Bottou, Bach, Usunier (2020). "A simple convergence proof of Adam and Adagrad."

# Randomized coordinate descent

# (One possible) motivation: back to supervised learning

$$\underset{\theta \in \mathbb{R}^d}{\mathsf{minimize}} \frac{1}{n} \sum_{i=1}^n h_i(\langle \theta, \mathsf{x}_i \rangle) + \frac{\lambda}{2} \|\theta\|_2^2.$$

What did we do with stochastic methods?

- $\rightarrow$  update parameter estimation, one sample at a time.
- → Other ways to do that? One possibility: artificially augmented problem:

$$\underset{\substack{\theta \in \mathbb{R}^d \\ \beta_1, \dots, \beta_n \in \mathbb{R}}}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n h_i(\beta_i) + \frac{\lambda}{2} \|\theta\|_2^2 \quad \text{s.t. } \beta_i = \langle \theta, x_i \rangle \text{ for } i = 1, \dots, n.$$

Introduce dual variables  $\omega_1, \ldots, \omega_n$ ; Lagrangian dual is:

ightarrow Use coordinate-based methods on dual. <sup>22</sup>

<sup>&</sup>lt;sup>22</sup>Shalev-Shwartz, Zhang (2013). "Stochastic dual coordinate ascent methods for regularized loss minimization."

#### Randomized block-coordinate methods

$$\underset{\omega \in \mathbb{R}^d}{\mathsf{minimize}} \ D(\omega)$$

where f is L-smooth and convex. Decompose decision space into n blocks:

$$\omega = \sum_{i=1}^n \mathbf{U}_i \omega$$
 with  $[\mathbf{U}_1 \mathbf{U}_2 \dots \mathbf{U}_n] = I_d$ .

#### Algorithm: RBCD

end

$$\begin{split} \text{Set } \omega^0 \in \mathbb{R}^d, \ \alpha > 0. \\ \text{for } t = 0, 1, \dots, T-1 \text{ do} \\ \mid \text{ sample } i_t \sim \mathcal{U}[[1, n]] \\ \omega^{t+1} = \omega^t - \alpha \mathbf{U}_{i_t} \nabla D(\omega^t) \end{split}$$

update rule corresponds to

Example: what  $\{\mathbf{U}_i\}_{i=1}^n$  corresponds to single coordinate decomposition?

# Randomized block-coordinate methods: convergence

Let 
$$\omega^t \in \mathbb{R}^d$$
,  $\omega^{t+1} = x^t - \alpha \mathbf{U}_{i_t} \nabla D(\omega^t)$  with  $\alpha \in (0, \frac{1}{L}]$ ,  $i_t \sim \mathcal{U}[[1, n]]$ . One can show: 
$$A_{t+1} \mathbb{E}[D(\omega^{t+1}) - D(\omega^\star)] + \frac{L}{2} \mathbb{E}[\|\omega^{t+1} - \omega^\star\|_2^2] \leqslant A_t(D(\omega^t) - D(\omega^\star)) + \frac{L}{2} \|\omega^t - \omega^\star\|_2^2$$
 for any  $A_t \geqslant 1$  and  $A_{t+1} = A_t + \frac{\alpha L}{n}$ .

- Many results, variants, etc. Easily fall into additional technical difficulties.
- More conventional to assume Lipschitz by block (simpler to compute and more aggressive step size strategies), but this result is simple.
- $\diamond$  Guarantee:  $\mathbb{E}[D(\omega^t) D(\omega^\star)] \leqslant \frac{n}{t} (D(\omega^0) D(\omega^\star) + \frac{L}{2} ||\omega^0 \omega^\star||_2^2)$  with  $\alpha = \frac{1}{L}$ .
- $\diamond$  Recall gradient descent:  $\mathbb{E}[D(\omega^t) D(\omega^\star)] \leqslant \frac{L}{2t} \|\omega^0 \omega^\star\|_2^2$ .

# Randomized block-coordinate methods – improvement

$$\underset{\omega \in \mathbb{R}^d}{\mathsf{minimize}} D(\omega)$$

where f is convex. Decompose decision space into n blocks:  $\omega = \sum_{i=1}^n \mathbf{U}_i \omega$  with  $[\mathbf{U}_1 \, \mathbf{U}_2 \, \dots \, \mathbf{U}_n] = I_d$ . Further assume  $\forall i \in \{1, 2, \dots, n\}$  and  $\forall \Delta \in \mathbb{R}^d$ :

$$D(x + \mathbf{U}_i \Delta) \leq D(x) + \langle \nabla D(x), \mathbf{U}_i \Delta \rangle + \frac{L_i}{2} \|\mathbf{U}_i \Delta\|_2^2.$$

#### Algorithm: RBCD

Set 
$$\omega^0 \in \mathbb{R}^d$$
.

$$\begin{array}{l} \textbf{for} \,\, t = 0, 1, \dots, T-1 \,\, \textbf{do} \\ \quad \text{sample} \,\, i_t \sim \mathcal{U}[[1, n]] \\ \quad \omega^{t+1} = \omega^t - \frac{1}{L_{l_t}} \textbf{U}_{l_t} \nabla D(\omega^t) \end{array}$$

end

- $\diamond L_i$  usually simpler to compute than L
- $\diamond$   $L_i$  often (much) smaller than L.

#### On RBCD

#### Questions:

- 1. Is the gradient estimate  $\mathbf{U}_{i_t} \nabla D(\omega^t)$  biased?
- 2. Consider the quadratic problem

$$\underset{\omega \in \mathbb{R}^d}{\mathsf{minimize}} \ \tfrac{1}{2} \omega^T A \omega$$

and the decomposition  $U_i = e_i$  (unit vector whose *i*th component is one).

- What do the  $L_i$ 's (i = 1, ..., d) correspond to?
- Show that the global Lipschitz constant L satisfies:  $\max_{1 \leqslant i \leqslant d} L_i \leqslant L \leqslant \sum_{i=1}^d L_i$ .
- Consider the matrix  $A = c \mathbf{1} \mathbf{1}^T$ . What are  $L_i$ 's? and L?

# Randomized block-coordinate methods — improvement

Denote  $\|\omega\|_{\{L_i\}}^2 \triangleq \sum_{i=1}^n L_i \|\mathbf{U}_i \omega\|_2^2$ .

Let 
$$\omega^t \in \mathbb{R}^d$$
,  $\omega^{t+1} = \omega^t - \frac{1}{L_{i_t}} \mathbf{U}_{i_t} \nabla D(\omega^t)$  with  $i_t \sim \mathcal{U}\{1, \dots, n\}$ . It holds:

Let 
$$\omega^t \in \mathbb{R}^d$$
,  $\omega^{t+1} = \omega^t - \frac{1}{L_{i_t}} \mathbf{U}_{i_t} \nabla D(\omega^t)$  with  $i_t \sim \mathcal{U}\{1,\ldots,n\}$ . It holds: 
$$A_{t+1} \mathbb{E}[D(\omega^{t+1}) - D(\omega^\star)] + \frac{1}{2} \mathbb{E}[\|\omega^{t+1} - \omega^\star\|_{\{L_i\}}^2] \leqslant A_t(D(\omega^t) - D(\omega^\star)) + \frac{1}{2} \|\omega^t - \omega^\star\|_{\{L_i\}}^2$$
 for any  $A_t \geqslant 1$  and  $A_{t+1} = A_t + \frac{1}{n}$ .

- ♦ Usually simpler to compute and allows for larger step-sizes.
- More conventional to assume Lipschitz by block.
- $\diamond \text{ Guarantee: } \mathbb{E}[D(\omega^t) D(\omega^\star)] \leqslant \frac{n}{t} \left( D(\omega^0) D(\omega^\star) + \frac{1}{2} \|\omega^0 \omega^\star\|_{L^{1/3}}^2 \right).$
- Possible to extend results to linear convergence (strong convexity-type assumptions).<sup>23</sup>

<sup>&</sup>lt;sup>23</sup>See, e.g., Nesterov (2012). "Efficiency of coordinate descent methods on huge-scale optimization problems."

#### Randomized block-coordinate methods

#### **Proof sketch.** Weighted sum of inequalities:

 $\diamond$  convexity of F between  $\omega^t$  and  $\omega^*$ , with weight  $A_{t+1} - A_t$ :

$$0 \geqslant D(\omega^t) - D(\omega^*) + \langle \nabla D(\omega^t), \omega^* - \omega^t \rangle,$$

 $\diamond$  expectation of the "block" descent lemma with weight  $A_{t+1}$ :

$$\mathbb{E}_{i_t}[D(\omega^{t+1})] \leqslant D(\omega^t) - \mathbb{E}_i\left[\frac{1}{2L_i}\|\mathbf{U}_i\nabla D(\omega^t)\|_2^2\right].$$

Weighted sum yields:

$$\begin{split} \mathbb{E}_{i_t}[V^{t+1}] \leqslant V^t - \frac{A_{t+1}-1}{2n} \|\nabla D(\omega^t)\|_{\{L_i^{-1}\}}^2 \\ + \left(A_{t+1} - A_t - \frac{1}{n}\right) \langle \nabla D(\omega^t), \omega^t - \omega^\star \rangle, \end{split}$$

with 
$$V^t = A_t(D(\omega^t) - D(\omega^*)) + \frac{1}{2} \|\omega^t - \omega^*\|_{\{L_i\}}^2$$
.

# Example - support vector machine

Soft-margin support vector machine (SVM):

$$\min_{ heta \in \mathbb{R}^d} \sum_{t=1}^{1} \| heta\|_2^2 + 
u \sum_{i=1}^n \max\left\{0, 1 - y_i \langle heta, x_i 
angle
ight\}$$

Reformulate:

$$egin{aligned} & \min_{ heta \in \mathbb{R}^d, \, s \in \mathbb{R}^n} rac{1}{2} \| heta \|_2^2 + 
u \sum_{i=1}^n s_i \ & ext{s.t. } y_i \langle heta, x_i 
angle \geqslant 1 - s_i \ & s_i \geqslant 0 \end{aligned}$$

Lagrange dual?

# Example – support vector machine

Denote by  $X = [y_1x_1 \mid y_2x_2 \mid \dots \mid y_nx_n] \in \mathbb{R}^{d \times n}$ . Lagrange duality yields:

$$\underset{0 \leqslant \lambda \leqslant \nu}{\text{maximize}} \left\{ D(\lambda) \triangleq -\frac{1}{2} \lambda^T X^T X \lambda + \sum_{i=1}^n \lambda_i \right\}$$

and a natural estimate of the primal variable  $\theta = \sum_{i=1}^{n} \lambda_i x_i y_i = X \lambda$ . Algorithm?

Algorithm: RBCD for dual SVM Set  $\lambda^0 \in \mathbb{R}^n$ . for t = 0, 1, ..., T - 1 do  $\begin{vmatrix} \text{ sample } i_t \sim \mathcal{U}[[1,n]] \\ \lambda_{(i_t)}^{t+1} = \operatorname{Proj}_{[0,\alpha]} \left[ \omega_{(i_t)}^t - \frac{1}{L_{i_t}} \nabla_{i_t} D(\lambda^t) \right] \end{vmatrix}$ end

- $\diamond \ \lambda_{(i)}^t$  denotes ith component.  $\diamond$  Projection OK within BCD for separable constraints.
- ⋄ L<sub>i</sub>'s?⋄ Exact 1-D optimization.

# **Example – support vector machine**

Denote by  $X = [y_1x_1 \mid y_2x_2 \mid \dots \mid y_nx_n] \in \mathbb{R}^{d \times n}$ . Lagrange duality yields:

$$\underset{0 \leqslant \lambda \leqslant \nu}{\text{maximize}} \left\{ D(\lambda) \triangleq -\frac{1}{2} \lambda^T X^T X \lambda + \sum_{i=1}^n \lambda_i \right\}$$

and a natural estimate of the primal variable  $\theta = \sum_{i=1}^{n} \lambda_i x_i y_i = X \lambda$ . Algorithm?

$$\begin{aligned} &\textbf{Algorithm:} \ \mathsf{RBCD} \ \mathsf{for} \ \mathsf{dual} \ \mathsf{SVM} \\ &\mathsf{Set} \ \lambda^0 = 0 \in \mathbb{R}^n, \ \theta^0 = 0 \in \mathbb{R}^d. \\ &\textbf{for} \ t = 0, 1, \dots, \mathcal{T} - 1 \ \textbf{do} \\ & | \ \mathsf{sample} \ i_t \sim \mathcal{U}[[1, n]] \\ & | \ \bar{\lambda} = \lambda^t_{(i_t)} \\ & | \ \lambda^{t+1}_{(i_t)} = \mathrm{Proj}_{[0, \alpha]} \left( \lambda^t_{(i_t)} + \frac{1 - y_{i_t} \langle \theta^t, \mathsf{x}_{i_t} \rangle}{||x_{i_t}||_2^2} \right) \\ & | \ \theta^{t+1} = \theta^t + y_{i_t} \mathsf{x}_{i_t} \left( \lambda^{t+1}_{(i_t)} - \bar{\lambda} \right) \end{aligned}$$
 
$$\mathbf{end}$$

### Table of contents

- 1. Stochastic optimization problems
- 2. Plain gradient methods
- 3. Stochastic gradient methods
- 4. Finite sums
- 5. Popular stochastic algorithms
- 6. Randomized coordinate descent
- 7. Conclusion



**Conclusion** 

## **Concluding remarks**

#### What did we do?

- exploit problem structures (finite sums/expectations).
- cheaper iterations vs. slower convergence per iteration.
- Different stochastic/randomized strategies.

#### Methods of extreme practical use, particularly when:

- even computing a gradient is too expensive,
- updates without accouting for full dataset,
- ♦ accurate solution not needed (no need to go beyond statistical accuracy).

