

# Stochastic Approximation

Francis Bach, Aymeric Dieuleveut, Alain Durmus, Eric Moulines

Ecole Polytechnique, Centre de Mathematiques Appliquees

July 21, 2021

# Context

## Machine learning for “big data”

- Large-scale machine learning: large  $d$ , large  $n$ 
  - $d$ : dimension of each observation (input)
  - $n$ : number of observations
- **Examples:** computer vision, bioinformatics, advertising
- **Ideal running-time complexity:**  $O(dn)$

# Context

## Machine learning for “big data”

- Large-scale machine learning: large  $d$ , large  $n$ 
  - $d$ : dimension of each observation (input)
  - $n$ : number of observations
- Examples: computer vision, bioinformatics, advertising
- Ideal running-time complexity:  $O(dn)$
- Going back to simple methods
  - Stochastic gradient methods (Robbins, Monro, 1951)

# Context

## Machine learning for “big data”

- Large-scale machine learning: large  $d$ , large  $n$ 
  - $d$ : dimension of each observation (input)
  - $n$ : number of observations
- Examples: computer vision, bioinformatics, advertising
- Ideal running-time complexity:  $O(dn)$
- Going back to simple methods
  - Stochastic gradient methods (Robbins, Monro, 1951)
  - Mixing statistics and optimization

# 1 Stochastic approximation

## 2 Proximal methods

## 3 Applications

- 1 Stochastic approximation
  - Finite-sum optimization
  - Online learning
  - Smooth strongly convex case
  - Stochastic subgradient descent/method
  - Stochastic Approximation for nonconvex optimization
- 2 Proximal methods
  - Proximal operator
  - Proximal gradient algorithm
  - Stochastic proximal gradient
- 3 Applications
  - Network structure estimation
  - High-dimensional logistic regression with random effect

# Finite-sum optimisation

## Empirical risk minimization

- Finite set of observations:  $Z_1, \dots, Z_n$  (typically,  $Z_i(Y_i, X_i)$ )
- Minimize the empirical risk:  $\hat{f}(\theta) = \frac{1}{n} \sum_{k=1}^n \ell(\theta, Z_i)$

## Batch stochastic gradient

- Let  $S \subset \{1, \dots, n\}$  be a mini-batch sampled with/without replacement in  $\{1, \dots, n\}$  with cardinal  $|S| = N$ .
- Define the mini-batch gradient

$$\nabla \hat{f}_S(\theta) = (1/p) \sum_{i \in S} \nabla_{\theta} \ell(\theta, Z_i),$$

where  $p = n/N$  or  $p = 1/\binom{N}{n}$ .

- Then,  $\nabla \hat{f}_S$  is an unbiased estimator of  $\nabla \hat{f}$ , i.e.

$$\mathbb{E}[\nabla \hat{f}_S(\theta) | (Z_i)_{i \in \{1, \dots, n\}}] = \nabla \hat{f}(\theta).$$

# Batch Stochastic Gradient

## Empirical risk minimization

- Minimize the empirical risk:  $\hat{f}(\theta) = \frac{1}{n} \sum_{k=1}^n \ell(\theta, Z_i)$

## Batch stochastic gradient

- Batch stochastic optimization consists in replacing  $\nabla \hat{f}(\theta_k)$  by the minibatch estimate  $\nabla \hat{f}_{S_{k+1}}(\theta_k)$  in the gradient descent scheme to define the iterates  $(\theta_k)_{k \in \mathbb{N}}$ ,

$$\theta_{k+1} = \theta_k - \gamma_{k+1} \nabla \hat{f}_{S_{k+1}}(\theta_k),$$

where  $(S_k)$  is an i.i.d. sequence of minibatches and  $(\gamma_k)_{k \in \mathbb{N}^*}$  is a sequence of stepsizes.



# Batch Stochastic Gradient

- $(S_k)_{k \in \mathbb{N}^*}$  uniform with/without replacement non necessary the best choice.
- $(\gamma_k)_{k \in \mathbb{N}^*}$  is either held constant or decreasing going to 0:
  - **constant stepsize:** If  $\gamma_k \equiv \gamma$ , the scheme does not converge in general.  $\{\theta_k^\gamma\}$  is an ergodic Markov chain (under appropriate conditions).
  - **decreasing stepsize:** If  $\lim_{k \rightarrow +\infty} \gamma_k = 0$ , then  $\{\theta_k\}$  converges a.s. to  $\theta_*$  (also under appropriate conditions).
- This is a specific instance of **stochastic approximation** schemes.

# Online learning

## Expected risk minimization

- Minimize the expected risk:  $f(\theta) = \mathbb{E}[\ell(\theta, Z)]$

## Online stochastic gradient

- Let  $(Z_k)_{k \in \mathbb{N}^*}$  be an i.i.d. sequence.
- Define for any  $k \in \mathbb{N}^*$ ,

$$\nabla f_k(\theta) = \nabla_{\theta} \ell(\theta, Z_k) .$$

- Then,  $\nabla f_k$  is an unbiased estimator of  $\nabla f$ , i.e.

$$\mathbb{E}[\nabla \hat{f}_k(\theta)] = \nabla \hat{f}(\theta)$$

where the expectation is taken over the data  $(Z_k)_{k \in \mathbb{N}^*}$ .

- 1 Stochastic approximation
  - Finite-sum optimization
  - Online learning
  - Smooth strongly convex case
  - Stochastic subgradient descent/method
  - Stochastic Approximation for nonconvex optimization
- 2 Proximal methods
  - Proximal operator
  - Proximal gradient algorithm
  - Stochastic proximal gradient
- 3 Applications
  - Network structure estimation
  - High-dimensional logistic regression with random effect

# Online learning

- Minimize the expected risk:  $f(\theta) = \mathbb{E}[\ell(\theta, Z)]$

## Online stochastic gradient

- Online stochastic gradient defines the iterates  $(\theta_k)_{k \in \mathbb{N}}$ ,

$$\theta_{n+1} = \theta_n - \gamma_{n+1} \nabla f_{n+1}(\theta_n),$$

where  $(\gamma_k)_{k \in \mathbb{N}^*}$  is a sequence of stepsizes.

## Remarks

- $(\gamma_k)_{k \in \mathbb{N}^*}$  is either constant or decrease to 0.
- This scheme also belongs to the class of stochastic approximation/optimization schemes.

# Stochastic gradient descent

- Minimize a function  $f$  defined on  $\mathbb{R}^d$
- given only unbiased estimates  $\nabla f_n$  of  $\nabla f$ ,
- or  $\partial f_n$  of its subgradients  $\partial f$ .

## Online learning

- loss for a single pair of observations:  $f_n(\theta) = \ell(Y_n, \langle \theta, \Phi(X_n) \rangle)$
- $f(\theta) = \mathbb{E}[f_n(\theta)] = \mathbb{E}[\ell(Y_n, \langle \theta, \Phi(X_n) \rangle)] =$  generalization error
- Expected gradient:

$$\nabla f(\theta) = \mathbb{E}[\nabla f_n(\theta)] = \mathbb{E}[\dot{\ell}(Y_n, \langle \theta, \Phi(X_n) \rangle) \Phi(X_n)]$$

- Non-asymptotic results

Number of iterations = number of observations

# Convex stochastic approximation

- Smoothness:  $f$   $B$ -Lipschitz continuous,  $\nabla f$   $L$ -Lipschitz continuous
- Strong convexity:  $f$   $\mu$ -strongly convex

Key algorithm: Stochastic (sub)gradient descent (a.k.a. Robbins-Monro)

$$\theta_n = \theta_{n-1} - \gamma_n \nabla f_n(\theta_{n-1}), \quad \theta_n = \theta_{n-1} - \gamma_n \partial f_n(\theta_{n-1})$$

- 1 Stochastic approximation
  - Finite-sum optimization
  - Online learning
  - Smooth strongly convex case
  - Stochastic subgradient descent/method
  - Stochastic Approximation for nonconvex optimization
- 2 Proximal methods
  - Proximal operator
  - Proximal gradient algorithm
  - Stochastic proximal gradient
- 3 Applications
  - Network structure estimation
  - High-dimensional logistic regression with random effect

# Convex stochastic approximation

Key properties of  $f$  and/or  $f_n$

- Smoothness:  $f$   $B$ -Lipschitz continuous,  $\nabla f$   $L$ -smooth
- Strong convexity:  $f$   $\mu$ -strongly convex

Key algorithm: Stochastic (sub)gradient descent (a.k.a. Robbins-Monro)

$$\theta_n = \theta_{n-1} - \gamma_n \nabla f_n(\theta_{n-1}), \quad \theta_n = \theta_{n-1} - \gamma_n \partial f_n(\theta_{n-1})$$

- Polyak-Ruppert averaging:  $\bar{\theta}_n = n^{-1} \sum_{k=0}^{n-1} \theta_k$
- Which learning rate sequence  $\gamma_n$ ? Classical setting:  $\gamma_n = Cn^{-\alpha}$

Desirable practical behavior

- Applicable (at least) to classical supervised learning problems
- Robustness to (potentially unknown) constants ( $L, B, \mu$ )
- Adaptive to problem behavior (e.g., convex / strongly convex)



## Smoothness/convexity assumptions

Iteration  $\theta_n = \theta_{n-1} - \gamma_n \nabla f_n(\theta_{n-1})$ .

Polyak-Ruppert averaging  $\bar{\theta}_n = \frac{1}{n} \sum_{k=0}^{n-1} \theta_k$

$f_n$  Convex +  $L$ -Smooth : For each  $n \geq 1$  the function  $f_n$  satisfies a.s.:

- convex;
- differentiable with  $L$ -Lipschitz-continuous gradient  $\nabla f_n$ ;
- bounded variance (bounded data): almost surely

$$\mathbb{E}[\|\nabla f_{n+1}(\theta^*)\|^2 | \mathcal{F}_n] \leq \sigma^2 .$$

$f$  Strongly convex : The function  $f$  is strongly convex with respect to the norm  $\|\cdot\|_2$  with convexity constant  $\mu > 0$ :

- Invertible population covariance matrix or regularization by  $\frac{\mu}{2} \|\theta\|^2$
- $\Rightarrow$  there exists a unique minimizer  $\theta^*$

# Summary

## Assumptions

- Stochastic gradient descent with learning rate  $\gamma_n = Cn^{-\alpha}$ ,  $\alpha \in [0, 1]$
- Strongly convex smooth objective functions
- Bounded variance (bounded data): w.p. 1,  
 $\mathbb{E}[\|\nabla f_{n+1}(\theta^*)\|^2 | \mathcal{F}_n] \leq \sigma^2$ .

## Results

- Old:  $O(n^{-1})$  rate achieved without averaging for  $\alpha = 1$
- New:  $O(n^{-1})$  rate achieved with averaging for  $\alpha \in [1/2, 1]$
- Non-asymptotic analysis with explicit constants
- Robust to the choice of  $C$

# Summary

## Assumptions

- Stochastic gradient descent with learning rate  $\gamma_n = Cn^{-\alpha}$ ,  $\alpha \in [0, 1]$
- Strongly convex smooth objective functions
- Bounded variance (bounded data): w.p. 1,  
 $\mathbb{E}[\|\nabla f_{n+1}(\theta^*)\|^2 | \mathcal{F}_n] \leq \sigma^2$ .

## Results

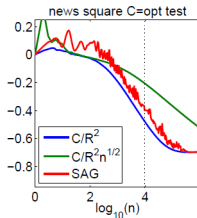
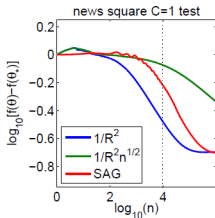
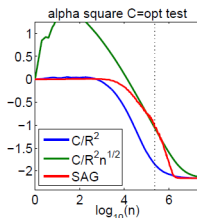
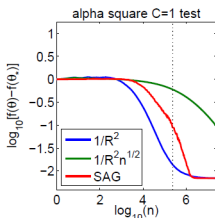
- Old:  $O(n^{-1})$  rate achieved without averaging for  $\alpha = 1$
- New:  $O(n^{-1})$  rate achieved with averaging for  $\alpha \in [1/2, 1]$
- Non-asymptotic analysis with explicit constants
- Robust to the choice of  $C$

Convergence rate for  $\mathbb{E}[\|\theta_n - \theta^*\|^2]$  and  $\mathbb{E}[\|\bar{\theta}_n - \theta^*\|^2]$ .

- without averaging:  $O(\gamma_n) + O(e^{-\mu n \gamma_n}) \|\theta_0 - \theta^*\|^2$
- with averaging:  $O(n^{-1}) + O(n^{-2\alpha}) + \mu^{-2} \|\theta_0 - \theta^*\|^2 O(n^{-2})$

# Examples

- *alpha* ( $d = 500$ ,  $n = 500\,000$ ), *news* ( $d = 1\,300\,000$ ,  $n = 20\,000$ )



# Sketch of proof

$f$  strongly convex,  $f_n$  smooth, bounded variance

- Consider  $\delta_n = \|\theta_n - \theta^*\|^2$ .
- Then, we have almost surely

$$\delta_{n+1} = \delta_n - \gamma_{n+1} \langle \nabla f_{n+1}(\theta_n), \theta_n - \theta^* \rangle + \gamma_{n+1}^2 \|\nabla f_{n+1}(\theta_n)\|^2 .$$

- $f$  is strongly convex:

$$\begin{aligned} \mathbb{E}[\delta_{n+1} | \mathcal{F}_n] &= \delta_n - \gamma_{n+1} \langle \nabla f(\theta_n), \theta_n - \theta^* \rangle + \gamma_{n+1}^2 \mathbb{E}[\|\nabla f_{n+1}(\theta_n)\|^2 | \mathcal{F}_n] \\ &\leq (1 - \mu\gamma_{n+1})\delta_n + \gamma_{n+1}^2 \mathbb{E}[\|\nabla f_{n+1}(\theta_n) - \nabla f(\theta^*)\|^2 | \mathcal{F}_n] . \end{aligned}$$

## Sketch of proof

$f$  strongly convex,  $f_n$  smooth, bounded variance

- Consider  $\delta_n = \|\theta_n - \theta^*\|^2$ .
- $f$  is strongly convex:

$$\begin{aligned}\mathbb{E}[\delta_{n+1} | \mathcal{F}_n] &= \delta_n - \gamma_{n+1} \langle \nabla f(\theta_n), \theta_n - \theta^* \rangle + \gamma_{n+1}^2 \mathbb{E}[\|\nabla f_{n+1}(\theta_n)\|^2 | \mathcal{F}_n] \\ &\leq (1 - \mu\gamma_{n+1})\delta_n + \gamma_{n+1}^2 \mathbb{E}[\|\nabla f_{n+1}(\theta_n) - \nabla f(\theta^*)\|^2 | \mathcal{F}_n].\end{aligned}$$

- Since  $\nabla f_{n+1}$  is a.s. Lipschitz with bounded variance at  $\theta^*$ ,

$$\begin{aligned}&\mathbb{E} \left[ \|\nabla f_{n+1}(\theta_n) - \nabla f(\theta^*)\|^2 \mid \mathcal{F}_n \right] \\ &\leq \mathbb{E} \left[ \|\nabla f_{n+1}(\theta_n) - \nabla f_{n+1}(\theta_*) + \nabla f_{n+1}(\theta_*) - \nabla f(\theta^*)\|^2 \mid \mathcal{F}_n \right] \\ &\leq 2(L^2\delta_n + \sigma^2).\end{aligned}$$

## Sketch of proof

$f$  strongly convex,  $f_n$  smooth, bounded variance

- Consider  $\delta_n = \|\theta_n - \theta^*\|^2$ .
- Since  $\nabla f_{n+1}$  is a.s. Lipschitz with bounded variance at  $\theta^*$ ,

$$\begin{aligned} & \mathbb{E} \left[ \|\nabla f_{n+1}(\theta_n) - \nabla f(\theta^*)\|^2 \mid \mathcal{F}_n \right] \\ & \leq \mathbb{E} \left[ \|\nabla f_{n+1}(\theta_n) - \nabla f_{n+1}(\theta_*) + \nabla f_{n+1}(\theta_*) - \nabla f(\theta^*)\|^2 \mid \mathcal{F}_n \right] \\ & \leq 2(L^2 \delta_n + \sigma^2). \end{aligned}$$

- Conclusion:

$$\mathbb{E}[\delta_{n+1} | \mathcal{F}_n] \leq (1 - \mu\gamma_{n+1} + 2L^2\gamma_{n+1}^2)\delta_n + 2\sigma^2\gamma_{n+1}^2.$$

# Sketch of proof

$f$  strongly convex,  $f_n$  smooth, bounded variance

- Consider  $\delta_n = \|\theta_n - \theta^*\|^2$ .
- Conclusion:

$$\mathbb{E}[\delta_{n+1} | \mathcal{F}_n] \leq (1 - \mu\gamma_{n+1} + 2L^2\gamma_{n+1}^2)\delta_n + 2\sigma^2\gamma_{n+1}^2.$$



# Convex Stochastic Approximation: take home message

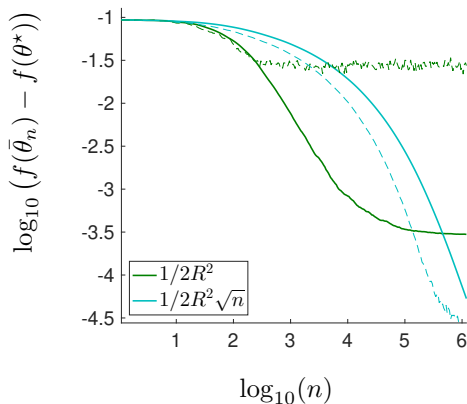
## ■ Pros

- Simple to implement
- Cheap
- No regularization needed
- Convergence guarantees

## ■ Cons:

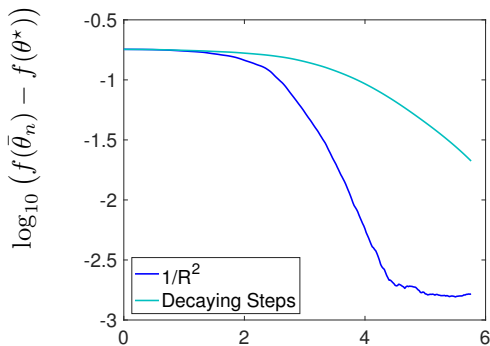
- Initial conditions can be forgotten slowly: could we use even larger/fixed step sizes?
- For fixed step sizes, the previous bounds do not show that  $\mathbb{E}[\|\theta_n - \theta^*\|^2] \not\rightarrow 0$  or  $\mathbb{E}[\|\bar{\theta}_n - \theta^*\|^2] \not\rightarrow 0$ .
- We only have  $\mathbb{E}[\|\theta_n - \theta^*\|^2] = O(\gamma)$  and  $\mathbb{E}[\|\bar{\theta}_n - \theta^*\|^2] = O(\gamma)$ .
- We illustrate these two facts using numerical simulations

## Motivation 1/ 2. Large step sizes!



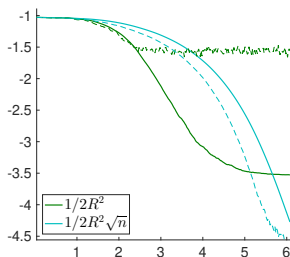
Logistic regression. Final iterate (dashed), and averaged recursion (plain).

## Motivation 1/ 2. Large step sizes, real data



Logistic regression, Covertypе dataset,  $n = 581012$ ,  $d = 54$ . Comparison between a constant learning rate and decaying learning rate as  $\frac{1}{\sqrt{n}}$ .

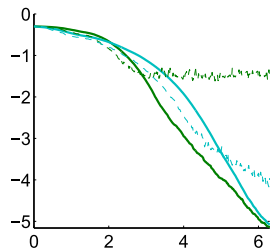
## Motivation 2/ 2. Difference between quadratic and logistic loss



Logistic Regression

$$\mathbb{E}f(\bar{\theta}_n) - f(\theta^*) = O(\gamma^2)$$

$$\text{with } \gamma = 1/(2R^2)$$



Least-Squares Regression

$$\mathbb{E}f(\bar{\theta}_n) - f(\theta^*) = O\left(\frac{1}{n}\right)$$

$$\text{with } \gamma = 1/(2R^2)$$

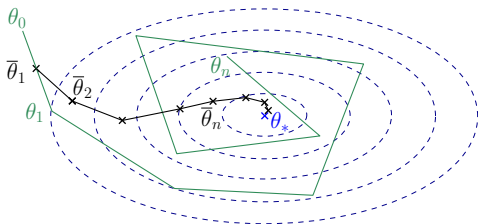
# Constant learning rate SGD: convergence in the quadratic case

**Least-squares:**  $f(\theta) = \frac{1}{2}\mathbb{E}[(Y - \langle \Phi(X), \theta \rangle)^2]$  with  $\theta \in \mathbb{R}^d$

- SGD = least-mean-square algorithm
- With strong convexity assumption  $\mathbb{E}[\Phi(X) \otimes \Phi(X)] = H \succcurlyeq \mu \cdot \text{Id}$

$$\theta^* = H^{-1}\mathbb{E}[Y\Phi(X)]$$

- $\bar{\theta}_n \rightarrow \theta^*$  as  $n \rightarrow +\infty$



# Constant learning rate SGD: convergence in the quadratic case

- Key identity:

$$\theta_{n+1} - \theta^* = (\text{Id} - \gamma H)(\theta_n - \theta^*) + \gamma \eta_{n+1}(\theta_n), \quad \mathbb{E}[\eta_{n+1}(\theta_n) | \mathcal{F}_n] = 0,$$

$$\eta_{n+1}(\theta) = H\theta - \mathbb{E}[Y\Phi(X)] - \Phi(X_{n+1})\Phi(X_{n+1})^\top \theta + Y_{n+1}\Phi(X_{n+1}).$$

- Therefore,

$$\theta_{n+1} - \theta^* = (\text{Id} - \gamma H)^{n+1}(\theta_0 - \theta^*) + \gamma \sum_{k=0}^n (\text{Id} - \gamma H)^{n-k} \eta_{k+1}(\theta_k),$$

and

$$\bar{\theta}_n - \theta^* = (n+1)^{-1} \sum_{k=0}^n (\theta_k - \theta^*) \approx (n+1)^{-1} H^{-1} \sum_{k=0}^n \eta_k(\theta_k).$$

# Constant learning rate SGD: convergence in the quadratic case

**Least-squares:**  $f(\theta) = \frac{1}{2}\mathbb{E}[(Y - \langle \Phi(X), \theta \rangle)^2]$  with  $\theta \in \mathbb{R}^d$

$$\theta_{n+1} - \theta^* = (\text{Id} - \gamma H)(\theta_n - \theta^*) + \gamma \eta_{n+1}(\theta_n),$$

- The sequence  $(\theta_n)_{n \geq 0}$  is a homogeneous Markov chain
  - 1 The distribution of  $(\theta_n)_{n \geq 0}$  converges to a stationary distribution  $\pi_\gamma$
  - 2  $\bar{\theta}_n$  converges to  $\bar{\theta}_\gamma = \int_{\mathbb{R}^d} \vartheta d\pi_\gamma(\vartheta)$  (Birkhoff theorem)
- Identification of  $\bar{\theta}_\gamma$ 
  - If  $\theta_0 \sim \pi_\gamma$ , then  $\theta_1 \sim \pi_\gamma$ .
  - Taking expectation, and using  $\mathbb{E}[\eta_1(\theta)] = 0$  for any  $\theta \in \mathbb{R}^d$ ,

$$\int_{\mathbb{R}^d} H(\vartheta - \theta^*) d\pi_\gamma(\vartheta) = 0 \Rightarrow \bar{\theta}_\gamma = \theta^* .$$

- Conclusion  $\bar{\theta}_n \rightarrow \theta^*$  as  $n \rightarrow +\infty$  if ergodic
- **Question:** What happens in the general case?

# SGD: an homogeneous Markov chain

- Consider a  $L$ -smooth and  $\mu$ -strongly convex function  $f$ .
- SGD with a step-size  $\gamma > 0$  is an **homogeneous Markov chain**:

$$\begin{aligned}\theta_{k+1}^\gamma &= \theta_k^\gamma - \gamma \nabla f_{k+1}(\theta_k^\gamma) = \theta_k^\gamma - \gamma [\nabla f(\theta_k^\gamma) + \eta_{k+1}(\theta_k^\gamma)] , \\ \eta_{k+1}(\theta_k^\gamma) &= \nabla f_{k+1}(\theta_k^\gamma) - \nabla f(\theta_k^\gamma) , \mathbb{E}[\eta_{k+1}(\theta_k^\gamma) | \mathcal{F}_k] = 0 .\end{aligned}$$

## Assumptions

- $\nabla f_k$  is almost surely  $L$ -co-coercive: for any  $\theta_1, \theta_2 \in \mathbb{R}^d$ ,
$$\langle \nabla f_k(\theta_1) - \nabla f_k(\theta_2), \theta_1 - \theta_2 \rangle \geq L^{-1} \|\nabla f_k(\theta_1) - \nabla f_k(\theta_2)\|^2 .$$
- Bounded moments for  $p$  large enough,

$$\mathbb{E}[\|\eta_k(\theta^*)\|^p] < \infty .$$



# Stochastic gradient descent as a Markov Chain: Analysis framework<sup>2</sup>

- Let  $R_\gamma$  be the Markov kernel associated with  $(\theta_n^\gamma)_{n \in \mathbb{N}}$ .
- Existence of a stationary distribution  $\pi_\gamma$  for  $R_\gamma$ , and convergence to this distribution.
- Behavior under the limit distribution ( $\gamma \rightarrow 0$ ):  $\bar{\theta}_\gamma = \theta^* + ?$   
↪ Provable convergence improvement with extrapolation tricks used for numerical integration and applied probability.
- Analysis of the convergence of  $\bar{\theta}_n^\gamma$  to  $\bar{\theta}_\gamma = \int_{\mathbb{R}^d} \vartheta d\pi_\gamma(\vartheta)$  through its MSE.

<sup>2</sup>Bach, Dieuleveut, Durmus, AOS, 2020.

# Existence and convergence to a stationary distribution

## Definition

Wasserstein distance:  $\nu$  and  $\lambda$  probability measures on  $\mathbb{R}^d$

$$W_2(\lambda, \nu) := \inf_{\xi \in \Pi(\lambda, \nu)} \left( \int \|\theta - \eta\|^2 \xi(d\theta \cdot d\eta) \right)^{1/2}$$

$\Pi(\lambda, \nu)$  is the set of probability measure  $\xi$  s.t.  $A \in \mathcal{B}(\mathbb{R}^d)$ ,  
 $\xi(A \times \mathbb{R}^d) = \lambda(A)$ ,  $\xi(\mathbb{R}^d \times A) = \nu(A)$ .

## Theorem

For  $\gamma < L^{-1}$ , the chain  $(\theta_k^\gamma)_{k \geq 0}$  admits a unique stationary distribution  $\pi_\gamma$  and for all  $\theta \in \mathbb{R}^d$ ,  $n \in \mathbb{N}$ :

$$W_2^2(\delta_\theta R_\gamma^n, \pi_\gamma) \leq (1 - 2\mu\gamma(1 - \gamma L))^n \int_{\mathbb{R}^d} \|\theta - \vartheta\|^2 d\pi_\gamma(\vartheta).$$

## Existence of a limit distribution: proof I /III

- **Coupling:**  $\theta^1, \theta^2$  be independent and distributed according to  $\lambda_1, \lambda_2$  respectively, and  $(\theta_{k,\gamma}^{(1)})_{k \geq 0}, (\theta_{k,\gamma}^{(2)})_{k \geq 0}$  SGD iterates:

$$\begin{cases} \theta_{k+1,\gamma}^{(1)} &= \theta_{k,\gamma}^{(1)} - \gamma [\nabla f(\theta_{k,\gamma}^{(1)}) + \eta_{k+1}(\theta_{k,\gamma}^{(1)})] \\ \theta_{k+1,\gamma}^{(2)} &= \theta_{k,\gamma}^{(2)} - \gamma [\nabla f(\theta_{k,\gamma}^{(2)}) + \eta_{k+1}(\theta_{k,\gamma}^{(2)})] \end{cases} .$$

- for all  $k \geq 0$ , the distribution of  $(\theta_{k,\gamma}^{(1)}, \theta_{k,\gamma}^{(2)})$  is in  $\Pi(\lambda_1 R_\gamma^k, \lambda_2 R_\gamma^k)$

# Existence of a limit distribution: proof II/III

$$\begin{aligned} & \mathbb{E} \left[ \|\theta_{k+1,\gamma}^{(1)} - \theta_{k+1,\gamma}^{(2)}\|^{(2)} \right] \\ & \leq \mathbb{E} \left[ \|\theta_{k,\gamma}^{(1)} - \gamma \nabla f_{k+1}(\theta_{k,\gamma}^{(1)}) - (\theta_{k,\gamma}^{(2)} - \gamma \nabla f_{k+1}(\theta_{k,\gamma}^{(2)}))\|^2 \right] \end{aligned}$$

## Existence of a limit distribution: proof II/III

$$\begin{aligned}
 & \mathbb{E} \left[ \left\| \theta_{k+1,\gamma}^{(1)} - \theta_{k+1,\gamma}^{(2)} \right\|^2 \right] \\
 & \leq \mathbb{E} \left[ \left\| \theta_{k,\gamma}^{(1)} - \gamma \nabla f_{k+1}(\theta_{k,\gamma}^{(1)}) - (\theta_{k,\gamma}^{(2)} - \gamma \nabla f_{k+1}(\theta_{k,\gamma}^{(2)})) \right\|^2 \right] \\
 & \leq \mathbb{E} \left[ \left\| \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\|^2 - 2\gamma \left\langle \nabla f_{k+1}(\theta_{k,\gamma}^{(1)}) - \nabla f_{k+1}(\theta_{k,\gamma}^{(2)}), \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\rangle \right] \\
 & \quad + \gamma^2 \mathbb{E} \left[ \left\| \nabla f_{k+1}(\theta_{k,\gamma}^{(1)}) - \nabla f_{k+1}(\theta_{k,\gamma}^{(2)}) \right\|^2 \right]
 \end{aligned}$$

## Existence of a limit distribution: proof II/III

$$\begin{aligned}
 & \mathbb{E} \left[ \left\| \theta_{k+1,\gamma}^{(1)} - \theta_{k+1,\gamma}^{(2)} \right\|^{(2)} \right] \\
 & \leq \mathbb{E} \left[ \left\| \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\|^2 - 2\gamma \left\langle \nabla f_{k+1}(\theta_{k,\gamma}^{(1)}) - \nabla f_{k+1}(\theta_{k,\gamma}^{(2)}), \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\rangle \right] \\
 & \quad + \gamma^2 \mathbb{E} \left[ \left\| \nabla f_{k+1}(\theta_{k,\gamma}^{(1)}) - \nabla f_{k+1}(\theta_{k,\gamma}^{(2)}) \right\|^2 \right] \\
 & \stackrel{\text{coco}}{\leq} \mathbb{E} \left[ \left\| \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\|^2 \right] \\
 & \quad - 2\gamma(1 - \gamma L) \mathbb{E} \left[ \left\langle \nabla f_{k+1}(\theta_{k,\gamma}^{(1)}) - \nabla f_{k+1}(\theta_{k,\gamma}^{(2)}), \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\rangle \right]
 \end{aligned}$$

# Existence of a limit distribution: proof II/III

$$\begin{aligned}
 & \mathbb{E} \left[ \left\| \theta_{k+1,\gamma}^{(1)} - \theta_{k+1,\gamma}^{(2)} \right\|^{(2)} \right] \\
 & \stackrel{\text{coco}}{\leq} \mathbb{E} \left[ \left\| \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\|^2 \right] \\
 & \quad - 2\gamma(1 - \gamma L) \mathbb{E} \left[ \left\langle \nabla f_{k+1}(\theta_{k,\gamma}^{(1)}) - \nabla f_{k+1}(\theta_{k,\gamma}^{(2)}), \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\rangle \right] \\
 & \stackrel{\text{unbiased}}{\leq} \mathbb{E} \left[ \left\| \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\|^2 \right] \\
 & \quad - 2\gamma(1 - \gamma L) \mathbb{E} \left[ \left\langle \nabla f(\theta_{k,\gamma}^{(1)}) - \nabla f(\theta_{k,\gamma}^{(2)}), \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\rangle \right]
 \end{aligned}$$

## Existence of a limit distribution: proof II/III

$$\begin{aligned}
 & \mathbb{E} \left[ \left\| \theta_{k+1,\gamma}^{(1)} - \theta_{k+1,\gamma}^{(2)} \right\|^2 \right] \\
 & \stackrel{\text{unbiased}}{\leq} \mathbb{E} \left[ \left\| \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\|^2 \right] \\
 & \quad - 2\gamma(1 - \gamma L) \mathbb{E} \left[ \left\langle \nabla f(\theta_{k,\gamma}^{(1)}) - \nabla f(\theta_{k,\gamma}^{(2)}), \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\rangle \right] \\
 & \stackrel{\text{s.cvx.}}{\leq} (1 - 2\mu\gamma(1 - \gamma L)) \mathbb{E} \left[ \left\| \theta_{k,\gamma}^{(1)} - \theta_{k,\gamma}^{(2)} \right\|^2 \right].
 \end{aligned}$$



## Existence of a limit distribution: proof III/III

- By induction:

$$\begin{aligned} W_2^2(\lambda_1 R_\gamma^n, \lambda_2 R_\gamma^n) &\leq \mathbb{E} \left[ \|\theta_{n,\gamma}^{(1)} - \theta_{n,\gamma}^{(2)}\|^2 \right] \\ &\leq (1 - 2\mu\gamma(1 - \gamma L))^n \int_{x,y} \|\theta_1 - \theta_2\|^2 d\lambda_1(\theta_1)d\lambda_2(\theta_2). \end{aligned}$$

- Thus  $W_2(\delta_{\theta_1} R_\gamma^n, \delta_{\theta_2} R_\gamma^n) \leq (1 - 2\mu\gamma(1 - \gamma L))^n \|\theta_1 - \theta_2\|^2$ .
- Uniqueness, invariance, and Theorem follow:

$$W_2^2(\delta_\theta R_\gamma^n, \pi_\gamma) \leq (1 - 2\mu\gamma(1 - \gamma L))^n \int_{\mathbb{R}^d} \|\theta - \vartheta\|^2 d\pi_\gamma(\vartheta).$$

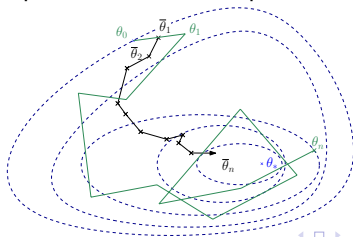
## Behavior under limit distribution.

- Then we have  $\mathbb{E}[\bar{\theta}_n] \rightarrow \bar{\theta}_\gamma$ . Where is  $\bar{\theta}_\gamma$ ? Close to  $\theta^*$ ?
- In the quadratic case  $\bar{\theta}_\gamma = \theta^*$
- In the general case, we show that

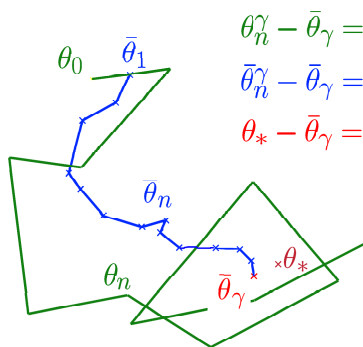
$$\bar{\theta}_\gamma = \theta^* + \gamma \Delta(\theta^*) + O(\gamma^2)$$

$$\Delta(\theta^*) = f''(\theta^*)^{-1} f'''(\theta^*) \left( [f''(\theta^*) \otimes I + I \otimes f''(\theta^*)]^{-1} \mathbb{E}[\eta(\theta^*)^{\otimes 2}] \right).$$

- Linearization of the proof for the least-square



# Richardson extrapolation



$$\theta_n^\gamma - \bar{\theta}_\gamma = O_p(\gamma^{1/2})$$

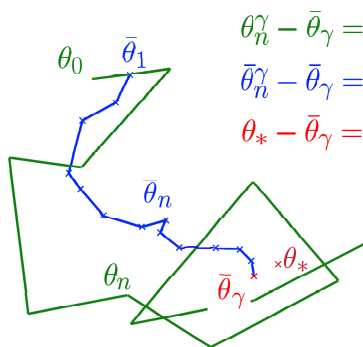
$$\bar{\theta}_n^\gamma - \bar{\theta}_\gamma = O_p(n^{-1/2})$$

$$\theta_* - \bar{\theta}_\gamma = O(\gamma)$$

$\bullet \theta_*$

$\bullet \leftarrow \theta_* + \gamma \Delta$

# Richardson extrapolation



$$\theta_n^\gamma - \bar{\theta}_\gamma = O_p(\gamma^{1/2})$$

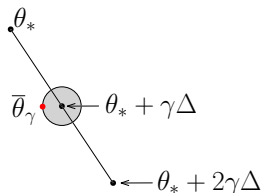
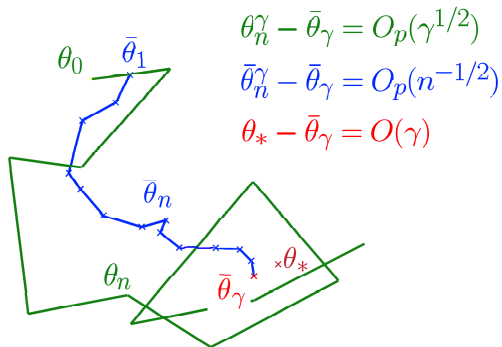
$$\bar{\theta}_n^\gamma - \bar{\theta}_\gamma = O_p(n^{-1/2})$$

$$\theta_* - \bar{\theta}_\gamma = O(\gamma)$$

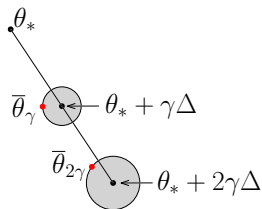
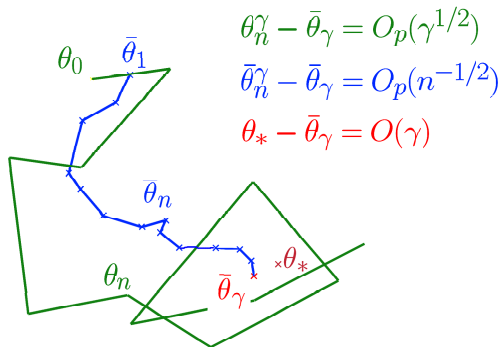
$\bullet \theta_*$

$$\bar{\theta}_\gamma \bullet \leftarrow \theta_* + \gamma \Delta$$

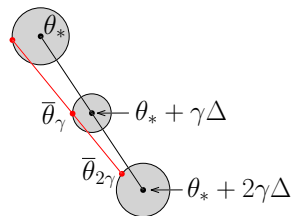
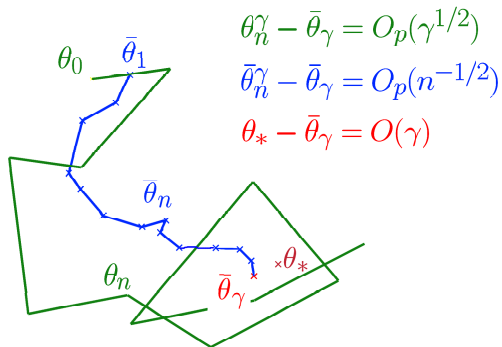
# Richardson extrapolation



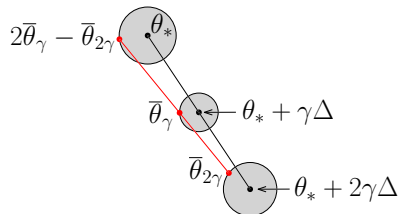
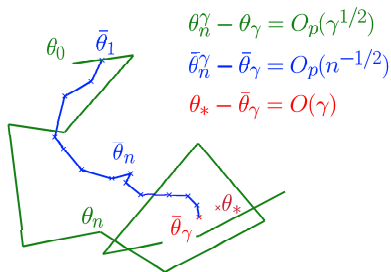
# Richardson extrapolation



# Richardson extrapolation



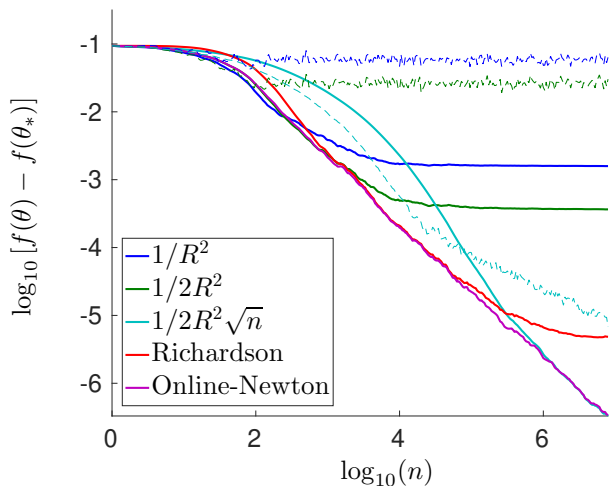
# Richardson extrapolation



Recovering convergence closer to  $\theta_*$  by **Richardson extrapolation**  
 $2\bar{\theta}_n^\gamma - \bar{\theta}_n^{2\gamma}$

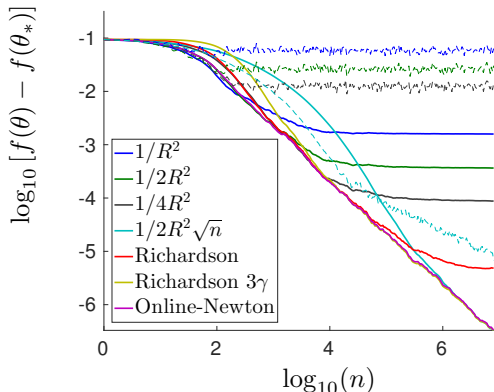


# Experiments



Synthetic data, logistic regression,  $n = 8.10^6$

## Experiments: Double Richardson



Synthetic data, logistic regression,  $n = 8.10^6$

“Richardson  $3\gamma$ ”: estimator built using *Richardson on 3 different*

*sequences*:  $\tilde{\theta}_n^3 = \frac{8}{3}\bar{\theta}_n^\gamma - 2\bar{\theta}_n^{2\gamma} + \frac{1}{3}\bar{\theta}_n^{4\gamma}$

## Real data

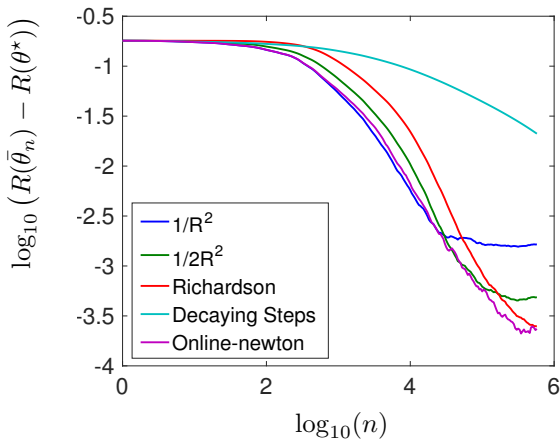


Figure: Logistic regression, Covertyp dataset.  $n = 581012$ ,  $d = 54$ .

- 1 Stochastic approximation
  - Finite-sum optimization
  - Online learning
  - Smooth strongly convex case
  - Stochastic subgradient descent/method
  - Stochastic Approximation for nonconvex optimization
- 2 Proximal methods
  - Proximal operator
  - Proximal gradient algorithm
  - Stochastic proximal gradient
- 3 Applications
  - Network structure estimation
  - High-dimensional logistic regression with random effect

# Stochastic subgradient descent/method

## Assumptions

- $f_n$  convex and  $B$ -Lipschitz-continuous on  $\{\|\theta\|_2 \leq D\}$
- $(f_n)$  i.i.d. functions such that  $\mathbb{E}[f_n(\theta)] = f(\theta)$
- $\theta_*$  global optimum of  $f$  on  $\{\|\theta\|_2 \leq D\}$

Algorithm:  $\theta_n = \Pi_D \left( \theta_{n-1} - \frac{2D}{B\sqrt{n}} \partial f_n(\theta_{n-1}) \right)$

Risk Bound:

$$\mathbb{E} \left[ f \left( \frac{1}{n} \sum_{k=0}^{n-1} \theta_k \right) \right] - f(\theta_*) \leq \frac{2DB}{\sqrt{n}}.$$

- Minimax convergence rate
- Running-time complexity:  $O(dn)$  after  $n$  iterations

## Stochastic subgradient method - proof - I

$\theta_n = \Pi_D(\theta_{n-1} - \gamma_n \partial f_n(\theta_{n-1}))$  where  $\mathcal{F}_n = \sigma((Y_k, X_k), j \leq n)$ .

$$\begin{aligned} \|\theta_n - \theta_*\|_2^2 &\leq \|\theta_{n-1} - \theta_* - \gamma_n \partial f_n(\theta_{n-1})\|_2^2 && \text{contractivity of projections} \\ &\leq \|\theta_{n-1} - \theta_*\|_2^2 + B^2 \gamma_n^2 - 2\gamma_n \langle \theta_{n-1} - \theta_*, \partial f_n(\theta_{n-1}) \rangle && \|\partial f_n(\theta_{n-1})\|_2 \leq B \end{aligned}$$

Taking the conditional expectations of the both sides

$$\begin{aligned} \mathbb{E}[\|\theta_n - \theta_*\|_2^2 | \mathcal{F}_{n-1}] &\leq \|\theta_{n-1} - \theta_*\|_2^2 + B^2 \gamma_n^2 - 2\gamma_n \langle (\theta_{n-1} - \theta_*), \partial f(\theta_{n-1}) \rangle \\ &\leq \|\theta_{n-1} - \theta_*\|_2^2 + B^2 \gamma_n^2 - 2\gamma_n [f(\theta_{n-1}) - f(\theta^*)] \quad (\text{subgradient property}) \end{aligned}$$

## Stochastic subgradient method - proof - I

$\theta_n = \Pi_D(\theta_{n-1} - \gamma_n \partial f_n(\theta_{n-1}))$  where  $\mathcal{F}_n = \sigma((Y_k, X_k), j \leq n)$ .

From

$$\mathbb{E}[\|\theta_n - \theta_*\|_2^2 | \mathcal{F}_{n-1}] \leq \|\theta_{n-1} - \theta_*\|_2^2 + B^2 \gamma_n^2 - 2\gamma_n [f(\theta_{n-1}) - f(\theta^*)]$$

the tower property of conditional expectation implies

$$\mathbb{E}[\|\theta_n - \theta_*\|_2^2] \leq \mathbb{E}[\|\theta_{n-1} - \theta_*\|_2^2] + B^2 \gamma_n^2 - 2\gamma_n [\mathbb{E}[f(\theta_{n-1})] - f(\theta^*)]$$

leading to

$$\mathbb{E}[f(\theta_{n-1})] - f(\theta^*) \leq \frac{B^2 \gamma_n}{2} + \frac{1}{2\gamma_n} \left\{ \mathbb{E}[\|\theta_{n-1} - \theta_*\|_2^2] - \mathbb{E}[\|\theta_n - \theta_*\|_2^2] \right\}$$

# Stochastic subgradient

$$\mathbb{E}[f(\theta_{n-1})] - f(\theta^*) \leq \frac{B^2\gamma_n}{2} + \frac{1}{2\gamma_n} [\mathbb{E}\|\theta_{n-1} - \theta^*\|_2^2 - \mathbb{E}\|\theta_n - \theta^*\|_2^2]$$

Constant step size

$$\begin{aligned} \sum_{u=1}^n [\mathbb{E}[f(\theta_{u-1})] - f(\theta^*)] &\leq \sum_{u=1}^n \frac{B^2\gamma}{2} + \sum_{u=1}^n \frac{1}{2\gamma} \{ \mathbb{E} [\|\theta_{u-1} - \theta^*\|_2^2] - \mathbb{E} [\|\theta_u - \theta^*\|_2^2] \} \\ &\leq \frac{nB^2\gamma}{2} + \frac{4D^2}{2\gamma}. \end{aligned}$$

Optimum stepsize  $\gamma = 2D/(\sqrt{n}B)$  (depends on the horizon).



# Stochastic subgradient

$$\mathbb{E}[f(\theta_{n-1})] - f(\theta^*) \leq \frac{B^2\gamma_n}{2} + \frac{1}{2\gamma_n} [\mathbb{E}\|\theta_{n-1} - \theta_*\|_2^2 - \mathbb{E}\|\theta_n - \theta_*\|_2^2]$$

Constant step size

$$\begin{aligned} \sum_{u=1}^n [\mathbb{E}[f(\theta_{u-1})] - f(\theta^*)] &\leq \sum_{u=1}^n \frac{B^2\gamma}{2} + \sum_{u=1}^n \frac{1}{2\gamma} \{ \mathbb{E} [\|\theta_{u-1} - \theta^*\|_2^2] - \mathbb{E} [\|\theta_u - \theta^*\|_2^2] \} \\ &\leq \frac{nB^2\gamma}{2} + \frac{4D^2}{2\gamma}. \end{aligned}$$

Convexity [fixed horizon]:

$$\mathbb{E} \left[ f \left( \frac{1}{n} \sum_{k=0}^{n-1} \theta_k \right) \right] - f(\theta_*) \leq \frac{2DB}{\sqrt{n}}$$

## Beyond convergence in expectation

Convergence in expectation:  $\mathbb{E} \left[ f \left( n^{-1} \sum_{k=0}^{n-1} \theta_k \right) - f(\theta^*) \right] \leq \frac{2DB}{\sqrt{n}}$

High-probability bounds

- Markov inequality:  $\mathbb{P} \left( f \left( n^{-1} \sum_{k=0}^{n-1} \theta_k \right) - f(\theta^*) \geq \epsilon \right) \leq \frac{2DB}{\sqrt{n}\epsilon}$
- Concentration inequality (Nemirovski et al., 2009; Nesterov and Vial, 2008)

$$\mathbb{P} \left( f \left( \frac{1}{n} \sum_{k=0}^{n-1} \theta_k \right) - f(\theta^*) \geq \frac{2DB}{\sqrt{n}} (2 + 4t) \right) \leq 2 \exp(-t^2)$$

# Stochastic subgradient method - proof - I

$$\theta_n = \Pi_D(\theta_{n-1} - \gamma_n \partial f_n(\theta_{n-1})) \text{ with } \mathcal{F}_n = \sigma((Y_k, X_k), j \leq n).$$

$$\begin{aligned} \|\theta_n - \theta_*\|_2^2 &\leq \|\theta_{n-1} - \theta_* - \gamma_n \partial f_n(\theta_{n-1})\|_2^2 && \text{contractivity of projections} \\ &\leq \|\theta_{n-1} - \theta_*\|_2^2 + B^2 \gamma_n^2 - 2\gamma_n \langle \theta_{n-1} - \theta_*, \partial f_n(\theta_{n-1}) \rangle && \|\partial f_n(\theta_{n-1})\|_2 \leq B \end{aligned}$$

Define by  $Z_n$  the error (approximation of the "true" subgradient by its noisy version)

$$Z_n = -2 \langle \theta_{n-1} - \theta^*, \partial f_n(\theta_{n-1}) - \partial f(\theta_{n-1}) \rangle$$

and using the convexity we get

$$\|\theta_n - \theta^*\|_2^2 \leq \|\theta_{n-1} - \theta^*\|_2^2 + B^2 \gamma_n^2 - 2\gamma_n [f(\theta_{n-1}) - f(\theta^*)] + 2\gamma_n Z_n$$

## Stochastic subgradient method - proof - II

$$Z_n = -\langle \theta_{n-1} - \theta^*, \partial f_n(\theta_{n-1}) - \partial f(\theta_{n-1}) \rangle$$

From the inequality

$$\|\theta_n - \theta^*\|_2^2 \leq \|\theta_{n-1} - \theta^*\|_2^2 + B^2 \gamma_n^2 - 2\gamma_n [f(\theta_{n-1}) - f(\theta^*)] + 2\gamma_n Z_n$$

we get

$$f(\theta_{n-1}) - f(\theta^*) \leq \frac{1}{2\gamma_n} \{ \|\theta_{n-1} - \theta^*\|_2^2 - \|\theta_n - \theta^*\|_2^2 \} + \frac{B^2 \gamma_n}{2} + Z_n$$

Summing up this identity

$$\sum_{u=1}^n [f(\theta_{u-1}) - f(\theta^*)] \leq \sum_{u=1}^n \frac{B^2 \gamma_u}{2} + \sum_{u=1}^n \frac{1}{2\gamma_u} \{ \|\theta_{u-1} - \theta^*\|_2^2 - \|\theta_u - \theta^*\|_2^2 \} + \sum_{u=1}^n Z_u$$

## Stochastic subgradient method - proof - II

$$Z_n = -\langle \theta_{n-1} - \theta^*, \partial f_n(\theta_{n-1}) - \partial f(\theta_{n-1}) \rangle$$

Setting  $\gamma_u = 2D/(B\sqrt{n})$  [depending on the horizon  $n$ ] in

$$\sum_{u=1}^n [f(\theta_{u-1}) - f(\theta^*)] \leq \sum_{u=1}^n \frac{B^2 \gamma_u}{2} + \sum_{u=1}^n \frac{1}{2\gamma_u} \{ \|\theta_{u-1} - \theta^*\|_2^2 - \|\theta_u - \theta^*\|_2^2 \} + \sum_{u=1}^n Z_u$$

we get

$$\frac{1}{n} \sum_{u=1}^n \{ f(\theta_{u-1}) - f(\theta^*) \} \leq \frac{2DB}{\sqrt{n}} + \frac{1}{n} \sum_{u=1}^n Z_u$$

## Stochastic subgradient method - proof - II

$$Z_n = -\langle \theta_{n-1} - \theta^*, \partial f_n(\theta_{n-1}) - \partial f(\theta_{n-1}) \rangle$$

Setting  $\gamma_u = 2D/(B\sqrt{n})$  [depending on the horizon  $n$ ] in

$$\sum_{u=1}^n [f(\theta_{u-1}) - f(\theta^*)] \leq \sum_{u=1}^n \frac{B^2 \gamma_u}{2} + \sum_{u=1}^n \frac{1}{2\gamma_u} \{ \|\theta_{u-1} - \theta^*\|_2^2 - \|\theta_u - \theta^*\|_2^2 \} + \sum_{u=1}^n Z_u$$

we get

$$\frac{1}{n} \sum_{u=1}^n \{ f(\theta_{u-1}) - f(\theta^*) \} \leq \frac{2DB}{\sqrt{n}} + \frac{1}{n} \sum_{u=1}^n Z_u$$

Require to study  $n^{-1} \sum_{k=1}^n Z_k$  where  $(Z_k)_{k \geq 1}$  is a bounded martingale increment sequence:  $|Z_k| \leq 4DB$ .

## Stochastic subgradient method - proof - II

$$Z_n = -\langle \theta_{n-1} - \theta^*, \partial f_n(\theta_{n-1}) - \partial f(\theta_{n-1}) \rangle$$

Setting  $\gamma_u = 2D/(B\sqrt{n})$  [depending on the horizon  $n$ ] in

$$\sum_{u=1}^n [f(\theta_{u-1}) - f(\theta^*)] \leq \sum_{u=1}^n \frac{B^2 \gamma_u}{2} + \sum_{u=1}^n \frac{1}{2\gamma_u} \{ \|\theta_{u-1} - \theta^*\|_2^2 - \|\theta_u - \theta^*\|_2^2 \} + \sum_{u=1}^n Z_u$$

we get

$$\frac{1}{n} \sum_{u=1}^n \{ f(\theta_{u-1}) - f(\theta^*) \} \leq \frac{2DB}{\sqrt{n}} + \frac{1}{n} \sum_{u=1}^n Z_u$$

Azuma-Hoeffding inequality for bounded martingale increments:

$$\mathbb{P} \left( \frac{1}{n} \sum_{u=1}^n Z_u \geq \frac{4DBt}{\sqrt{n}} \right) \leq \exp(-t^2/2)$$

## Stochastic subgradient method - proof - II

$$Z_n = -\langle \theta_{n-1} - \theta^*, \partial f_n(\theta_{n-1}) - \partial f(\theta_{n-1}) \rangle$$

Setting  $\gamma_u = 2D/(B\sqrt{n})$  [depending on the horizon  $n$ ] in

$$\sum_{u=1}^n [f(\theta_{u-1}) - f(\theta^*)] \leq \sum_{u=1}^n \frac{B^2 \gamma_u}{2} + \sum_{u=1}^n \frac{1}{2\gamma_u} \{ \|\theta_{u-1} - \theta^*\|_2^2 - \|\theta_u - \theta^*\|_2^2 \} + \sum_{u=1}^n Z_u$$

we get

$$\frac{1}{n} \sum_{u=1}^n \{f(\theta_{u-1}) - f(\theta^*)\} \leq \frac{2DB}{\sqrt{n}} + \frac{1}{n} \sum_{u=1}^n Z_u$$

Moment bounds can be deduced from Burkholder-Rosenthal-Pinelis inequality (Pinelis, 1994)



- 1 Stochastic approximation**
  - Finite-sum optimization
  - Online learning
  - Smooth strongly convex case
  - Stochastic subgradient descent/method
  - Stochastic Approximation for nonconvex optimization
- 2 Proximal methods**
  - Proximal operator
  - Proximal gradient algorithm
  - Stochastic proximal gradient
- 3 Applications**
  - Network structure estimation
  - High-dimensional logistic regression with random effect

# Stochastic approximation beyond convex optimization

- Stochastic approximation goes far beyond convex optimization.
- **Problem:** find the roots of the **mean field** function  $h$ , i.e. solve  $h(\theta) = 0$ .
- **Stochastic gradient:**  $h = \nabla f$ .
- The function  $h$  is not known in closed form, but

$$h(\theta) = \int H(\theta, x) \nu(dx)$$

where  $H : \Theta \times X \rightarrow \Theta$  is a known function and  $\nu$  is a probability distribution over  $X$ .

## Robbins Monro set up

- Assume that there is an i.i.d. sequence  $\{X_n, n \in \mathbb{N}\}$  distributed according to  $\nu$
- The **stochastic approximation** procedure:

$$\theta_n = \theta_{n-1} + \gamma_n H(\theta_{n-1}, X_n) \text{ with } \mathbb{E}[h_n(\theta_{n-1}) | \mathcal{F}_{n-1}] = h(\theta_{n-1})$$

where  $\mathcal{F}_{n-1}$  is the  $\sigma$ -algebra of summarizing "past" observations.

- Can alternatively be written

$$\theta_n = \theta_{n-1} + \gamma_n h(\theta_{n-1}) + \gamma_n M_n$$

where  $M_n = H(\theta_{n-1}, X_n) - h(\theta_n)$ .

- Under the stated assumptions,  $\mathbb{E}[M_n | \mathcal{F}_{n-1}] = 0$ , i.e. the sequence  $\{M_n, n \in \mathbb{N}\}$  is a **martingale increment** sequence.

## Limiting ODE

- The limiting ODE which the SA procedure might be expected to track is  $\dot{\theta} = h(\theta)$
- In absence of noise ( $M_n \equiv 0$ ), the recursion

$$\theta_n = \theta_{n-1} + \gamma_n h(\theta_n)$$

is the Euler discretization of the ODE  $\dot{\theta} = h(\theta)$  with stepsize  $\{\gamma_n, n \in \mathbb{N}\}$ .

- Many asymptotic convergence results (see Kushner and Yin (2003), Borkar (2008)) but few **quantitative** results.

## Randomized Stochastic Gradient (RGSD) Method

Stochastic oracle: for  $\theta \in \mathbb{R}^d$ ,

- Unbiasedness  $\mathbb{E}[G(\theta, \xi)] = \nabla f(\theta)$
- Bounded variance  $\mathbb{E}[\|G(\theta, \xi) - \nabla f(\theta)\|^2] \leq \sigma^2$

Stochastic gradient:

- Initial point  $\theta_0$ , iteration limit  $N$ , stepsizes  $\{\gamma_k\}_{k=0}^{N-1}$  and probability over  $\Pi$  on  $\{0, \dots, N\}$
- step 0. Draw  $R$  from  $\Pi$
- step 1. for  $k \in \{1, \dots, R\}$ , call the stochastic oracle  $G(\theta_{k-1}, \xi_k)$  and set

$$\theta_k = \theta_{k-1} - \gamma_k G(\theta_{k-1}, \xi_k)$$

# RSGD convergence

## Theorem

Suppose that the stepsizes  $\{\gamma_k\}$  and the probability  $\Pi$  satisfies,  $\gamma_k \leq 1/2L$  and,

$$\Pi(k) := \frac{2\gamma_{k+1} - L\gamma_{k+1}^2}{\sum_{k=1}^N (2\gamma_k - L\gamma_k^2)}, k = 0, \dots, N - 1$$

For any  $N \geq 1$ , we have

$$\frac{1}{L} \mathbb{E} \left[ \|\nabla f(\theta_R)\|^2 \right] \leq \frac{D_f^2 + \sigma^2 \sum_{k=1}^N \gamma_k^2}{\sum_{k=0}^{N-1} (2\gamma_k - L\gamma_k^2)}$$

where, denoting  $f^*$  denotes the optimal value,

$$D_f := [2(f(\theta_1) - f^*) / L]^{1/2}$$

# RSGD convergence

## Theorem

Suppose that the stepsizes  $\{\gamma_k\}$  and the probability  $\Pi$  satisfies,  
 $\gamma_k \leq 1/2L$  and,

$$\Pi(k) := \frac{2\gamma_{k+1} - L\gamma_{k+1}^2}{\sum_{k=1}^N (2\gamma_k - L\gamma_k^2)}, k = 0, \dots, N - 1$$

If in addition  $f$  is convex with an optimal solution  $\theta^*$ , then for any  
 $N \geq 1$ ,

$$\mathbb{E}[f(\theta_R) - f^*] \leq \frac{D_Y^2 + \sigma^2 \sum_{k=1}^N \gamma_k^2}{\sum_{k=1}^N (2\gamma_k - L\gamma_k^2)^2}$$

where

$$D_X := \|\theta_1 - \theta^*\|$$

# RSGD convergence: Proof 1

Denote  $\delta_k \equiv G(\theta_{k-1}, \xi_k) - \nabla f(\theta_{k-1})$ ,  $k \geq 1$ . Then

$$\begin{aligned} f(\theta_k) &\leq f(\theta_{k-1}) + \langle \nabla f(\theta_{k-1}), \theta_k - \theta_{k-1} \rangle + \frac{L}{2} \gamma_k^2 \|G(\theta_{k-1}, \xi_k)\|^2 \\ &= f(\theta_{k-1}) - \gamma_k \langle \nabla f(\theta_k), G(\theta_{k-1}, \xi_k) \rangle + \frac{L}{2} \gamma_k^2 \|G(\theta_{k-1}, \xi_k)\|^2 \end{aligned}$$



# RSGD convergence: Proof 1

Denote  $\delta_k \equiv G(\theta_{k-1}, \xi_k) - \nabla f(\theta_{k-1})$ ,  $k \geq 1$ . Then

$$\begin{aligned} f(\theta_k) &\leq f(\theta_{k-1}) + \langle \nabla f(\theta_{k-1}), \theta_k - \theta_{k-1} \rangle + \frac{L}{2} \gamma_k^2 \|G(\theta_{k-1}, \xi_k)\|^2 \\ &= f(\theta_{k-1}) - \gamma_k \|\nabla f(\theta_{k-1})\|^2 - \gamma_k \langle \nabla f(\theta_{k-1}), \delta_k \rangle \\ &\quad + \frac{L}{2} \gamma_k^2 \left[ \|\nabla f(\theta_{k-1})\|^2 + 2 \langle \nabla f(\theta_{k-1}), \delta_k \rangle + \|\delta_k\|^2 \right] \end{aligned}$$

# RSGD convergence: Proof 1

Denote  $\delta_k \equiv G(\theta_{k-1}, \xi_k) - \nabla f(\theta_{k-1})$ ,  $k \geq 1$ . Then

$$\begin{aligned} f(\theta_k) &\leq f(\theta_{k-1}) + \langle \nabla f(\theta_{k-1}), \theta_k - \theta_{k-1} \rangle + \frac{L}{2} \gamma_k^2 \|G(\theta_{k-1}, \xi_k)\|^2 \\ &= f(\theta_{k-1}) - \left( \gamma_k - \frac{L}{2} \gamma_k^2 \right) \|\nabla f(\theta_{k-1})\|^2 \\ &\quad - (\gamma_k - L\gamma_k^2) \langle \nabla f(\theta_{k-1}), \delta_k \rangle + \frac{L}{2} \gamma_k^2 \|\delta_k\|^2 \end{aligned}$$

## RSGD convergence: Proof 2

Summing up the above inequality and rearranging terms

$$\begin{aligned} & \sum_{k=1}^N \left( \gamma_k - \frac{L}{2} \gamma_k^2 \right) \|\nabla f(\theta_{k-1})\|^2 \\ & \leq f(\theta_0) - f(\theta_N) - \sum_{k=1}^N (\gamma_k - L\gamma_k^2) \langle \nabla f(\theta_{k-1}), \delta_k \rangle + \frac{L}{2} \sum_{k=1}^N \gamma_k^2 \|\delta_k\|^2 \end{aligned}$$

## RSGD convergence: Proof 2

Summing up the above inequality and rearranging terms

$$\begin{aligned} & \sum_{k=1}^N \left( \gamma_k - \frac{L}{2} \gamma_k^2 \right) \|\nabla f(\theta_{k-1})\|^2 \\ & \leq f(\theta_0) - f^* - \sum_{k=1}^N (\gamma_k - L\gamma_k^2) \langle \nabla f(\theta_{k-1}), \delta_k \rangle + \frac{L}{2} \sum_{k=1}^N \gamma_k^2 \|\delta_k\|^2 \end{aligned}$$

## RS GD convergence: Proof 2

Summing up the above inequality and rearranging terms

$$\begin{aligned} & \sum_{k=1}^N \left( \gamma_k - \frac{L}{2} \gamma_k^2 \right) \|\nabla f(\theta_{k-1})\|^2 \\ & \leq f(\theta_0) - f^* - \sum_{k=1}^N (\gamma_k - L\gamma_k^2) \langle \nabla f(\theta_{k-1}), \delta_k \rangle + \frac{L}{2} \sum_{k=1}^N \gamma_k^2 \|\delta_k\|^2 \end{aligned}$$

Taking expectation and using  $\mathbb{E}[\|\delta_k\|^2] \leq \sigma^2$  we get

$$\sum_{k=1}^N \left( \gamma_k - \frac{L}{2} \gamma_k^2 \right) \mathbb{E} \|\nabla f(\theta_{k-1})\|^2 \leq f(\theta_0) - f^* + \frac{L\sigma^2}{2} \sum_{k=1}^N \gamma_k^2$$

## RSGD convergence: Proof 2

Summing up the above inequality and rearranging terms

$$\begin{aligned} & \sum_{k=1}^N \left( \gamma_k - \frac{L}{2} \gamma_k^2 \right) \|\nabla f(\theta_{k-1})\|^2 \\ & \leq f(\theta_0) - f^* - \sum_{k=1}^N (\gamma_k - L\gamma_k^2) \langle \nabla f(\theta_{k-1}), \delta_k \rangle + \frac{L}{2} \sum_{k=1}^N \gamma_k^2 \|\delta_k\|^2 \end{aligned}$$

Dividing both sides by  $L \sum_{k=1}^N (\gamma_k - L\gamma_k^2/2)$  we conclude

$$\frac{1}{L} \mathbb{E} \left[ \|\nabla f(\theta_R)\|^2 \right] \leq \frac{1}{\sum_{k=1}^N (2\gamma_k - L\gamma_k^2)} \left[ \frac{2(f(\theta_0) - f^*)}{L} + \sigma^2 \sum_{k=1}^N \gamma_k^2 \right]$$

1 Stochastic approximation

**2** Proximal methods

3 Applications

- 1 Stochastic approximation
  - Finite-sum optimization
  - Online learning
  - Smooth strongly convex case
  - Stochastic subgradient descent/method
  - Stochastic Approximation for nonconvex optimization
- 2 Proximal methods
  - Proximal operator
  - Proximal gradient algorithm
  - Stochastic proximal gradient
- 3 Applications
  - Network structure estimation
  - High-dimensional logistic regression with random effect



# Definition

## Definition (Proximal mapping)

$g$ : closed convex function;  $\gamma$ : stepsize

$$\text{prox}_{\gamma,g}(\theta) = \underset{\eta \in \Theta}{\text{argmin}} (g(\eta) + (2\gamma)^{-1} \|\eta - \theta\|_2^2)$$

- The **uniqueness** of the minimizer stems from the strong convexity of the function  $\eta \mapsto g(\eta) + 1/(2\gamma)\|\eta - \theta\|_2^2$
- If  $g = \mathbb{I}_{\mathcal{K}}$ , where  $\mathcal{K}$  is a closed convex set, then  $\text{prox}_{\gamma,g}$  is the Euclidean projection on  $\mathcal{K}$

$$\text{prox}_{\gamma,g}(\theta) = \underset{\eta \in \mathcal{K}}{\text{argmin}} \|\eta - \theta\|_2^2 = P_{\mathcal{K}}(\theta)$$

- The proximal operator may be seen as a generalisation of the projection on closed convex sets.

# Proximal operator

## Lemma

If  $\theta = (\theta_1, \theta_2, \dots, \theta_p)$  and  $g(\theta) = \sum_{i=1}^p g_i(\theta_i)$ , then

$$\text{prox}_{\gamma, g}(\theta) = (\text{prox}_{\gamma, g_1}(\theta_1), \text{prox}_{\gamma, g_2}(\theta_2), \dots, \text{prox}_{\gamma, g_p}(\theta_p))$$

# Proximal operator

## Lemma

If  $\theta = (\theta_1, \theta_2, \dots, \theta_p)$  and  $g(\theta) = \sum_{i=1}^p g_i(\theta_i)$ , then

$$\text{prox}_{\gamma, g}(\theta) = (\text{prox}_{\gamma, g_1}(\theta_1), \text{prox}_{\gamma, g_2}(\theta_2), \dots, \text{prox}_{\gamma, g_p}(\theta_p))$$

$$\begin{aligned} \operatorname{argmin}_{(\eta_1, \dots, \eta_p)} \left\{ \sum_{i=1}^p g_i(\eta_i) + 2\gamma^{-1} \sum_{i=1}^p \|\eta_i - \theta_i\|^2 \right\} \\ = \sum_{i=1}^p \operatorname{argmin}_{\eta_i} \{ g_i(\eta_i) + (2\gamma)^{-1} \|\eta_i - \theta_i\|^2 \} \end{aligned}$$

# A characterization of the proximal operator

## Theorem

Let  $g$  be a convex function on  $\Theta$ ,  $(\theta, p) \in \Theta^2$ ,

$$p = \text{prox}_{\gamma, g}(\theta) \iff \text{for all } \eta \in \Theta, \quad g(p) + \gamma^{-1} \langle \eta - p, \theta - p \rangle \leq g(\eta)$$

i.e.  $p$  is the unique element of  $\Theta$  satisfying  $\gamma^{-1}(\theta - p) \in \partial g(p)$ .

# A characterization of the proximal operator

## Theorem

Let  $g$  be a convex function on  $\Theta$ ,  $(\theta, p) \in \Theta^2$ ,

$$p = \text{prox}_{\gamma, g}(\theta) \iff \text{for all } \eta \in \Theta, \quad g(p) + \gamma^{-1} \langle \eta - p, \theta - p \rangle \leq g(\eta)$$

i.e.  $p$  is the unique element of  $\Theta$  satisfying  $\gamma^{-1}(\theta - p) \in \partial g(p)$ .

Follows also from the characterization of the subdifferential

$$p \text{ is the minimizer of } \eta \mapsto g(\eta) + (2\gamma)^{-1} \|\eta - \theta\|_2^2$$

$$\iff$$

$$0 \in \partial g(p) + \gamma^{-1}(p - \theta).$$

# Proximal operator: LASSO and Elastic net

- If  $g(\theta) = \sum_{i=1}^p \lambda_i |\theta_i|$  then  $\text{prox}_{\gamma, g}$  is shrinkage (soft threshold) operation

$$[S_{\lambda, \gamma}(\theta)]_i = \begin{cases} \theta_i - \gamma \lambda_i & \theta_i \geq \gamma \lambda_i \\ 0 & |\theta_i| \leq \gamma \lambda_i \\ \theta_i + \gamma \lambda_i & \theta_i \leq -\gamma \lambda_i \end{cases}$$

- If  $g(\theta) = \lambda ((1 - \alpha)/2 \|\theta\|_2^2 + \alpha \|\theta\|_1)$

$$(\text{Prox}_{\gamma, g}(\tau))_i = \frac{1}{1 + \gamma \lambda (1 - \alpha)} \begin{cases} \tau_i - \gamma \lambda \alpha & \text{if } \tau_i \geq \gamma \lambda \alpha \\ \tau_i + \gamma \lambda \alpha & \text{if } \tau_i \leq -\gamma \lambda \alpha \\ 0 & \text{otherwise} \end{cases}$$

# Fixed points of the proximal operator

## Theorem

Let  $g$  be a proper convex function on  $\Theta$ . The set of fixed points

$$\{\theta \in \Theta, \text{prox}_{\gamma, g}(\theta) = \theta\}$$

coincide with the set of global minimum of  $g$ .

# Fixed points of the proximal operator

## Theorem

Let  $g$  be a proper convex function on  $\Theta$ . The set of fixed points

$$\{\theta \in \Theta, \text{prox}_{\gamma, g}(\theta) = \theta\}$$

coincide with the set of global minimum of  $g$ .

- Characterization of the proximal point

$$\gamma^{-1}(\theta - \text{prox}_{\gamma, g}(\theta)) \in \partial g(\text{prox}_{\gamma, g}(\theta)).$$

- Sub-gradient: for all  $\eta \in \Theta$ ,

$$\gamma^{-1}\langle \eta - \text{prox}_{\gamma, g}(\theta), \theta - \text{prox}_{\gamma, g}(\theta) \rangle + g(\text{prox}_{\gamma, g}(\theta)) \leq g(\eta)$$

## Conclusion

$$\theta = \text{prox}_{\gamma, g}(\theta) \iff \text{for all } \eta \in \Theta, g(\text{prox}_{\gamma, g}(\theta)) \leq g(\eta).$$



# Firm non-expansiveness

## Theorem

If  $g$  is a proper convex function, then  $\text{prox}_{\gamma,g}$  and  $(\text{Id} - \text{prox}_{\gamma,g})$  are *firmly non-expansive* (or *co-coercive* with constant 1), i.e. for all  $\theta, \eta \in \Theta$ ,

$$\begin{aligned} \|p - q\|^2 + \|(\theta - p) - (\eta - q)\|^2 &\leq \|\theta - \eta\|^2, \\ \iff \langle p - q, \theta - \eta \rangle &\geq \|p - q\|^2. \end{aligned}$$

where  $p = \text{prox}_{\gamma,g}(\theta)$  and  $q = \text{prox}_{\gamma,g}(\eta)$ .

## Firm non-expansiveness

### Theorem

If  $g$  is a proper convex function, then  $\text{prox}_{\gamma,g}$  and  $(\text{Id} - \text{prox}_{\gamma,g})$  are *firmly non-expansive* (or *co-coercive with constant 1*), i.e. for all  $\theta, \eta \in \Theta$ ,

$$\begin{aligned} \|p - q\|^2 + \|(\theta - p) - (\eta - q)\|^2 &\leq \|\theta - \eta\|^2, \\ \iff \langle p - q, \theta - \eta \rangle &\geq \|p - q\|^2. \end{aligned}$$

where  $p = \text{prox}_{\gamma,g}(\theta)$  and  $q = \text{prox}_{\gamma,g}(\eta)$ .

$$\gamma^{-1} \langle q - p, \theta - p \rangle + g(p) \leq g(q) \quad \gamma^{-1} \langle p - q, \eta - q \rangle + g(q) \leq g(p)$$

Adding these two equations yield

$$\langle p - q, (\theta - p) - (\eta - q) \rangle \geq 0.$$

Conclude by writing  $\|\theta - \eta\|^2 = \|\theta - p + p - \eta + \eta - q + q - p\|^2 = \|\theta - p - (\eta - q) + (p - q)\|^2$

# Assumptions

$$(P) \quad \min_{\theta \in \mathbb{R}^d} F(\theta) \quad F(\theta) = f(\theta) + g(\theta),$$

## Assumptions

- $g : \mathbb{R}^d \rightarrow (-\infty, +\infty]$  closed convex
- $f : \Theta \rightarrow \mathbb{R}$  is convex continuously differentiable and  $\nabla f$  is gradient Lipshitz: for all  $\theta, \theta' \in \Theta$ ,

$$\|\nabla f(\theta) - \nabla f(\theta')\| \leq L\|\theta - \theta'\| ,$$

- 1 Stochastic approximation
  - Finite-sum optimization
  - Online learning
  - Smooth strongly convex case
  - Stochastic subgradient descent/method
  - Stochastic Approximation for nonconvex optimization
- 2 Proximal methods
  - Proximal operator
  - Proximal gradient algorithm
  - Stochastic proximal gradient
- 3 Applications
  - Network structure estimation
  - High-dimensional logistic regression with random effect

# Proximal gradient algorithm

$$\theta_{n+1} = \text{Prox}_{\gamma_{n+1}, g}(\theta_n - \gamma_{n+1} \nabla f(\theta_n))$$

where

$$\text{Prox}_{\gamma, g}(\tau) = \min_{\theta \in \Theta} \left( g(\theta) + \frac{1}{2\gamma} \|\theta - \tau\|^2 \right)$$

# Majorization-Minimization interpretation

- Since  $f$  is gradient Lipschitz, for all  $\gamma \in (0, 1/L]$

$$F(\eta) = f(\eta) + g(\eta) \leq f(\theta) + \langle \nabla f(\theta), \eta - \theta \rangle + \frac{1}{2\gamma} \|\theta - \eta\|^2 + g(\eta)$$

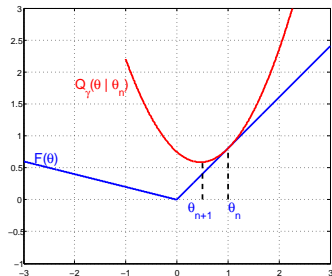
- Consider the following **surrogate function**

$$Q_\gamma(\eta|\theta) = f(\theta) + \langle \nabla f(\theta), \eta - \theta \rangle + \frac{1}{2\gamma} \|\theta - \eta\|^2 + g(\eta)$$

- For all  $\theta \in \Theta$ ,  $\eta \mapsto Q_\gamma(\eta|\theta)$  is strongly convex and has a **unique** minimum and

$$F(\eta) \leq Q_\gamma(\eta|\theta)$$

$$F(\theta) = Q_\gamma(\theta|\theta)$$



$$F(\eta) \leq Q_\gamma(\eta|\theta_n)$$

$$F(\theta_n) = Q_\gamma(\theta_n|\theta_n)$$

# Majorization-Minimization interpretation

$$\begin{aligned} Q_\gamma(\eta|\theta) &\stackrel{\text{def}}{=} f(\theta) + \langle \nabla f(\theta), \eta - \theta \rangle + \frac{1}{2\gamma} \|\eta - \theta\|^2 + g(\eta) \\ &= f(\theta) + \frac{1}{2\gamma} \|\eta - (\theta - \gamma \nabla f(\theta))\|^2 - \frac{\gamma}{2} \|\nabla f(\theta)\|^2 + g(\eta), \end{aligned}$$

The iterates of the proximal gradient algorithms may be rewritten as  $\theta_{n+1} = T_{\gamma_{n+1}}(\theta_n)$  with the point-to-point map  $T_\gamma$  defined by

$$\begin{aligned} T_\gamma(\theta) &\stackrel{\text{def}}{=} \text{Prox}_{\gamma, d}(\theta - \gamma \nabla f(\theta)) \\ &= \operatorname{argmin}_{\eta \in \text{Dom}(g)} Q_\gamma(\eta|\theta). \end{aligned}$$



# Proximal gradient

- If  $g(\theta) \equiv 0$ ,  $\Leftrightarrow$  gradient proximal = classical stochastic gradient

$$\theta_n = \theta_{n-1} - \gamma_n \nabla f(\theta_{n-1})$$

# Proximal gradient

- If  $g(\theta) \equiv 0$ ,  $\Leftrightarrow$  gradient proximal = classical stochastic gradient

$$\theta_n = \theta_{n-1} - \gamma_n \nabla f(\theta_{n-1})$$

- If  $g(\theta) \equiv 0$  if  $\theta \in \mathcal{C}$  and  $g(\theta) = +\infty$  otherwise where  $\mathcal{C}$  is a closed convex set,

$$\text{Prox}_{\gamma, g}(\tau) = \min_{\theta \in \mathcal{C}} \|\tau - \theta\|^2$$

$\Leftrightarrow$  gradient proximal = projected gradient

$$\theta_n = \Pi_{\mathcal{C}}(\theta_{n-1} - \gamma_n \nabla f(\theta_{n-1}))$$

# Proximal gradient for the elastic net penalty

$$\text{If } g(\theta) = \lambda \left( \frac{1-\alpha}{2} \|\theta\|_2^2 + \alpha \|\theta\|_1 \right)$$

$$(\text{Prox}_{\gamma, g}(\tau))_i = \frac{1}{1 + \gamma\lambda(1 - \alpha)} \begin{cases} \tau_i - \gamma\lambda\alpha & \text{if } \tau_i \geq \gamma\lambda\alpha \\ \tau_i + \gamma\lambda\alpha & \text{if } \tau_i \leq -\gamma\lambda\alpha \\ 0 & \text{otherwise} \end{cases}$$

↔ Proximal gradient = soft-thresholded gradient

$$\theta_{n+1} = \mathcal{S}_{\alpha, \lambda, \gamma_{n+1}} (\theta_n - \gamma_{n+1} \nabla f(\theta_n))$$

# Stationary points of the proximal gradient

$$\theta_{n+1} = \text{Prox}_{\gamma, g}(\theta_n - \gamma \nabla f(\theta_n)) = T_\gamma(\theta_n),$$

where  $T_\gamma$  is the proximal map,

$$T_\gamma(\theta) \stackrel{\text{def}}{=} \text{Prox}_{\gamma, g}(\theta - \gamma \nabla f(\theta)) = \underset{\eta \in \text{Dom}(g)}{\text{argmin}} Q_\gamma(\eta | \theta).$$

## Theorem

*The fixed points of the proximal map are the global minimizers of  $F(\theta) = f(\theta) + g(\theta)$ :*

$$\mathbf{L} = \{\theta : \theta = \text{Prox}_{\gamma, g}(\theta - \gamma \nabla f(\theta))\} = \{\theta \in \text{Dom}(g) : 0 \in \nabla f(\theta) + \partial g(\theta)\}.$$

## Fixed points of the proximal map

Since

$$F(\theta) = f(\theta) + g(\theta) ,$$

we get

$$\begin{aligned} 0 \in \partial F(\theta) &\iff 0 \in \partial \gamma F(\theta) \\ &\iff 0 \in \gamma \nabla f(\theta) + \partial \gamma g(\theta) \\ &\iff \theta - \gamma \nabla f(\theta) \in (\theta + \gamma \partial g(\theta)) \end{aligned}$$

## Fixed points of the proximal map

Since

$$F(\theta) = f(\theta) + g(\theta),$$

we get

$$\begin{aligned} 0 \in \partial F(\theta) &\iff 0 \in \partial \gamma F(\theta) \\ &\iff 0 \in \gamma \nabla f(\theta) + \partial \gamma g(\theta) \\ &\iff \theta - \gamma \nabla f(\theta) \in (\theta + \gamma \partial g(\theta)) \end{aligned}$$

Recall that, for any  $\eta$

$$p = \text{prox}_{\gamma g}(\eta) \iff (\eta - p) \in \gamma \partial g(p) \iff \eta \in p + \gamma \partial g(p).$$

## Fixed points of the proximal map

Since

$$F(\theta) = f(\theta) + g(\theta),$$

we get

$$\begin{aligned} 0 \in \partial F(\theta) &\iff 0 \in \partial \gamma F(\theta) \\ &\iff 0 \in \gamma \nabla f(\theta) + \partial \gamma g(\theta) \\ &\iff \theta - \gamma \nabla f(\theta) \in (\theta + \gamma \partial g(\theta)) \end{aligned}$$

Recall that, for any  $\eta$

$$p = \text{prox}_{\gamma g}(\eta) \iff (\eta - p) \in \gamma \partial g(p) \iff \eta \in p + \gamma \partial g(p).$$

Hence, taking  $p \leftarrow \theta$  and  $\eta \leftarrow \theta - \gamma \nabla f(\theta)$

$$0 \in \partial F(\theta) \iff \theta = T_\gamma(\theta)$$

# Lyapunov function

$$Q_\gamma(\eta|\theta) = f(\theta) + \langle \nabla f(\theta), \eta - \theta \rangle + \frac{1}{2\gamma} \|\theta - \eta\|^2 + g(\eta)$$

- For all  $\theta \in \Theta$ ,  $F \circ T_\gamma(\theta) \leq F(\theta)$ :

$$F \circ T_\gamma(\theta) \leq Q_\gamma(T_\gamma(\theta)|\theta) \leq Q_\gamma(\theta|\theta) = F(\theta)$$

Moreover, the inequality is strict unless  $\theta$  is a fixed point of the mapping  $T_\gamma$ .

- $F$  is a **Lyapunov function** for the proximal map  $T_\gamma$ .



# Convergence result

$$(P) \quad (\arg)\min_{\theta \in \Theta} \{f(\theta) + g(\theta)\},$$

- the objective function always converge  $\{F(\theta_n), n \geq 0\}$
- $f$  is convex: then  $\{\theta_n, n \in \mathbb{N}\}$  converges to  $\theta_*$ , where  $\theta_*$  is a minimizer of  $F$ .
- $F(\theta_n) - F(\theta_*) = O(1/n)$ .
- Results similar to smooth optimization ( $O(1/n)$  where  $n$  is the number of iterations)
- Acceleration methods: Nesterov, 2007; Beck and Teboulle, 2009. ( $O(1/n^2)$ ) [algorithm FISTA]

- 1 Stochastic approximation
  - Finite-sum optimization
  - Online learning
  - Smooth strongly convex case
  - Stochastic subgradient descent/method
  - Stochastic Approximation for nonconvex optimization
- 2 Proximal methods
  - Proximal operator
  - Proximal gradient algorithm
  - Stochastic proximal gradient
- 3 Applications
  - Network structure estimation
  - High-dimensional logistic regression with random effect

# Stochastic proximal gradient

## Objective

- Exact algorithm :

$$\theta_{n+1} = \text{Prox}_{\gamma_{n+1},g} (\theta_n - \gamma_{n+1} \nabla f(\theta_n))$$

- Perturbed algorithm :

$$\theta_{n+1} = \text{Prox}_{\gamma_{n+1},g} (\theta_n - \gamma_{n+1} H_{n+1})$$

where  $H_{n+1}$  is a noisy approximation of the true gradient  $\nabla f(\theta_n)$ .

- Problem** find sufficient conditions on the **stochastic error**

$$\eta_{n+1} = H_{n+1} - \nabla f(\theta_n)$$

to preserve convergence (closely related to SA).

# Convergence of the parameter

## Theorem

Assume  $f$  is  $L$ -smooth and the set  $\mathbf{L} = \operatorname{argmin}_{\theta \in \Theta} F(\theta)$  is non-empty. Assume in addition that  $\gamma_n \in (0, 1/L]$  for any  $n \geq 1$  and  $\sum_n \gamma_n = +\infty$ . If the following series converge

$$\sum_{n \geq 0} \gamma_{n+1} \langle T_{\gamma_{n+1}}(\theta_n), \eta_{n+1} \rangle, \quad \sum_{n \geq 0} \gamma_{n+1} \eta_{n+1}, \quad \sum_{n \geq 0} \gamma_{n+1}^2 \|\eta_{n+1}\|^2,$$

then there exists  $\theta_\infty \in \mathbf{L}$  such that  $\lim_n \theta_n = \theta_\infty$ .

# Convergence of the function

## Theorem

Assume  $f$  is  $L$ -smooth and the set  $\mathbf{L} = \operatorname{argmin}_{\theta \in \Theta} F(\theta)$  is non-empty. Assume that  $\gamma_n \in (0, 1/L]$  and let  $\{a_0, \dots, a_n\}$  be nonnegative weights. Then, for any  $\theta_\star \in \mathbf{L}$  and  $n \geq 1$ ,

$$\sum_{k=1}^n a_k \{F(\theta_k) - \min F\} \leq U_n(\theta_\star)$$

where

$$U_n(\theta_\star) \stackrel{\text{def}}{=} \frac{1}{2} \sum_{k=1}^n \left( \frac{a_k}{\gamma_k} - \frac{a_{k-1}}{\gamma_{k-1}} \right) \|\theta_{k-1} - \theta_\star\|^2 + \frac{a_0}{2\gamma_0} \|\theta_0 - \theta_\star\|^2 \\ - \sum_{k=1}^n a_k \langle T_{\gamma_k}(\theta_{k-1}) - \theta_\star, \eta_k \rangle + \sum_{k=1}^n a_k \gamma_k \|\eta_k\|^2 .$$

## Sanity check

- Assume that the gradient is exact, i.e.  $\eta_n = 0$ . Set  $A_n = \sum_{k=1}^n a_k$   
Then

$$\begin{aligned} F\left(A_n^{-1} \sum_{j=1}^n \theta_j\right) - \min F &\leq A_n^{-1} \sum_{j=1}^n a_j F(\theta_j) - \min F \\ &\leq \frac{1}{2} \sum_{k=1}^n \left(\frac{a_k}{\gamma_k} - \frac{a_{k-1}}{\gamma_{k-1}}\right) \|\theta_{k-1} - \theta_\star\|^2 + \frac{a_0}{2\gamma_0} \|\theta_0 - \theta_\star\|^2 \end{aligned}$$

- Setting  $a_k \equiv 1$  and  $\gamma_k \equiv 1/L$

$$\begin{aligned} F\left(n^{-1} \sum_{j=1}^n \theta_j\right) - \min F &\leq n^{-1} \sum_{j=1}^n F(\theta_j) - \min F \\ &\leq \frac{L}{2} \|\theta_0 - \theta_\star\|^2 \end{aligned}$$

- Up to constant, this is the same bound than the gradient algorithm for smooth convex function.

# Perturbed gradient

- Take  $a_k = \gamma_k$ , for  $k \in \{1, \dots, n\}$ . Then, for any  $\theta_\star \in \mathbf{L}$  and  $n \geq 1$ ,

$$F\left(\Gamma_n^{-1} \sum_{k=1}^n \gamma_k \theta_k\right) - \min F \leq \frac{1}{2\Gamma_n} \|\theta_0 - \theta_\star\|^2 \\ - \Gamma_n^{-1} \sum_{k=1}^n \gamma_k \langle T_{\gamma_k}(\theta_{k-1}) - \theta_\star, \eta_k \rangle + \Gamma_n^{-1} \sum_{k=1}^n \gamma_k^2 \|\eta_k\|^2.$$

- Problem:** Control the sequences  $\sum_{k=1}^n \gamma_k \langle T_{\gamma_k}(\theta_{k-1}) - \theta_\star, \eta_k \rangle$  and  $\sum_{k=1}^n \gamma_k^2 \|\eta_k\|^2$  in expectation or using high-probability bounds.

# Robbins-Monro setting

$$\nabla f(\theta) = \int_{\mathcal{X}} H_{\theta}(x) \pi(dx)$$

- Set

$$H_{n+1} = m_{n+1}^{-1} \sum_{j=1}^{m_{n+1}} H_{\theta_n}(X_{n+1}^{(j)})$$

where  $m_{n+1}$  is the size of the batch and  $\{X_{n+1}^{(j)}, 1 \leq j \leq m_{n+1}\}$  is a sample from  $\pi$  independent of  $\sigma(\theta_{\ell}, \ell \leq n)$ .

- In such case,

$$\mathbb{E}[H_{n+1} | \mathcal{F}_n] = m_{n+1}^{-1} \sum_{j=1}^{m_{n+1}} \mathbb{E}[H_{\theta_n}(X_{n+1}^{(j)}) | \mathcal{F}_n] = \nabla f(\theta_n) \text{ and}$$

$\eta_{n+1} = H_{n+1} - \nabla f(\theta_n)$  is a martingale increment.



## Bounded case / Constant stepsizes - Risk Bounds

- Assume that  $\|H_\theta(x)\| \leq B$ , then  $\|\eta_{n+1}\| \leq 2B$  and the stepsizes are constant  $\gamma_k \equiv 1/B\sqrt{n}$  for  $k \in \{1, \dots, n\}$ .
- On one hand

$$\Gamma_n^{-1} \sum_{k=1}^n \gamma_k^2 \|\eta_{k+1}\|^2 \leq \frac{4B}{\sqrt{n}}$$

- Risk bound:** since  $\mathbb{E}[\langle T_{\gamma_k}(\theta_{k-1}) - \theta_*, \eta_k \rangle \mid \mathcal{F}_{k-1}] = 0$  (since  $\mathbb{E}[\eta_k \mid \mathcal{F}_{k-1} = 0] = 0$ ), the **risk bound** is

$$\mathbb{E} \left[ F \left( n^{-1} \sum_{k=1}^n \theta_k \right) \right] - \min F \leq \frac{B}{2\sqrt{n}} \|\theta_0 - \theta_*\|^2 + \frac{4B}{\sqrt{n}}.$$

- Same risk bound than the Stochastic subgradient method (minimax rate)

# Bounded case / Constant stepsizes - Concentration

- Azuma-Hoeffding inequality for bounded martingale increments  $\{Z_k, k \in \mathbb{N}^*\}$ :

$$\mathbb{P} \left( \frac{1}{n} \sum_{k=1}^n Z_k \geq \frac{Ct}{\sqrt{n}} \right) \leq \exp(-t^2/2)$$

- Apply it to

$$Z_k = \langle T_{\gamma_k}(\theta_{k-1}) - \theta_*, \eta_k \rangle .$$

1 Stochastic approximation

2 Proximal methods

**3 Applications**

- 1 Stochastic approximation
  - Finite-sum optimization
  - Online learning
  - Smooth strongly convex case
  - Stochastic subgradient descent/method
  - Stochastic Approximation for nonconvex optimization
- 2 Proximal methods
  - Proximal operator
  - Proximal gradient algorithm
  - Stochastic proximal gradient
- 3 Applications
  - Network structure estimation
  - High-dimensional logistic regression with random effect

# Network structure estimation

- **Problem** fitting a discrete graphical models in a setting where the number of nodes in the graph is large compared to the sample size.
- **Formalization** Let  $A$  be a nonempty finite set, and  $p \geq 1$  an integer. Consider a graphical model on  $X = A^p$  with p.m.f.

$$f_{\theta}(x_1, \dots, x_p) = \frac{1}{Z_{\theta}} \exp \left\{ \sum_{k=1}^p \theta_{kk} B_0(x_k) + \sum_{1 \leq j < k \leq p} \theta_{kj} B(x_k, x_j) \right\},$$

for a non-zero function  $B_0 : A \rightarrow \mathbb{R}$  and a symmetric non-zero function  $B : A \times A \rightarrow \mathbb{R}$ .

- The term  $Z_{\theta}$  is the normalizing constant of the distribution (the partition function), which cannot (in general) be computed explicitly.

# Network structure estimation

- **Problem** fitting a discrete graphical models in a setting where the number of nodes in the graph is large compared to the sample size.
- **Formalization** Let  $A$  be a nonempty finite set, and  $p \geq 1$  an integer. Consider a graphical model on  $X = A^p$  with p.m.f.

$$f_{\theta}(x_1, \dots, x_p) = \frac{1}{Z_{\theta}} \exp \left\{ \sum_{k=1}^p \theta_{kk} B_0(x_k) + \sum_{1 \leq j < k \leq p} \theta_{kj} B(x_k, x_j) \right\},$$

for a non-zero function  $B_0 : A \rightarrow \mathbb{R}$  and a symmetric non-zero function  $B : A \times A \rightarrow \mathbb{R}$ .

- The real-valued symmetric matrix  $\theta$  defines the graph structure and is the parameter of interest. Same interpretation as the precision matrix in a multivariate Gaussian distribution.

# Network structure estimation

- **Problem:** Estimate  $\theta$  from  $N$  realizations  $\{x^{(i)}, 1 \leq i \leq N\}$  where  $x^{(i)} = (x_1^{(i)}, \dots, x_p^{(i)}) \in A^p$  under sparsity constraint.
- **Applications** biology, social sciences,
- **Main difficulty:** the log-partition function  $\log Z_\theta$  is intractable in general.
  - Most of the existing results use a pseudo-likelihood function.
  - One exception is [hoefling09], using an active set strategy (to preserve sparsity), and the junction tree algorithm for computing the partial derivatives of the log-partition function. However, this algorithm does not scale

# Model

- Penalized likelihood  $F(\theta) = -\ell(\theta) + g(\theta)$  where

$$\ell(\theta) = \frac{1}{N} \sum_{i=1}^N \langle \theta, \bar{B}(x^{(i)}) \rangle - \log Z_{\theta} \quad \text{and} \quad g(\theta) = \lambda \sum_{1 \leq k \leq j \leq p} |\theta_{jk}| ;$$

the matrix-valued function  $\bar{B} : \mathcal{X} \rightarrow \mathbb{R}^{p \times p}$  is defined by

$$\bar{B}_{kk}(x) = B_0(x_k) \quad \bar{B}_{kj}(x) = B(x_k, x_j), k \neq j .$$

- **Intractable** canonical exponential model.



# Model

- Penalized likelihood  $F(\theta) = -\ell(\theta) + g(\theta)$  where

$$\ell(\theta) = \frac{1}{N} \sum_{i=1}^N \langle \theta, \bar{B}(x^{(i)}) \rangle - \log Z_{\theta} \text{ and } g(\theta) = \lambda \sum_{1 \leq k \leq j \leq p} |\theta_{jk}| ;$$

the matrix-valued function  $\bar{B} : \mathcal{X} \rightarrow \mathbb{R}^{p \times p}$  is defined by

$$\bar{B}_{kk}(x) = B_0(x_k) \quad \bar{B}_{kj}(x) = B(x_k, x_j), k \neq j .$$

- $\theta \mapsto -\ell(\theta)$  is convex and

$$\nabla \ell(\theta) = \frac{1}{N} \sum_{i=1}^N \bar{B}(x^{(i)}) - \int_{\mathcal{X}} \bar{B}(z) f_{\theta}(z) \mu(dz) ,$$

# Implementation

- Direct simulation from the distribution  $f_{\theta}$  is not feasible.
- If  $X$  is not too large, then a Gibbs sampler that samples from the full conditional distributions of  $f_{\theta}$  can be easily implemented.
- Gibbs sampler is a generic algorithm that in some cases is known to mix poorly. Whenever possible we recommend the use of specialized problem-specific MCMC algorithms with better mixing properties...

# Set up

- $X = \{1, \dots, M\}$ ,  $B_0(x) = 0$ , and  $B(x, y) = \mathbf{1}_{\{x=y\}}$ , which corresponds to the Potts model.
- We use  $M = 20$ ,  $B_0(x) = x$ ,  $N = 250$  and for  $p \in \{50, 100, 200\}$ .
- We generate the 'true' matrix  $\theta_{\text{true}}$  such that it has on average  $p$  non-zero elements off-diagonal which are simulated from a uniform distribution on  $(-4, -1) \cup (1, 4)$ .
- All the diagonal elements are set to 0.

# Algorithms

- Two versions of the stochastic proximal gradient are considered
  - 1 Solver 1: A version with a fixed Monte Carlo batch size  $m_n = 500$ , and decreasing step size  $\gamma_n = \frac{25}{p} \frac{1}{n^{0.7}}$ .
  - 2 Solver 2: A version with increasing Monte Carlo batch size  $m_n = 500 + n^{1.2}$ , and fixed step size  $\gamma_n = \frac{25}{p} \frac{1}{\sqrt{50}}$ .
- The set-up is such that both solvers draw approximately the same number of Monte Carlo samples.

# Algorithms

- Two versions of the stochastic proximal gradient are considered
  - 1 Solver 1: A version with a fixed Monte Carlo batch size  $m_n = 500$ , and decreasing step size  $\gamma_n = \frac{25}{p} \frac{1}{n^{0.7}}$ .
  - 2 Solver 2: A version with increasing Monte Carlo batch size  $m_n = 500 + n^{1.2}$ , and fixed step size  $\gamma_n = \frac{25}{p} \frac{1}{\sqrt{50}}$ .
- We evaluate the convergence of each solver by computing the relative error  $\|\theta_n - \theta_\infty\| / \|\theta_\infty\|$ , along the iterations, where  $\theta_\infty$  denotes the value returned by the solver on its last iteration.

# Algorithms

- Two versions of the stochastic proximal gradient are considered
  - 1 Solver 1: A version with a fixed Monte Carlo batch size  $m_n = 500$ , and decreasing step size  $\gamma_n = \frac{25}{p} \frac{1}{n^{0.7}}$ .
  - 2 Solver 2: A version with increasing Monte Carlo batch size  $m_n = 500 + n^{1.2}$ , and fixed step size  $\gamma_n = \frac{25}{p} \frac{1}{\sqrt{50}}$ .
- We compare the optimizer output to  $\theta_\infty$ , not  $\theta_{\text{true}}$ . Ideally, we would like to compare the iterates to the solution of the optimization problem. However in the present setting a solution is not available in closed form (and there could be more than one solution).

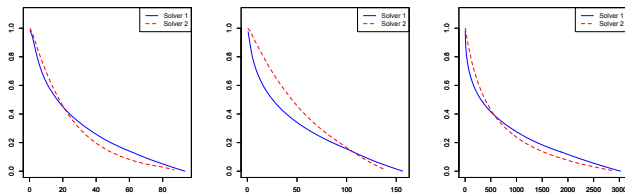


Figure: Relative errors plotted as function of computing time for Solver 1 and Solver 2.

When measured as function of resource used, Solver 1 and Solver 2 have roughly the same convergence rate.

# Sensitivity and Precision

- We also compute the statistic  $F_n \stackrel{\text{def}}{=} \frac{2\text{Sen}_n \text{Prec}_n}{\text{Sen}_n + \text{Prec}_n}$  which measures the recovery of the sparsity structure of  $\theta_\infty$  along the iteration.
- In this definition  $\text{Sen}_n$  is the sensitivity, and  $\text{Prec}_n$  is the precision defined as

$$\text{Sen}_n = \frac{\sum_{j < i} \mathbf{1}_{\{|\theta_{n,ij}| > 0\}} \mathbf{1}_{\{|\theta_{\infty,ij}| > 0\}}}{\sum_{j < i} \mathbf{1}_{\{|\theta_{\infty,ij}| > 0\}}}$$
$$\text{Prec}_n = \frac{\sum_{j < i} \mathbf{1}_{\{|\theta_{n,ij}| > 0\}} \mathbf{1}_{\{|\theta_{\infty,ij}| > 0\}}}{\sum_{j < i} \mathbf{1}_{\{|\theta_{n,ij}| > 0\}}}.$$



# Sensitivity and Precision

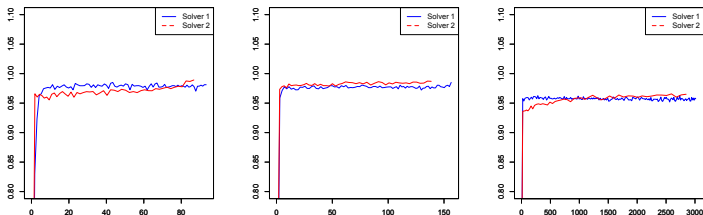


Figure: Statistic  $F_n$  plotted as function of computing time for Solver 1 and Solver 2.

- 1 Stochastic approximation
  - Finite-sum optimization
  - Online learning
  - Smooth strongly convex case
  - Stochastic subgradient descent/method
  - Stochastic Approximation for nonconvex optimization
- 2 Proximal methods
  - Proximal operator
  - Proximal gradient algorithm
  - Stochastic proximal gradient
- 3 Applications
  - Network structure estimation
  - High-dimensional logistic regression with random effect

# High-dimensional logistic regression with random effects

- Observations :  $N$  observations  $\mathbf{Y} \in \{0, 1\}^N$
- Random effect : Conditionally to  $\mathbf{U}$ , for all  $i = 1, \dots, N$ ,

$$Y_i \stackrel{\text{ind.}}{\sim} \mathcal{B} \left( \frac{\exp(\eta_i)}{1 + \exp(\eta_i)} \right)$$

where

$$\begin{bmatrix} \eta_1 \\ \vdots \\ \eta_N \end{bmatrix} = \mathbf{X}\beta_* + \sigma_*\mathbf{Z}\mathbf{U}$$

- The regressors  $\mathbf{X} \in \mathbb{R}^{N \times p}$  and the factor loadings  $\mathbf{Z} \in \mathbb{R}^{N \times q}$ , known.
- Objective: estimate  $\beta_* \in \mathbb{R}^p, \sigma_* > 0$ .

# Penalized likelihood

- log-likelihood : Taking  $\mathbf{U} \sim \mathcal{N}_q(0, I)$ , setting

$$\theta = (\beta, \sigma) \quad F(\eta) = \frac{\exp(\eta)}{1 + \exp(\eta)}$$

the log-likelihood of the observations  $\mathbf{Y}$  (with respect to  $\theta$ ) is

$$\ell(\theta) = \log \int \prod_{i=1}^N \{F(\mathbf{X}_i \cdot \beta + \sigma(\mathbf{Z}\mathbf{U})_i)\}^{Y_i} \{1 - F(\mathbf{X}_i \cdot \beta + \sigma(\mathbf{Z}\mathbf{U})_i)\}^{1-Y_i} \phi(\mathbf{u}) d\mathbf{u}$$

- Elastic net penalty

$$g_{\lambda, \theta}(\theta) = \lambda \left( \frac{1 - \alpha}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 \right)$$
$$\tilde{g}_{\mathcal{C}}(\theta) = \begin{cases} 0 & \text{si } \theta \in \mathcal{C} \\ +\infty & \text{otherwise} \end{cases}$$

# Penalized likelihood

$$\min_{\theta \in \Theta} (f(\theta) + g(\theta)) \quad , \quad f(\theta) = -\ell(\theta) \quad ,$$

with

$$\ell(\theta) = \log \int \exp(\ell_c(\theta|\mathbf{u})) \phi(\mathbf{u}) d\mathbf{u}$$

$$\ell_c(\theta|\mathbf{u}) = \sum_{i=1}^N \{Y_i (\mathbf{X}_i \cdot \beta + \sigma(\mathbf{Z}\mathbf{U})_i) - \ln(1 + \exp(\mathbf{X}_i \cdot \beta + \sigma(\mathbf{Z}\mathbf{U})_i))\}$$

Gradient :

$$\nabla \ell(\theta) = \int \nabla \ell_c(\theta|\mathbf{u}) \pi_\theta(\mathbf{u}) d\mathbf{u}$$

where  $\pi_\theta(\mathbf{u})$  is the **posterior distribution** of the random effect given the observations

$$\pi_\theta(\mathbf{u}) = \exp(\ell_c(\theta|\mathbf{u}) - \ell(\theta)) \phi(\mathbf{u})$$

# Penalized likelihood

$$\min_{\theta \in \Theta} (f(\theta) + g(\theta)) , \quad f(\theta) = -\ell(\theta)$$

where

$$g_{\lambda, \theta}(\theta) = \lambda \left( \frac{1 - \alpha}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 \right) + \mathbb{I}_{\mathcal{C}}(\theta)$$

$$\mathbb{I}_{\mathcal{C}}(\theta) = \begin{cases} 0 & \text{if } \theta \in \mathcal{C} \\ +\infty & \text{otherwise} \end{cases} \quad \mathcal{C} \text{ compact convex set}$$

$\hookrightarrow$  proper convex,  
lower-semi continuous, not differentiable.

# MCMC algorithm

- The distribution  $\pi_\theta$  is sampled using the MCMC sampler proposed in (Polson et al, 2012) based on data-augmentation.
- We write  $-\nabla\ell(\theta) = \int_{\mathbb{R}^q \times \mathbb{R}^N} H_\theta(\mathbf{u}) \tilde{\pi}_\theta(\mathbf{u}, \mathbf{w}) \, d\mathbf{u} d\mathbf{w}$  where  $\tilde{\pi}_\theta(\mathbf{u}, \mathbf{w})$  is defined for  $\mathbf{u} \in \mathbb{R}^q$  and  $\mathbf{w} = (w_1, \dots, w_N) \in \mathbb{R}^N$  by

$$\tilde{\pi}_\theta(\mathbf{u}, \mathbf{w}) = \left( \prod_{i=1}^N \bar{\pi}_{\text{PG}}(w_i; x'_i \beta + \sigma z'_i \mathbf{u}) \right) \pi_\theta(\mathbf{u}) ;$$

- in this expression,  $\bar{\pi}_{\text{PG}}(\cdot; c)$  is the density of the Polya-Gamma distribution on the positive real line with parameter  $c$  given by

$$\bar{\pi}_{\text{PG}}(w; c) = \cosh(c/2) \exp(-wc^2/2) \rho(w) \mathbb{1}_{\mathbb{R}^+}(w) ,$$

where  $\rho(w) \propto \sum_{k \geq 0} (-1)^k (2k+1) \exp(-(2k+1)^2/(8w)) w^{-3/2}$

# MCMC algorithm

- The distribution  $\pi_\theta$  is sampled using the MCMC sampler proposed in (Polson et al, 2012) based on data-augmentation.
- We write  $-\nabla\ell(\theta) = \int_{\mathbb{R}^q \times \mathbb{R}^N} H_\theta(\mathbf{u}) \tilde{\pi}_\theta(\mathbf{u}, \mathbf{w}) \, d\mathbf{u} d\mathbf{w}$  where  $\tilde{\pi}_\theta(\mathbf{u}, \mathbf{w})$  is defined for  $\mathbf{u} \in \mathbb{R}^q$  and  $\mathbf{w} = (w_1, \dots, w_N) \in \mathbb{R}^N$  by

$$\tilde{\pi}_\theta(\mathbf{u}, \mathbf{w}) = \left( \prod_{i=1}^N \tilde{\pi}_{\text{PG}}(w_i; x_i' \beta + \sigma z_i' \mathbf{u}) \right) \pi_\theta(\mathbf{u}) ;$$

- Thus, we have

$$\tilde{\pi}_\theta(\mathbf{u}, \mathbf{w}) = C_\theta \phi(\mathbf{u}) \prod_{i=1}^N \exp(\sigma(Y_i - 1/2)z_i' \mathbf{u} - w_i(x_i' \beta + \sigma z_i' \mathbf{u})^2 / 2) \rho(w_i) \mathbb{1}_{\cdot}$$

where  $\ln C_\theta = -N \ln 2 - \ell(\theta) + \sum_{i=1}^N (Y_i - 1/2)x_i' \beta$ .



# MCMC algorithm

- The distribution  $\pi_\theta$  is sampled using the MCMC sampler proposed in (Polson et al, 2012) based on data-augmentation.
- We write  $-\nabla\ell(\theta) = \int_{\mathbb{R}^q \times \mathbb{R}^N} H_\theta(\mathbf{u}) \tilde{\pi}_\theta(\mathbf{u}, \mathbf{w}) \, d\mathbf{u} d\mathbf{w}$  where  $\tilde{\pi}_\theta(\mathbf{u}, \mathbf{w})$  is defined for  $\mathbf{u} \in \mathbb{R}^q$  and  $\mathbf{w} = (w_1, \dots, w_N) \in \mathbb{R}^N$  by

$$\tilde{\pi}_\theta(\mathbf{u}, \mathbf{w}) = \left( \prod_{i=1}^N \tilde{\pi}_{\text{PG}}(w_i; x'_i \beta + \sigma z'_i \mathbf{u}) \right) \pi_\theta(\mathbf{u}) ;$$

- This target distribution can be sampled using a Gibbs algorithm

# Numerics

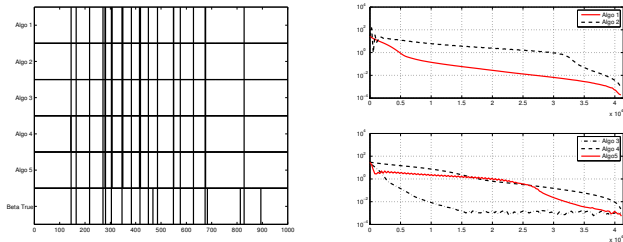
- We test the algorithms with  $N = 500$ ,  $p = 1,000$  and  $q = 5$ .
- We generate the  $N \times p$  covariates matrix  $X$  columnwise, by sampling a stationary  $\mathbb{R}^N$ -valued autoregressive model with parameter  $\rho = 0.8$  and Gaussian noise  $\sqrt{1 - \rho^2} \mathcal{N}_N(0, I)$ .
- We generate the vector of regressors  $\beta_{\text{true}}$  from the uniform distribution on  $[1, 5]$  and randomly set 98% of the coefficients to zero.
- The variance of the random effect is set to  $\sigma^2 = 0.1$ .

# Numerics

We first illustrate the ability of Monte Carlo Proximal Gradient algorithms to find a minimizer of  $F$ . We compare the Monte Carlo proximal gradient algorithm

- 1 with fixed batch size:  $\gamma_n = 0.01/\sqrt{n}$  and  $m_n = 275$  (Algo 1);  
 $\gamma_n = 0.5/n$  and  $m_n = 275$  (Algo 2).
- 2 with increasing batch size:  $\gamma_n = \gamma = 0.005$ ,  $m_n = 200 + n$  (Algo 3);  
 $\gamma_n = \gamma = 0.001$ ,  $m_n = 200 + n$  (Algo 4); and  $\gamma_n = 0.05/\sqrt{n}$  and  
 $m_n = 270 + \lceil \sqrt{n} \rceil$  (Algo 5).

# Results



**Figure:** [left] The support of the sparse vector  $\beta_\infty$  obtained by Algo 1 to Algo 5; for comparison, the support of  $\beta_{\text{true}}$  is on the bottom row. [right] Relative error along one path of each algorithm as a function of the total number of Monte Carlo samples.

# Results

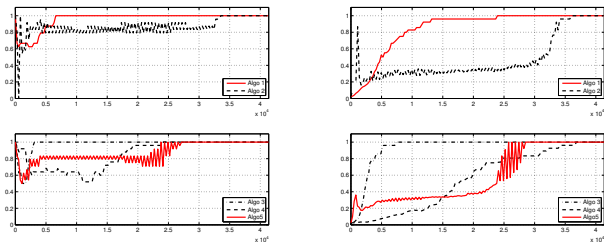


Figure: The sensitivity  $\text{Sen}_n$  [left] and the precision  $\text{Prec}_n$  [right] along a path, versus the total number of Monte Carlo samples up to time  $n$

# Bibliography I