# Large-scale machine learning and convex optimization

# **Francis Bach**

INRIA - Ecole Normale Supérieure, Paris, France



Hausdorff institute, Bonn - January 2016 Slides available at: www.di.ens.fr/~fbach/gradsto\_bonn\_2016.pdf

# "Big data" revolution? A new scientific context

- Data everywhere: size does not (always) matter
- Science and industry
- Size and variety
- Learning from examples
  - n observations in dimension d

# **Search engines - advertising**

e o o Sfete de la science - Recherch ×								
🗲 $\Rightarrow$ C 🔒 https://www.google.fr/search?hl=fr&safe=active&q=fete+de+la+science&oq=fete+de+la+sci&gs_l=serp.3.0.0i 🏠 🚍								
🔟 Le Monde 🛛 🔏 Intranet IM	NRIA 🖉 Francis Bach 📉 GMAIL 📥 Liberation 🗾 L'EC	QUIPE 🛛 👌 Google Scholar	PAMI 🚦 iGoogle	СР »				
+Francis Recherc	he Images Maps Play YouTube A	Actualités Gmail	Drive Agenda	Plus -				
Google	fete de la science							
0								
Recherche Environ 561 000 000 résultats (0,20 secondes)								
Web	Accueil - Fête de la science (site inte	rnet)						
Image	www.fetedelascience.fr/							
Images	Fête de la science 2012, du 10 au 14 octobre. La science vient à votre rencontre !							
Maps	Manipulez, jouez, experimentez, visitez des labo							
Vidéos	<u>Les programmes régionaux</u> <u>imprimable. Quel que soit votre</u> <u>Fête de la science 2012</u> <u>Villages des sciences, opérations</u>							
Actualités	choix, toutes les animations	d'envergure, man	ifestations					
Changing	Déposer un projet ? Le mode	20e édition en	2011					
Snopping	Déposer un projet ? Le mode d'emploi.	20e édition en 20	11. La Fête de la					
Plus	Bienvenue aux futurs	science se déroul	e du 12 au 16					
	Tout savoir sur la Fête de la	Les lauréats n	ationaux					

# **Search engines - Advertising**

000	tour de france - Bing ×					
← → C	https://www.bing.com/search?q=tour	+de+france&go=Submit&qs=n&	form=QBRE&	filt=all&p	q=tour+de+	-france≻=
🔛 Apps 🛛 G	MAIL 🔀 Intranet 🧁 Francis Bach - INRIA	JII Le Monde 📄 CP 🔣 Scholar	💈 Equipe 🛛 🔢	Agenda	Liberation	PAMI 🧠
	WEB IMAGES VIDEOS MAPS	NEWS MORE				Sign in
bing 🛛	tour de france		Q			
	121 000 000 RESULTS Narrow by langua	age • Narrow by region •				
	Tour de France 2014 Translate this	Translate this page		Relat	ed searche	es
	www.letour.fr - tour de picardie 2014 ag2r la mondiale; a racing team; bretagne - seche environnement	nc	Tracé Tour de France 2014 Regarder Tour de France Direct		2014 nce Direct	
	Parcours	Tour de France 2011		Itinéraire Tour de France		nce
	Du samedi 29 juin au dimanche 21 juillet 2013, le 100 e Tour de …	Tour de France 2014 - Site officiel de célèbre course cycliste Le Tour	e la	Etape Du Tour France 2 Tour de France Cyclisme		
	<u>Classements</u>	Étape 14				
	Classements - Tour de France 2013. Tour de France 2013 - Site officiel	Étape 14 - Saint-Pourçain-sur-Sioule Lyon - Tour de	>	Tour de France Online		B
	Nice 2013	Étape 18				
	Tour de France 2012 - Site officiel de la célèbre course cycliste Le Tour	Étape 18 - Gap > Alpe-d'Huez - Tour France 2013	de			
	Tour de France 2013 Translate this www.letour.fr/le-tour/2013/fr ▼ Tour de France 2013 - Site officiel de la célé Contient les itinéraires, coureurs, équipes et Tour de France (cyclisme) — W fr.wikipedia.org/wiki/Tour_de_France_(cyclis Le Tour de France est une compétition cycli Desgrange et Géo Lefèvre, chef de la rubriqu Histoire · Médiatisation du · Équipes et par	s page èbre course cycliste Le <b>Tour de Franc</b> les infos <b>des Tours</b> pas <mark>sés.</mark> <b>'<u>ikipédia</u> Translate this page sme) ▼ iste par étapes créée en 1903 par Henr le cyclisme du journal L'Auto. rticipation</b>	e. i			

# Marketing - Personalized recommendation



## **Visual object recognition**



#### **Personal photos**



### **Bioinformatics**



- Protein: Crucial elements of cell life
- Massive data: 2 millions for humans
- Complex data

# 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

# 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 and Monro, 1951)
  - Mixing statistics and optimization

# Outline

#### 1. Large-scale machine learning and optimization

- Traditional statistical analysis
- Classical methods for convex optimization

#### 2. Non-smooth stochastic approximation

- Stochastic (sub)gradient and averaging
- Non-asymptotic results and lower bounds
- Strongly convex vs. non-strongly convex

#### 3. Smooth stochastic approximation algorithms

- Asymptotic and non-asymptotic results
- 4. Beyond decaying step-sizes
- 5. Finite data sets

- Data: n observations  $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ ,  $i = 1, \ldots, n$ , i.i.d.
- Prediction as a linear function  $\theta^{\top} \Phi(x)$  of features  $\Phi(x) \in \mathbb{R}^d$
- (regularized) empirical risk minimization: find  $\hat{\theta}$  solution of

$$\min_{\theta \in \mathbb{R}^d} \quad \frac{1}{n} \sum_{i=1}^n \ell(y_i, \theta^\top \Phi(x_i)) + \mu \Omega(\theta)$$
  
convex data fitting term + regularizer

#### **Usual losses**

- **Regression**:  $y \in \mathbb{R}$ , prediction  $\hat{y} = \theta^{\top} \Phi(x)$ 
  - quadratic loss  $\frac{1}{2}(y-\hat{y})^2 = \frac{1}{2}(y-\theta^{\top}\Phi(x))^2$

#### **Usual losses**

• **Regression**:  $y \in \mathbb{R}$ , prediction  $\hat{y} = \theta^{\top} \Phi(x)$ 

– quadratic loss  $\frac{1}{2}(y-\hat{y})^2 = \frac{1}{2}(y-\theta^{\top}\Phi(x))^2$ 

- Classification :  $y \in \{-1, 1\}$ , prediction  $\hat{y} = \operatorname{sign}(\theta^{\top} \Phi(x))$ 
  - loss of the form  $\ell(y\,\theta^{\top}\Phi(x))$
  - "True" 0-1 loss:  $\ell(y \theta^{\top} \Phi(x)) = 1_{y \theta^{\top} \Phi(x) < 0}$
  - Usual convex losses:



## Main motivating examples

• Support vector machine (hinge loss): non-smooth

$$\ell(Y,\theta^{\top}\Phi(X)) = \max\{1 - Y\theta^{\top}\Phi(X), 0\}$$

• Logistic regression: smooth

$$\ell(Y, \theta^{\top} \Phi(X)) = \log(1 + \exp(-Y\theta^{\top} \Phi(X)))$$

• Least-squares regression

$$\ell(Y,\theta^{\top}\Phi(X)) = \frac{1}{2}(Y-\theta^{\top}\Phi(X))^2$$

# **Usual regularizers**

- Main goal: avoid overfitting
- (squared) Euclidean norm:  $\|\theta\|_2^2 = \sum_{j=1}^d |\theta_j|^2$ 
  - Numerically well-behaved
  - Representer theorem and kernel methods :  $\theta = \sum_{i=1}^{n} \alpha_i \Phi(x_i)$
  - See, e.g., Schölkopf and Smola (2001); Shawe-Taylor and Cristianini (2004)

# **Usual regularizers**

- Main goal: avoid overfitting
- (squared) Euclidean norm:  $\|\theta\|_2^2 = \sum_{j=1}^d |\theta_j|^2$ 
  - Numerically well-behaved
  - Representer theorem and kernel methods :  $\theta = \sum_{i=1}^{n} \alpha_i \Phi(x_i)$
  - See, e.g., Schölkopf and Smola (2001); Shawe-Taylor and Cristianini (2004)

#### • Sparsity-inducing norms

- Main example:  $\ell_1$ -norm  $\|\theta\|_1 = \sum_{j=1}^d |\theta_j|$
- Perform model selection as well as regularization
- Non-smooth optimization and structured sparsity
- See, e.g., Bach, Jenatton, Mairal, and Obozinski (2012b,a)

- Data: n observations  $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ ,  $i = 1, \ldots, n$ , i.i.d.
- Prediction as a linear function  $\theta^{\top} \Phi(x)$  of features  $\Phi(x) \in \mathbb{R}^d$
- (regularized) empirical risk minimization: find  $\hat{\theta}$  solution of

$$\min_{\theta \in \mathbb{R}^d} \quad \frac{1}{n} \sum_{i=1}^n \ell(y_i, \theta^\top \Phi(x_i)) + \mu \Omega(\theta)$$
  
convex data fitting term + regularizer

- Data: n observations  $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ ,  $i = 1, \ldots, n$ , i.i.d.
- Prediction as a linear function  $\theta^{\top} \Phi(x)$  of features  $\Phi(x) \in \mathbb{R}^d$
- (regularized) empirical risk minimization: find  $\hat{\theta}$  solution of

$$\min_{\theta \in \mathbb{R}^d} \quad \frac{1}{n} \sum_{i=1}^n \ell(y_i, \theta^\top \Phi(x_i)) \quad + \quad \mu \Omega(\theta)$$

convex data fitting term + regularizer

- Empirical risk:  $\hat{f}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \theta^{\top} \Phi(x_i))$  training cost
- Expected risk:  $f(\theta) = \mathbb{E}_{(x,y)} \ell(y, \theta^{\top} \Phi(x))$  testing cost
- Two fundamental questions: (1) computing  $\hat{\theta}$  and (2) analyzing  $\hat{\theta}$

- Data: n observations  $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ ,  $i = 1, \ldots, n$ , i.i.d.
- Prediction as a linear function  $\theta^{\top} \Phi(x)$  of features  $\Phi(x) \in \mathbb{R}^d$
- (regularized) empirical risk minimization: find  $\hat{\theta}$  solution of

$$\min_{\theta \in \mathbb{R}^d} \quad \frac{1}{n} \sum_{i=1}^n \ell(y_i, \theta^\top \Phi(x_i)) \quad + \quad \mu \Omega(\theta)$$

convex data fitting term + regularizer

- Empirical risk:  $\hat{f}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \theta^{\top} \Phi(x_i))$  training cost
- Expected risk:  $f(\theta) = \mathbb{E}_{(x,y)} \ell(y, \theta^{\top} \Phi(x))$  testing cost
- Two fundamental questions: (1) computing  $\hat{\theta}$  and (2) analyzing  $\hat{\theta}$ 
  - May be tackled simultaneously

- Data: n observations  $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ ,  $i = 1, \ldots, n$ , i.i.d.
- Prediction as a linear function  $\theta^{\top} \Phi(x)$  of features  $\Phi(x) \in \mathbb{R}^d$
- (regularized) empirical risk minimization: find  $\hat{\theta}$  solution of

$$\min_{\theta \in \mathbb{R}^d} \quad \frac{1}{n} \sum_{i=1}^n \ell(y_i, \theta^\top \Phi(x_i)) \text{ such that } \Omega(\theta) \leqslant D$$

convex data fitting term + constraint

- Empirical risk:  $\hat{f}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \theta^{\top} \Phi(x_i))$  training cost
- Expected risk:  $f(\theta) = \mathbb{E}_{(x,y)} \ell(y, \theta^{\top} \Phi(x))$  testing cost
- Two fundamental questions: (1) computing  $\hat{\theta}$  and (2) analyzing  $\hat{\theta}$ 
  - May be tackled simultaneously

#### **General assumptions**

- Data: n observations  $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ ,  $i = 1, \ldots, n$ , i.i.d.
- Bounded features  $\Phi(x) \in \mathbb{R}^d$ :  $\|\Phi(x)\|_2 \leq R$
- Empirical risk:  $\hat{f}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \theta^{\top} \Phi(x_i))$  training cost
- Expected risk:  $f(\theta) = \mathbb{E}_{(x,y)} \ell(y, \theta^{\top} \Phi(x))$  testing cost
- Loss for a single observation:  $f_i(\theta) = \ell(y_i, \theta^\top \Phi(x_i))$  $\Rightarrow \forall i, f(\theta) = \mathbb{E}f_i(\theta)$
- Properties of  $f_i, f, \hat{f}$ 
  - Convex on  $\mathbb{R}^d$
  - Additional regularity assumptions: Lipschitz-continuity, smoothness and strong convexity

## Lipschitz continuity

 Bounded gradients of f (⇔ Lipschitz-continuity): the function f if convex, differentiable and has (sub)gradients uniformly bounded by B on the ball of center 0 and radius D:

$$\forall \theta \in \mathbb{R}^d, \|\theta\|_2 \leqslant D \Rightarrow \|f'(\theta)\|_2 \leqslant B$$

#### • Machine learning

- with  $f(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \theta^{\top} \Phi(x_i))$
- G-Lipschitz loss and R-bounded data: B = GR

• A function  $f : \mathbb{R}^d \to \mathbb{R}$  is *L*-smooth if and only if it is differentiable and its gradient is *L*-Lipschitz-continuous

$$\forall \theta_1, \theta_2 \in \mathbb{R}^d, \| f'(\theta_1) - f'(\theta_2) \|_2 \leq L \| \theta_1 - \theta_2 \|_2$$

• If f is twice differentiable:  $\forall \theta \in \mathbb{R}^d, f''(\theta) \preccurlyeq L \cdot Id$ 



• A function  $f : \mathbb{R}^d \to \mathbb{R}$  is *L*-smooth if and only if it is differentiable and its gradient is *L*-Lipschitz-continuous

$$\forall \theta_1, \theta_2 \in \mathbb{R}^d, \ \|f'(\theta_1) - f'(\theta_2)\|_2 \leq L \|\theta_1 - \theta_2\|_2$$

• If f is twice differentiable:  $\forall \theta \in \mathbb{R}^d, f''(\theta) \preccurlyeq L \cdot Id$ 

#### • Machine learning

- with  $f(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \theta^{\top} \Phi(x_i))$
- Hessian  $\approx$  covariance matrix  $\frac{1}{n} \sum_{i=1}^{n} \Phi(x_i) \Phi(x_i)^{\top}$
- $\ell$ -smooth loss and R-bounded data:  $L = \ell R^2$

• A function  $f : \mathbb{R}^d \to \mathbb{R}$  is  $\mu$ -strongly convex if and only if

 $\forall \theta_1, \theta_2 \in \mathbb{R}^d, \ f(\theta_1) \ge f(\theta_2) + f'(\theta_2)^\top (\theta_1 - \theta_2) + \frac{\mu}{2} \|\theta_1 - \theta_2\|_2^2$ 

• If f is twice differentiable:  $\forall \theta \in \mathbb{R}^d, f''(\theta) \succcurlyeq \mu \cdot \mathrm{Id}$ 



• A function  $f : \mathbb{R}^d \to \mathbb{R}$  is  $\mu$ -strongly convex if and only if

 $\forall \theta_1, \theta_2 \in \mathbb{R}^d, \ f(\theta_1) \ge f(\theta_2) + f'(\theta_2)^\top (\theta_1 - \theta_2) + \frac{\mu}{2} \|\theta_1 - \theta_2\|_2^2$ 

• If f is twice differentiable:  $\forall \theta \in \mathbb{R}^d, f''(\theta) \succcurlyeq \mu \cdot \mathrm{Id}$ 



(large  $\mu$ )

(small  $\mu$ )

• A function  $f : \mathbb{R}^d \to \mathbb{R}$  is  $\mu$ -strongly convex if and only if

 $\forall \theta_1, \theta_2 \in \mathbb{R}^d, \ f(\theta_1) \ge f(\theta_2) + f'(\theta_2)^\top (\theta_1 - \theta_2) + \frac{\mu}{2} \|\theta_1 - \theta_2\|_2^2$ 

- If f is twice differentiable:  $\forall \theta \in \mathbb{R}^d, f''(\theta) \succcurlyeq \mu \cdot \mathrm{Id}$
- Machine learning
  - with  $f(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \theta^{\top} \Phi(x_i))$
  - Hessian  $\approx$  covariance matrix  $\frac{1}{n} \sum_{i=1}^{n} \Phi(x_i) \Phi(x_i)^{\top}$
  - Data with invertible covariance matrix (low correlation/dimension)

• A function  $f : \mathbb{R}^d \to \mathbb{R}$  is  $\mu$ -strongly convex if and only if

 $\forall \theta_1, \theta_2 \in \mathbb{R}^d, \ f(\theta_1) \ge f(\theta_2) + f'(\theta_2)^\top (\theta_1 - \theta_2) + \frac{\mu}{2} \|\theta_1 - \theta_2\|_3^2$ 

- If f is twice differentiable:  $\forall \theta \in \mathbb{R}^d, f''(\theta) \succcurlyeq \mu \cdot \mathrm{Id}$
- Machine learning
  - with  $f(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \theta^\top \Phi(x_i))$
  - Hessian  $\approx$  covariance matrix  $\frac{1}{n} \sum_{i=1}^{n} \Phi(x_i) \Phi(x_i)^{\top}$
  - Data with invertible covariance matrix (low correlation/dimension)
- Adding regularization by  $\frac{\mu}{2} \|\theta\|^2$ 
  - creates additional bias unless  $\mu$  is small

# Summary of smoothness/convexity assumptions

• Bounded gradients of f (Lipschitz-continuity): the function f if convex, differentiable and has (sub)gradients uniformly bounded by B on the ball of center 0 and radius D:

 $\forall \theta \in \mathbb{R}^d, \|\theta\|_2 \leqslant D \Rightarrow \|f'(\theta)\|_2 \leqslant B$ 

• Smoothness of f: the function f is convex, differentiable with L-Lipschitz-continuous gradient f':

$$\forall \theta_1, \theta_2 \in \mathbb{R}^d, \quad \|f'(\theta_1) - f'(\theta_2)\|_2 \leq L \|\theta_1 - \theta_2\|_2$$

• Strong convexity of f: The function f is strongly convex with respect to the norm  $\|\cdot\|$ , with convexity constant  $\mu > 0$ :

$$\forall \theta_1, \theta_2 \in \mathbb{R}^d, \ f(\theta_1) \ge f(\theta_2) + f'(\theta_2)^\top (\theta_1 - \theta_2) + \frac{\mu}{2} \|\theta_1 - \theta_2\|_2^2$$

### Analysis of empirical risk minimization

• Approximation and estimation errors:  $\Theta = \{\theta \in \mathbb{R}^d, \Omega(\theta) \leq D\}$ 

$$f(\hat{\theta}) - \min_{\theta \in \mathbb{R}^d} f(\theta) = \begin{bmatrix} f(\hat{\theta}) - \min_{\theta \in \Theta} f(\theta) \end{bmatrix} + \begin{bmatrix} \min_{\theta \in \Theta} f(\theta) - \min_{\theta \in \mathbb{R}^d} f(\theta) \end{bmatrix}$$
  
Estimation error Approximation error

– NB: may replace  $\min_{\theta \in \mathbb{R}^d} f(\theta)$  by best (non-linear) predictions

#### Analysis of empirical risk minimization

• Approximation and estimation errors:  $\Theta = \{\theta \in \mathbb{R}^d, \Omega(\theta) \leq D\}$ 

$$f(\hat{\theta}) - \min_{\theta \in \mathbb{R}^d} f(\theta) = \begin{bmatrix} f(\hat{\theta}) - \min_{\theta \in \Theta} f(\theta) \end{bmatrix} + \begin{bmatrix} \min_{\theta \in \Theta} f(\theta) - \min_{\theta \in \mathbb{R}^d} f(\theta) \end{bmatrix}$$
  
Estimation error Approximation error

**1**. Uniform deviation bounds, with  $\hat{\theta} \in \arg\min_{\theta \in \Theta} \hat{f}(\theta)$ 

$$\begin{split} f(\hat{\theta}) - \min_{\theta \in \Theta} f(\theta) &= \left[ f(\hat{\theta}) - \hat{f}(\hat{\theta}) \right] + \left[ \hat{f}(\hat{\theta}) - \hat{f}(\theta_{\Theta}^{*}) \right] + \left[ \hat{f}(\theta_{\Theta}^{*}) - f(\theta_{\Theta}^{*}) \right] \\ &\leqslant \sup_{\theta \in \Theta} |\hat{f}(\theta) - f(\theta)| + 0 \qquad + \sup_{\theta \in \Theta} |\hat{f}(\theta) - f(\theta)| \\ \end{split}$$

### Analysis of empirical risk minimization

• Approximation and estimation errors:  $\Theta = \{\theta \in \mathbb{R}^d, \Omega(\theta) \leq D\}$ 

$$f(\hat{\theta}) - \min_{\theta \in \mathbb{R}^d} f(\theta) = \begin{bmatrix} f(\hat{\theta}) - \min_{\theta \in \Theta} f(\theta) \end{bmatrix} + \begin{bmatrix} \min_{\theta \in \Theta} f(\theta) - \min_{\theta \in \mathbb{R}^d} f(\theta) \end{bmatrix}$$
  
Estimation error Approximation error  
Uniform deviation bounds, with  $\hat{\theta} \in \arg\min \hat{f}(\theta)$ 

**1**. Uniform deviation bounds, with  $\hat{\theta} \in \arg\min_{\theta \in \Theta} \hat{f}(\theta)$ 

$$f(\hat{\theta}) - \min_{\theta \in \Theta} f(\theta) \leq 2 \cdot \sup_{\theta \in \Theta} |\hat{f}(\theta) - f(\theta)|$$

– Typically slow rate  $O(1/\sqrt{n})$ 

**2**. More refined concentration results with faster rates O(1/n)

#### **Motivation from least-squares**

• For least-squares, we have  $\ell(y, \theta^{\top} \Phi(x)) = \frac{1}{2}(y - \theta^{\top} \Phi(x))^2$ , and

$$f(\theta) - \hat{f}(\theta) = \frac{1}{2} \theta^{\top} \left( \frac{1}{n} \sum_{i=1}^{n} \Phi(x_i) \Phi(x_i)^{\top} - \mathbb{E} \Phi(X) \Phi(X)^{\top} \right) \theta$$
$$-\theta^{\top} \left( \frac{1}{n} \sum_{i=1}^{n} y_i \Phi(x_i) - \mathbb{E} Y \Phi(X) \right) + \frac{1}{2} \left( \frac{1}{n} \sum_{i=1}^{n} y_i^2 - \mathbb{E} Y^2 \right),$$
$$\sup_{\|\theta\|_2 \leq D} |f(\theta) - \hat{f}(\theta)| \leq \frac{D^2}{2} \left\| \frac{1}{n} \sum_{i=1}^{n} \Phi(x_i) \Phi(x_i)^{\top} - \mathbb{E} \Phi(X) \Phi(X)^{\top} \right\|_{\text{op}}$$
$$+ D^{\left\| \frac{1}{2} \sum_{i=1}^{n} \Phi(x_i) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Phi(x_i) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Phi(x_i) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Phi(x_i) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} + \frac{1}{2} \left\| \frac{1}{2} \sum_{i=1}^{n} \Psi(X) - \mathbb{E} Y \Phi(X) \right\|_{\infty} +$$

$$+D\left\|\frac{1}{n}\sum_{i=1}^{n}y_{i}\Phi(x_{i}) - \mathbb{E}Y\Phi(X)\right\|_{2} + \frac{1}{2}\left|\frac{1}{n}\sum_{i=1}^{n}y_{i}^{2} - \mathbb{E}Y^{2}\right|,$$

 $\sup_{\|\theta\|_2 \leqslant D} |f(\theta) - \hat{f}(\theta)| \leqslant O(1/\sqrt{n}) \text{ with high probability}$ 

## Slow rate for supervised learning

- Assumptions (f is the expected risk,  $\hat{f}$  the empirical risk)
  - $\Omega(\theta) = \|\theta\|_2$  (Euclidean norm)
  - "Linear" predictors:  $\theta(x) = \theta^{\top} \Phi(x)$ , with  $\|\Phi(x)\|_2 \leq R$  a.s.
  - G-Lipschitz loss: f and  $\hat{f}$  are GR-Lipschitz on  $\Theta = \{ \|\theta\|_2 \leq D \}$
  - No assumptions regarding convexity
## Slow rate for supervised learning

- Assumptions (f is the expected risk,  $\hat{f}$  the empirical risk)
  - $\Omega(\theta) = \|\theta\|_2$  (Euclidean norm)
  - "Linear" predictors:  $\theta(x) = \theta^{\top} \Phi(x)$ , with  $\|\Phi(x)\|_2 \leq R$  a.s.
  - G-Lipschitz loss: f and  $\hat{f}$  are GR-Lipschitz on  $\Theta = \{ \|\theta\|_2 \leq D \}$
  - No assumptions regarding convexity
- $\bullet$  With probability greater than  $1-\delta$

$$\sup_{\theta \in \Theta} |\hat{f}(\theta) - f(\theta)| \leqslant \frac{GRD}{\sqrt{n}} \left[ 2 + \sqrt{2\log\frac{2}{\delta}} \right]$$

- Expectated estimation error:  $\mathbb{E}\left[\sup_{\theta\in\Theta}|\hat{f}(\theta) f(\theta)|\right] \leqslant \frac{4GRD}{\sqrt{n}}$
- Using Rademacher averages (see, e.g., Boucheron et al., 2005)
- Lipschitz functions  $\Rightarrow$  slow rate

### Symmetrization with Rademacher variables

• Let  $\mathcal{D}' = \{x'_1, y'_1, \dots, x'_n, y'_n\}$  an independent copy of the data  $\mathcal{D} = \{x_1, y_1, \dots, x_n, y_n\}$ , with corresponding loss functions  $f'_i(\theta)$ 

$$\begin{split} \mathbb{E}\left[\sup_{\theta\in\Theta}\left|f(\theta)-\hat{f}(\theta)\right|\right] &= \mathbb{E}\left[\sup_{\theta\in\Theta}\left(f(\theta)-\frac{1}{n}\sum_{i=1}^{n}f_{i}(\theta)\right)\right] \\ &= \mathbb{E}\left[\sup_{\theta\in\Theta}\left|\frac{1}{n}\sum_{i=1}^{n}\mathbb{E}\left(f'_{i}(\theta)-f_{i}(\theta)\right)\right|\right] \\ &\leqslant \mathbb{E}\left[\mathbb{E}\left[\left[\sup_{\theta\in\Theta}\left|\frac{1}{n}\sum_{i=1}^{n}\left(f'_{i}(\theta)-f_{i}(\theta)\right)\right|\right]\right] \\ &= \mathbb{E}\left[\left[\sup_{\theta\in\Theta}\left|\frac{1}{n}\sum_{i=1}^{n}\varepsilon_{i}\left(f'_{i}(\theta)-f_{i}(\theta)\right)\right|\right] \\ &= \mathbb{E}\left[\left[\sup_{\theta\in\Theta}\left|\frac{1}{n}\sum_{i=1}^{n}\varepsilon_{i}\left(f'_{i}(\theta)-f_{i}(\theta)\right)\right|\right] \\ &\leqslant 2\mathbb{E}\left[\left[\sup_{\theta\in\Theta}\left|\frac{1}{n}\sum_{i=1}^{n}\varepsilon_{i}f_{i}(\theta)\right|\right] = \mathsf{Rademacher complexity} \end{split}$$

## **Rademacher complexity**

• Define the Rademacher complexity of the class of functions  $(X,Y)\mapsto \ell(Y,\theta^{\top}\Phi(X))$  as

$$R_n = \mathbb{E}\left[\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f_i(\theta) \right| \right].$$

- Note two expectations, with respect to  ${\cal D}$  and with respect to  $\varepsilon$
- Main property:

$$\mathbb{E}\Big[\sup_{\theta\in\Theta}\left|f(\theta)-\hat{f}(\theta)\right|\Big]\leqslant 2R_n$$

### From Rademacher complexity to uniform bound

- Let  $Z = \sup_{\theta \in \Theta} \left| f(\theta) \hat{f}(\theta) \right|$
- By changing the pair  $(x_i, y_i)$ , Z may only change by

$$\frac{2}{n} \sup |\ell(Y, \theta^{\top} \Phi(X))| \leq \frac{2}{n} (\sup |\ell(Y, 0)| + GRD) \leq \frac{2}{n} (\ell_0 + GRD) = c$$
  
with  $\sup |\ell(Y, 0)| = \ell_0$ 

• MacDiarmid inequality: with probability greater than  $1-\delta$ ,

$$Z \leqslant \mathbb{E}Z + \sqrt{\frac{n}{2}}c \cdot \sqrt{\log\frac{1}{\delta}} \leqslant 2R_n + \frac{\sqrt{2}}{\sqrt{n}}(\ell_0 + GRD)\sqrt{\log\frac{1}{\delta}}$$

#### Bounding the Rademacher average - I

• We have, with  $\varphi_i(u) = \ell(y_i, u) - \ell(y_i, 0)$  is almost surely *B*-Lipschitz:

$$\begin{aligned} R_n &= \mathbb{E} \left[ \sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f_i(\theta) \right| \right] \\ &\leqslant \mathbb{E} \left[ \sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f_i(0) \right| \right] + \mathbb{E} \left[ \sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i [f_i(\theta) - f_i(0)] \right| \right] \\ &\leqslant \frac{\ell_0}{\sqrt{n}} + \mathbb{E} \left[ \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \varepsilon_i [f_i(\theta) - f_i(0)] \right] \\ &= \frac{\ell_0}{\sqrt{n}} + \mathbb{E} \left[ \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \varphi_i(\theta^\top \Phi(x_i)) \right] \end{aligned}$$

• Using Ledoux-Talagrand concentration results for Rademacher averages (since  $\varphi_i$  is G-Lipschitz, we get:

$$R_n \leqslant \frac{\ell_0}{\sqrt{n}} + 2G \cdot \mathbb{E} \left[ \sup_{\|\theta\|_2 \leqslant D} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \theta^\top \Phi(x_i) \right] \right|$$

# Bounding the Rademacher average - II

• We have:

$$\begin{split} R_n &\leqslant \frac{\ell_0}{\sqrt{n}} + 2G\mathbb{E}\bigg[\sup_{\|\theta\|_2 \leqslant D} \bigg| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \theta^\top \Phi(x_i) \bigg] \bigg| \\ &= \frac{\ell_0}{\sqrt{n}} + 2G\mathbb{E} \bigg\| D \frac{1}{n} \sum_{i=1}^n \varepsilon_i \Phi(x_i) \bigg\|_2 \\ &\leqslant \frac{\ell_0}{\sqrt{n}} + 2GD \sqrt{\mathbb{E} \bigg\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i \Phi(x_i) \bigg\|_2^2} \text{ by Jensen's inequality} \\ &\leqslant \frac{2(\ell_0 + GRD)}{\sqrt{n}} \text{ by using} \|\Phi(x)\|_2 \leqslant R \end{split}$$

• Overall, we get, with probability  $1 - \delta$ :

$$\sup_{\theta \in \Theta} \left| f(\theta) - \hat{f}(\theta) \right| \leq \frac{1}{\sqrt{n}} \left( \ell_0 + GRD \right) \left( 4 + \sqrt{2\log\frac{1}{\delta}} \right)$$

## Putting it all together

- $\bullet$  We have, with probability  $1-\delta$ 
  - For exact minimizer  $\hat{\theta} \in \arg \min_{\theta \in \Theta} \hat{f}(\theta)$ , we have

$$f(\hat{\theta}) - \min_{\theta \in \Theta} f(\theta) \leq 2 \cdot \sup_{\theta \in \Theta} |\hat{f}(\theta) - f(\theta)|$$
$$\leq \frac{2}{\sqrt{n}} \left(\ell_0 + GRD\right) \left(4 + \sqrt{2\log\frac{1}{\delta}}\right)$$

– For all  $\theta \in \Theta$ 

$$f(\theta) - \min_{\theta \in \Theta} f(\theta) \leq 2 \cdot \sup_{\theta \in \Theta} |\hat{f}(\theta) - f(\theta)| + \left[\hat{f}(\theta) - \hat{f}(\hat{\theta})\right]$$

• Only need to optimize with precision  $\frac{2}{\sqrt{n}}(\ell_0 + GRD)$ 

## Slow rate for supervised learning (summary)

- Assumptions (f is the expected risk,  $\hat{f}$  the empirical risk)
  - $\Omega(\theta) = \|\theta\|_2$  (Euclidean norm)
  - "Linear" predictors:  $\theta(x) = \theta^{\top} \Phi(x)$ , with  $\|\Phi(x)\|_2 \leq R$  a.s.
  - G-Lipschitz loss: f and  $\hat{f}$  are GR-Lipschitz on  $\Theta = \{ \|\theta\|_2 \leq D \}$
  - No assumptions regarding convexity
- $\bullet$  With probability greater than  $1-\delta$

$$\sup_{\theta \in \Theta} |\hat{f}(\theta) - f(\theta)| \leq \frac{(\ell_0 + GRD)}{\sqrt{n}} \left[ 2 + \sqrt{2\log\frac{2}{\delta}} \right]$$

- Expectated estimation error:  $\mathbb{E} \left[ \sup_{\theta \in \Theta} |\hat{f}(\theta) f(\theta)| \right] \leq \frac{4(\ell_0 + GRD)}{\sqrt{n}}$
- Using Rademacher averages (see, e.g., Boucheron et al., 2005)
- Lipschitz functions  $\Rightarrow$  slow rate

### **Motivation from mean estimation**

- Estimator  $\hat{\theta} = \frac{1}{n} \sum_{i=1}^{n} z_i = \arg\min_{\theta \in \mathbb{R}} \frac{1}{2n} \sum_{i=1}^{n} (\theta z_i)^2 = \hat{f}(\theta)$
- From before:

$$- f(\theta) = \frac{1}{2}\mathbb{E}(\theta - z)^2 = \frac{1}{2}(\theta - \mathbb{E}z)^2 + \frac{1}{2}\operatorname{var}(z) = \hat{f}(\theta) + O(1/\sqrt{n}) - f(\hat{\theta}) = \frac{1}{2}(\hat{\theta} - \mathbb{E}z)^2 + \frac{1}{2}\operatorname{var}(z) = f(\mathbb{E}z) + O(1/\sqrt{n})$$

### **Motivation from mean estimation**

- Estimator  $\hat{\theta} = \frac{1}{n} \sum_{i=1}^{n} z_i = \arg\min_{\theta \in \mathbb{R}} \frac{1}{2n} \sum_{i=1}^{n} (\theta z_i)^2 = \hat{f}(\theta)$
- From before:

$$-f(\theta) = \frac{1}{2}\mathbb{E}(\theta - z)^2 = \frac{1}{2}(\theta - \mathbb{E}z)^2 + \frac{1}{2}\operatorname{var}(z) = \hat{f}(\theta) + O(1/\sqrt{n}) \\ -f(\hat{\theta}) = \frac{1}{2}(\hat{\theta} - \mathbb{E}z)^2 + \frac{1}{2}\operatorname{var}(z) = f(\mathbb{E}z) + O(1/\sqrt{n})$$

• More refined/direct bound:

$$f(\hat{\theta}) - f(\mathbb{E}z) = \frac{1}{2}(\hat{\theta} - \mathbb{E}z)^2$$
$$\mathbb{E}[f(\hat{\theta}) - f(\mathbb{E}z)] = \frac{1}{2}\mathbb{E}\left(\frac{1}{n}\sum_{i=1}^n z_i - \mathbb{E}z\right)^2 = \frac{1}{2n}\operatorname{var}(z)$$

• Bound only at  $\hat{\theta}$  + strong convexity (instead of uniform bound)

## Fast rate for supervised learning

- Assumptions (f is the expected risk,  $\hat{f}$  the empirical risk)
  - Same as before (bounded features, Lipschitz loss)
  - Regularized risks:  $f^{\mu}(\theta) = f(\theta) + \frac{\mu}{2} \|\theta\|_2^2$  and  $\hat{f}^{\mu}(\theta) = \hat{f}(\theta) + \frac{\mu}{2} \|\theta\|_2^2$
  - Convexity
- For any a > 0, with probability greater than  $1 \delta$ , for all  $\theta \in \mathbb{R}^d$ ,

$$f^{\mu}(\theta) - \min_{\eta \in \mathbb{R}^d} f^{\mu}(\eta) \leqslant (1+a)(\hat{f}^{\mu}(\theta) - \min_{\eta \in \mathbb{R}^d} \hat{f}^{\mu}(\eta)) + \frac{8(1+\frac{1}{a})G^2R^2(32+\log\frac{1}{\delta})}{\mu n}$$

- Results from Sridharan, Srebro, and Shalev-Shwartz (2008)
  - see also Boucheron and Massart (2011) and references therein
- Strongly convex functions  $\Rightarrow$  fast rate
  - Warning:  $\mu$  should decrease with n to reduce approximation error

# Outline

### 1. Large-scale machine learning and optimization

- Traditional statistical analysis
- Classical methods for convex optimization

## 2. Non-smooth stochastic approximation

- Stochastic (sub)gradient and averaging
- Non-asymptotic results and lower bounds
- Strongly convex vs. non-strongly convex

### 3. Smooth stochastic approximation algorithms

- Asymptotic and non-asymptotic results
- 4. Beyond decaying step-sizes
- 5. Finite data sets

## **Complexity results in convex optimization**

- Assumption: f convex on  $\mathbb{R}^d$
- Classical generic algorithms
  - (sub)gradient method/descent
  - Accelerated gradient descent
  - Newton method
- $\bullet$  Key additional properties of f
  - Lipschitz continuity, smoothness or strong convexity
- Key insight from Bottou and Bousquet (2008)
  - In machine learning, no need to optimize below estimation error
- Key reference: Nesterov (2004)

# (smooth) gradient descent

#### • Assumptions

- f convex with L-Lipschitz-continuous gradient
- Minimum attained at  $\theta_*$
- Algorithm:

$$\theta_t = \theta_{t-1} - \frac{1}{L}f'(\theta_{t-1})$$



# (smooth) gradient descent

#### • Assumptions

- f convex with L-Lipschitz-continuous gradient
- Minimum attained at  $\theta_*$
- Algorithm:

$$\theta_t = \theta_{t-1} - \frac{1}{L}f'(\theta_{t-1})$$

• Bound:

$$f(\theta_t) - f(\theta_*) \leqslant \frac{2L\|\theta_0 - \theta_*\|^2}{t+4}$$

- Three-line proof
- Not best possible convergence rate after O(d) iterations

# (smooth) gradient descent - strong convexity

#### • Assumptions

- f convex with L-Lipschitz-continuous gradient
- $f \mu$ -strongly convex
- Algorithm:

$$\theta_t = \theta_{t-1} - \frac{1}{L}f'(\theta_{t-1})$$

• Bound:

$$f(\theta_t) - f(\theta_*) \leqslant (1 - \mu/L)^t [f(\theta_0) - f(\theta_*)]$$

- Three-line proof
- Adaptivity of gradient descent to problem difficulty
- Line search

### **Gradient descent - Proof for quadratic functions**

- Quadratic convex function:  $f(\theta) = \frac{1}{2}\theta^{\top}H\theta c^{\top}\theta$ 
  - $\mu$  and L are smallest largest eigenvalues of H
  - Global optimum  $\theta_* = H^{-1}c$  (or  $H^{\dagger}c$ )
- Gradient descent:

$$\theta_t = \theta_{t-1} - \frac{1}{L}(H\theta - c) = \theta_{t-1} - \frac{1}{L}(H\theta - H\theta_*)$$
  
$$\theta_t - \theta_* = (I - \frac{1}{L}H)(\theta_{t-1} - \theta_*) = (I - \frac{1}{L}H)^t(\theta_0 - \theta_*)$$

- Strong convexity  $\mu > 0$ : eigenvalues of  $(I \frac{1}{L}H)^t$  in  $[0, (1 \frac{\mu}{L})^t]$ 
  - Convergence of iterates:  $\|\theta_t \theta_*\|^2 \leq (1 \mu/L)^{2t} \|\theta_0 \theta_*\|^2$
  - Function values:  $f(\theta_t) f(\theta_*) \leq (1 \mu/L)^{2t} [f(\theta_0) f(\theta_*)]$

### **Gradient descent - Proof for quadratic functions**

- Quadratic convex function:  $f(\theta) = \frac{1}{2}\theta^{\top}H\theta c^{\top}\theta$ 
  - $\mu$  and L are smallest largest eigenvalues of H
  - Global optimum  $\theta_* = H^{-1}c$  (or  $H^{\dagger}c$ )
- Gradient descent:

$$\theta_t = \theta_{t-1} - \frac{1}{L}(H\theta - c) = \theta_{t-1} - \frac{1}{L}(H\theta - H\theta_*)$$
  
$$\theta_t - \theta_* = (I - \frac{1}{L}H)(\theta_{t-1} - \theta_*) = (I - \frac{1}{L}H)^t(\theta_0 - \theta_*)$$

- Convexity  $\mu = 0$ : eigenvalues of  $(I \frac{1}{L}H)^t$  in [0, 1]
  - No convergence of iterates:  $\|\theta_t \theta_*\|^2 \leq \|\theta_0 \theta_*\|^2$
  - Function values:  $f(\theta_t) f(\theta_*) \leq \max_{e \in [0,L]} e(1 e/L)^{2t} \|\theta_0 \theta_*\|^2$  $f(\theta_t) - f(\theta_*) \leq \frac{L}{t} \|\theta_0 - \theta_*\|^2$

## Accelerated gradient methods (Nesterov, 1983)

#### • Assumptions

– f convex with L-Lipschitz-cont. gradient , min. attained at  $\theta_*$ 

• Algorithm:

$$\theta_t = \eta_{t-1} - \frac{1}{L} f'(\eta_{t-1})$$
  
$$\eta_t = \theta_t + \frac{t-1}{t+2} (\theta_t - \theta_{t-1})$$

• Bound:

$$f(\theta_t) - f(\theta_*) \leqslant \frac{2L\|\theta_0 - \theta_*\|^2}{(t+1)^2}$$

- Ten-line proof (see, e.g., Schmidt, Le Roux, and Bach, 2011)
- Not improvable
- Extension to strongly convex functions

# **Optimization for sparsity-inducing norms** (see Bach, Jenatton, Mairal, and Obozinski, 2012b)

• Gradient descent as a **proximal method** (differentiable functions)

$$-\theta_{t+1} = \arg\min_{\theta \in \mathbb{R}^d} f(\theta_t) + (\theta - \theta_t)^\top \nabla f(\theta_t) + \frac{L}{2} \|\theta - \theta_t\|_2^2$$
$$-\theta_{t+1} = \theta_t - \frac{1}{L} \nabla f(\theta_t)$$

# **Optimization for sparsity-inducing norms** (see Bach, Jenatton, Mairal, and Obozinski, 2012b)

• Gradient descent as a **proximal method** (differentiable functions)

$$-\theta_{t+1} = \arg\min_{\theta \in \mathbb{R}^d} f(\theta_t) + (\theta - \theta_t)^\top \nabla f(\theta_t) + \frac{L}{2} \|\theta - \theta_t\|_2^2$$
$$-\theta_{t+1} = \theta_t - \frac{1}{L} \nabla f(\theta_t)$$

• Problems of the form:  $\left| \min_{\theta \in \mathbb{R}^d} f(\theta) + \mu \Omega(\theta) \right|$ 

$$-\theta_{t+1} = \arg\min_{\theta \in \mathbb{R}^d} f(\theta_t) + (\theta - \theta_t)^\top \nabla f(\theta_t) + \mu \Omega(\theta) + \frac{L}{2} \|\theta - \theta_t\|_2^2$$
$$-\Omega(\theta) = \|\theta\|_1 \Rightarrow \text{Thresholded gradient descent}$$

- Similar convergence rates than smooth optimization
  - Acceleration methods (Nesterov, 2007; Beck and Teboulle, 2009)

# Subgradient method/"descent" (Shor et al., 1985)

• Assumptions

- f convex and B-Lipschitz-continuous on  $\{\|\theta\|_2 \leq D\}$ 

- Algorithm:  $\theta_t = \Pi_D \left( \theta_{t-1} \frac{2D}{B\sqrt{t}} f'(\theta_{t-1}) \right)$ 
  - $\Pi_D$ : orthogonal projection onto  $\{\|\theta\|_2 \leq D\}$



# Subgradient method/"descent" (Shor et al., 1985)

- Assumptions
  - f convex and B-Lipschitz-continuous on  $\{\|\theta\|_2 \leq D\}$

• Algorithm: 
$$\theta_t = \Pi_D \left( \theta_{t-1} - \frac{2D}{B\sqrt{t}} f'(\theta_{t-1}) \right)$$

- $\Pi_D$ : orthogonal projection onto  $\{\|\theta\|_2 \leq D\}$
- Bound:

$$f\left(\frac{1}{t}\sum_{k=0}^{t-1}\theta_k\right) - f(\theta_*) \leqslant \frac{2DB}{\sqrt{t}}$$

- Three-line proof
- Best possible convergence rate after O(d) iterations

## Subgradient method/"descent" - proof - I

• Iteration: 
$$\theta_t = \prod_D (\theta_{t-1} - \gamma_t f'(\theta_{t-1}))$$
 with  $\gamma_t = \frac{2D}{B\sqrt{t}}$ 

• Assumption:  $||f'(\theta)||_2 \leq B$  and  $||\theta||_2 \leq D$ 

$$\begin{aligned} \|\theta_t - \theta_*\|_2^2 &\leqslant \|\theta_{t-1} - \theta_* - \gamma_t f'(\theta_{t-1})\|_2^2 \text{ by contractivity of projections} \\ &\leqslant \|\theta_{t-1} - \theta_*\|_2^2 + B^2 \gamma_t^2 - 2\gamma_t (\theta_{t-1} - \theta_*)^\top f'(\theta_{t-1}) \text{ because } \|f'(\theta_{t-1})\|_2 \leqslant B \\ &\leqslant \|\theta_{t-1} - \theta_*\|_2^2 + B^2 \gamma_t^2 - 2\gamma_t \big[f(\theta_{t-1}) - f(\theta_*)\big] \text{ (property of subgradients)} \end{aligned}$$

• leading to

$$f(\theta_{t-1}) - f(\theta_*) \leqslant \frac{B^2 \gamma_t}{2} + \frac{1}{2\gamma_t} \left[ \|\theta_{t-1} - \theta_*\|_2^2 - \|\theta_t - \theta_*\|_2^2 \right]$$

## Subgradient method/"descent" - proof - II

- Starting from  $f(\theta_{t-1}) f(\theta_*) \leq \frac{B^2 \gamma_t}{2} + \frac{1}{2\gamma_t} \left[ \|\theta_{t-1} \theta_*\|_2^2 \|\theta_t \theta_*\|_2^2 \right]$
- Constant step-size  $\gamma_t = \gamma$

$$\sum_{u=1}^{t} \left[ f(\theta_{u-1}) - f(\theta_{*}) \right] \leqslant \sum_{u=1}^{t} \frac{B^{2}\gamma}{2} + \sum_{u=1}^{t} \frac{1}{2\gamma} \left[ \|\theta_{u-1} - \theta_{*}\|_{2}^{2} - \|\theta_{u} - \theta_{*}\|_{2}^{2} \right]$$
$$\leqslant t \frac{B^{2}\gamma}{2} + \frac{1}{2\gamma} \|\theta_{0} - \theta_{*}\|_{2}^{2} \leqslant t \frac{B^{2}\gamma}{2} + \frac{2}{\gamma} D^{2}$$

• Optimized step-size  $\gamma_t = \frac{2D}{B\sqrt{t}}$  depends on "horizon"

– Leads to bound of  $2DB\sqrt{t}$ 

• Using convexity: 
$$f\left(\frac{1}{t}\sum_{k=0}^{t-1}\theta_k\right) - f(\theta_*) \leqslant \frac{2DB}{\sqrt{t}}$$

## Subgradient method/"descent" - proof - III

• Starting from 
$$f(\theta_{t-1}) - f(\theta_*) \leq \frac{B^2 \gamma_t}{2} + \frac{1}{2\gamma_t} \left[ \|\theta_{t-1} - \theta_*\|_2^2 - \|\theta_t - \theta_*\|_2^2 \right]$$

• Decreasing step-size

$$\begin{split} \sum_{u=1}^{t} \left[ f(\theta_{u-1}) - f(\theta_{*}) \right] &\leqslant \quad \sum_{u=1}^{t} \frac{B^{2} \gamma_{u}}{2} + \sum_{u=1}^{t} \frac{1}{2 \gamma_{u}} \left[ \|\theta_{u-1} - \theta_{*}\|_{2}^{2} - \|\theta_{u} - \theta_{*}\|_{2}^{2} \right] \\ &= \sum_{u=1}^{t} \frac{B^{2} \gamma_{u}}{2} + \sum_{u=1}^{t-1} \|\theta_{u} - \theta_{*}\|_{2}^{2} \left(\frac{1}{2 \gamma_{u+1}} - \frac{1}{2 \gamma_{u}}\right) + \frac{\|\theta_{0} - \theta_{*}\|_{2}^{2}}{2 \gamma_{1}} - \frac{\|\theta_{t} - \theta_{*}\|_{2}^{2}}{2 \gamma_{t}} \\ &\leqslant \quad \sum_{u=1}^{t} \frac{B^{2} \gamma_{u}}{2} + \sum_{u=1}^{t-1} 4D^{2} \left(\frac{1}{2 \gamma_{u+1}} - \frac{1}{2 \gamma_{u}}\right) + \frac{4D^{2}}{2 \gamma_{1}} \\ &= \quad \sum_{u=1}^{t} \frac{B^{2} \gamma_{u}}{2} + \frac{4D^{2}}{2 \gamma_{t}} \leqslant 2DB\sqrt{t} \text{ with } \gamma_{t} = \frac{2D}{B\sqrt{t}} \end{split}$$

• Using convexity:  $f(\frac{1}{t}\sum_{k=0}^{t-1}\theta_k) - f(\theta_*) \leq \frac{2DB}{\sqrt{t}}$ 

## Subgradient descent for machine learning

- Assumptions (f is the expected risk,  $\hat{f}$  the empirical risk)
  - "Linear" predictors:  $\theta(x) = \theta^{\top} \Phi(x)$ , with  $\|\Phi(x)\|_2 \leq R$  a.s. -  $\hat{f}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \Phi(x_i)^{\top} \theta)$
  - G-Lipschitz loss: f and  $\hat{f}$  are GR-Lipschitz on  $\Theta = \{ \|\theta\|_2 \leq D \}$
- Statistics: with probability greater than  $1 \delta$  $\sup_{\theta \in \Theta} |\hat{f}(\theta) - f(\theta)| \leqslant \frac{GRD}{\sqrt{n}} \left[ 2 + \sqrt{2\log\frac{2}{\delta}} \right]$
- **Optimization**: after t iterations of subgradient method

$$\hat{f}(\hat{\theta}) - \min_{\eta \in \Theta} \hat{f}(\eta) \leqslant \frac{GRD}{\sqrt{t}}$$

• t = n iterations, with total running-time complexity of  $O(n^2d)$ 

## Subgradient descent - strong convexity

#### • Assumptions

- f convex and B-Lipschitz-continuous on  $\{\|\theta\|_2 \leq D\}$ -  $f \mu$ -strongly convex

• Algorithm: 
$$\theta_t = \Pi_D \left( \theta_{t-1} - \frac{2}{\mu(t+1)} f'(\theta_{t-1}) \right)$$

• Bound:

$$f\left(\frac{2}{t(t+1)}\sum_{k=1}^{t}k\theta_{k-1}\right) - f(\theta_*) \leqslant \frac{2B^2}{\mu(t+1)}$$

- Three-line proof
- Best possible convergence rate after O(d) iterations

#### Subgradient method - strong convexity - proof - I

• Iteration: 
$$\theta_t = \prod_D(\theta_{t-1} - \gamma_t f'(\theta_{t-1}))$$
 with  $\gamma_t = \frac{2}{\mu(t+1)}$ 

• Assumption:  $||f'(\theta)||_2 \leq B$  and  $||\theta||_2 \leq D$  and  $\mu$ -strong convexity of f

• leading to

$$f(\theta_{t-1}) - f(\theta_{*}) \leqslant \frac{B^{2}\gamma_{t}}{2} + \frac{1}{2} \Big[ \frac{1}{\gamma_{t}} - \mu \Big] \|\theta_{t-1} - \theta_{*}\|_{2}^{2} - \frac{1}{2\gamma_{t}} \|\theta_{t} - \theta_{*}\|_{2}^{2}$$
$$\leqslant \frac{B^{2}}{\mu(t+1)} + \frac{\mu}{2} \Big[ \frac{t-1}{2} \Big] \|\theta_{t-1} - \theta_{*}\|_{2}^{2} - \frac{\mu(t+1)}{4} \|\theta_{t} - \theta_{*}\|_{2}^{2}$$

### Subgradient method - strong convexity - proof - II

• From 
$$f(\theta_{t-1}) - f(\theta_*) \leq \frac{B^2}{\mu(t+1)} + \frac{\mu}{2} \left[\frac{t-1}{2}\right] \|\theta_{t-1} - \theta_*\|_2^2 - \frac{\mu(t+1)}{4} \|\theta_t - \theta_*\|_2^2$$

$$\sum_{u=1}^{t} u \left[ f(\theta_{u-1}) - f(\theta_{*}) \right] \leq \sum_{t=1}^{u} \frac{B^{2}u}{\mu(u+1)} + \frac{1}{4} \sum_{u=1}^{t} \left[ u(u-1) \|\theta_{u-1} - \theta_{*}\|_{2}^{2} - u(u+1) \|\theta_{u} - \theta_{*}\|_{2}^{2} \right]$$
$$\leq \frac{B^{2}t}{\mu} + \frac{1}{4} \left[ 0 - t(t+1) \|\theta_{t} - \theta_{*}\|_{2}^{2} \right] \leq \frac{B^{2}t}{\mu}$$

• Using convexity: 
$$f\left(\frac{2}{t(t+1)}\sum_{u=1}^{t}u\theta_{u-1}\right) - f(\theta_*) \leqslant \frac{2B^2}{t+1}$$

• NB: with step-size  $\gamma_n = 1/(n\mu)$ , extra logarithmic factor

## Summary: minimizing convex functions

- Assumption: *f* convex
- Gradient descent:  $\theta_t = \theta_{t-1} \gamma_t f'(\theta_{t-1})$ 
  - $-O(1/\sqrt{t})$  convergence rate for non-smooth convex functions -O(1/t) convergence rate for smooth convex functions  $-O(e^{-\rho t})$  convergence rate for strongly smooth convex functions
- Newton method:  $\theta_t = \theta_{t-1} f''(\theta_{t-1})^{-1}f'(\theta_{t-1})$ 
  - $O(e^{-\rho 2^t})$  convergence rate

## **Summary: minimizing convex functions**

- Assumption: *f* convex
- Gradient descent:  $\theta_t = \theta_{t-1} \gamma_t f'(\theta_{t-1})$ 
  - $\begin{array}{l} \ O(1/\sqrt{t}) \ {\rm convergence} \ {\rm rate} \ {\rm for} \ {\rm non-smooth} \ {\rm convex} \ {\rm functions} \\ \ O(1/t) \ \ {\rm convergence} \ {\rm rate} \ {\rm for} \ {\rm smooth} \ {\rm convex} \ {\rm functions} \\ \ O(e^{-\rho t}) \ \ {\rm convergence} \ {\rm rate} \ {\rm for} \ {\rm strongly} \ {\rm smooth} \ {\rm convex} \ {\rm functions} \end{array}$
- Newton method:  $\theta_t = \theta_{t-1} f''(\theta_{t-1})^{-1}f'(\theta_{t-1})$ 
  - $-O(e^{-\rho 2^t})$  convergence rate
- Key insights from Bottou and Bousquet (2008)
  - In machine learning, no need to optimize below statistical error
    In machine learning, cost functions are averages

### $\Rightarrow$ Stochastic approximation

# Outline

### 1. Large-scale machine learning and optimization

- Traditional statistical analysis
- Classical methods for convex optimization

## 2. Non-smooth stochastic approximation

- Stochastic (sub)gradient and averaging
- Non-asymptotic results and lower bounds
- Strongly convex vs. non-strongly convex

### 3. Smooth stochastic approximation algorithms

- Asymptotic and non-asymptotic results
- 4. Beyond decaying step-sizes
- 5. Finite data sets

## **Stochastic approximation**

- **Goal**: Minimizing a function f defined on  $\mathbb{R}^d$ 
  - given only unbiased estimates  $f_n'(\theta_n)$  of its gradients  $f'(\theta_n)$  at certain points  $\theta_n\in\mathbb{R}^d$

## **Stochastic approximation**

- **Goal**: Minimizing a function f defined on  $\mathbb{R}^d$ 
  - given only unbiased estimates  $f_n'(\theta_n)$  of its gradients  $f'(\theta_n)$  at certain points  $\theta_n\in\mathbb{R}^d$
- Machine learning statistics
  - loss for a single pair of observations:

$$f_n(\theta) = \ell(y_n, \theta^\top \Phi(x_n))$$

- $-f(\theta) = \mathbb{E}f_n(\theta) = \mathbb{E}\ell(y_n, \theta^{\top}\Phi(x_n)) = \mathbf{g}$ eneralization error
- Expected gradient:  $f'(\theta) = \mathbb{E}f'_n(\theta) = \mathbb{E}\left\{\ell'(y_n, \theta^\top \Phi(x_n)) \Phi(x_n)\right\}$
- Non-asymptotic results

#### • Number of iterations = number of observations

## **Stochastic approximation**

- **Goal**: Minimizing a function f defined on  $\mathbb{R}^d$ 
  - given only unbiased estimates  $f'_n(\theta_n)$  of its gradients  $f'(\theta_n)$  at certain points  $\theta_n \in \mathbb{R}^d$
- Stochastic approximation
  - (much) broader applicability beyond convex optimization

$$\theta_n = \theta_{n-1} - \gamma_n h_n(\theta_{n-1})$$
 with  $\mathbb{E}[h_n(\theta_{n-1})|\theta_{n-1}] = h(\theta_{n-1})$ 

- Beyond convex problems, i.i.d assumption, finite dimension, etc.
- Typically asymptotic results
- See, e.g., Kushner and Yin (2003); Benveniste et al. (2012)
## **Relationship to online learning**

- Stochastic approximation
  - Minimize  $f(\theta) = \mathbb{E}_z \ell(\theta, z) =$  generalization error of  $\theta$
  - Using the gradients of single i.i.d. observations

## **Relationship to online learning**

#### • Stochastic approximation

- Minimize  $f(\theta) = \mathbb{E}_z \ell(\theta, z) =$  generalization error of  $\theta$
- Using the gradients of single i.i.d. observations

## • Batch learning

- Finite set of observations:  $z_1, \ldots, z_n$
- Empirical risk:  $\hat{f}(\theta) = \frac{1}{n} \sum_{k=1}^{n} \ell(\theta, z_i)$
- Estimator  $\hat{\theta}$  = Minimizer of  $\hat{f}(\theta)$  over a certain class  $\Theta$
- Generalization bound using uniform concentration results

## **Relationship to online learning**

#### • Stochastic approximation

- Minimize  $f(\theta) = \mathbb{E}_z \ell(\theta, z) =$  generalization error of  $\theta$
- Using the gradients of single i.i.d. observations

### • Batch learning

- Finite set of observations:  $z_1, \ldots, z_n$
- Empirical risk:  $\hat{f}(\theta) = \frac{1}{n} \sum_{k=1}^{n} \ell(\theta, z_i)$
- Estimator  $\hat{\theta}$  = Minimizer of  $\hat{f}(\theta)$  over a certain class  $\Theta$
- Generalization bound using uniform concentration results

## • Online learning

- Update  $\hat{\theta}_n$  after each new (potentially adversarial) observation  $z_n$
- Cumulative loss:  $\frac{1}{n} \sum_{k=1}^{n} \ell(\hat{\theta}_{k-1}, z_k)$
- Online to batch through averaging (Cesa-Bianchi et al., 2004)

## **Convex stochastic approximation**

- Key properties of f and/or  $f_n$ 
  - Smoothness: f B-Lipschitz continuous, f' L-Lipschitz continuous
  - Strong convexity:  $f \mu$ -strongly convex

## **Convex stochastic approximation**

- Key properties of f and/or  $f_n$ 
  - Smoothness: f B-Lipschitz continuous, f' L-Lipschitz continuous
  - Strong convexity:  $f \mu$ -strongly convex
- Key algorithm: Stochastic gradient descent (a.k.a. Robbins-Monro)

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

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

$$\gamma_n = C n^{-\alpha}$$

## **Convex stochastic approximation**

- Key properties of f and/or  $f_n$ 
  - Smoothness: f B-Lipschitz continuous, f' L-Lipschitz continuous
  - Strong convexity:  $f \mu$ -strongly convex
- Key algorithm: Stochastic gradient descent (a.k.a. Robbins-Monro)

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

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

$$\gamma_n = C n^{-\alpha}$$

- Desirable practical behavior
  - Applicable (at least) to classical supervised learning problems
  - Robustness to (potentially unknown) constants (L,B, $\mu$ )
  - Adaptivity to difficulty of the problem (e.g., strong convexity)

## 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  $\mathbb{E}f_n = f$
- $\theta_*$  global optimum of f on  $\{\|\theta\|_2 \leq D\}$

• Algorithm: 
$$\theta_n = \prod_D \left( \theta_{n-1} - \frac{2D}{B\sqrt{n}} f'_n(\theta_{n-1}) \right)$$

## 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  $\mathbb{E}f_n = f$
- $\theta_*$  global optimum of f on  $\{\|\theta\|_2 \leq D\}$

• Algorithm: 
$$\theta_n = \Pi_D \left( \theta_{n-1} - \frac{2D}{B\sqrt{n}} f'_n(\theta_{n-1}) \right)$$

• Bound:

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

- "Same" three-line proof as in the deterministic case
- Minimax rate (Nemirovsky and Yudin, 1983; Agarwal et al., 2012)
- Running-time complexity: O(dn) after n iterations

#### Stochastic subgradient method - proof - I

• Iteration: 
$$\theta_n = \prod_D(\theta_{n-1} - \gamma_n f'_n(\theta_{n-1}))$$
 with  $\gamma_n = \frac{2D}{B\sqrt{n}}$ 

- $\mathcal{F}_n$  : information up to time n
- $||f'_n(\theta)||_2 \leq B$  and  $||\theta||_2 \leq D$ , unbiased gradients/functions  $\mathbb{E}(f_n|\mathcal{F}_{n-1}) = f$

$$\begin{aligned} \|\theta_n - \theta_*\|_2^2 &\leqslant \|\theta_{n-1} - \theta_* - \gamma_n f'_n(\theta_{n-1})\|_2^2 \text{ by contractivity of projections} \\ &\leqslant \|\theta_{n-1} - \theta_*\|_2^2 + B^2 \gamma_n^2 - 2\gamma_n (\theta_{n-1} - \theta_*)^\top f'_n(\theta_{n-1}) \text{ because } \|f'_n(\theta_{n-1})\|_2 \leqslant B \end{aligned}$$

$$\begin{split} \mathbb{E}\left[\|\theta_n - \theta_*\|_2^2 |\mathcal{F}_{n-1}\right] &\leqslant \|\theta_{n-1} - \theta_*\|_2^2 + B^2 \gamma_n^2 - 2\gamma_n (\theta_{n-1} - \theta_*)^\top f'(\theta_{n-1}) \\ &\leqslant \|\theta_{n-1} - \theta_*\|_2^2 + B^2 \gamma_n^2 - 2\gamma_n \left[f(\theta_{n-1}) - f(\theta_*)\right] \text{ (subgradient property)} \\ \mathbb{E}\|\theta_n - \theta_*\|_2^2 &\leqslant \mathbb{E}\|\theta_{n-1} - \theta_*\|_2^2 + B^2 \gamma_n^2 - 2\gamma_n \left[\mathbb{E}f(\theta_{n-1}) - f(\theta_*)\right] \end{split}$$

• 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 method - proof - II

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

$$\sum_{u=1}^{n} \left[ \mathbb{E}f(\theta_{u-1}) - f(\theta_{*}) \right] \leqslant \sum_{u=1}^{n} \frac{B^{2} \gamma_{u}}{2} + \sum_{u=1}^{n} \frac{1}{2\gamma_{u}} \left[ \mathbb{E} \|\theta_{u-1} - \theta_{*}\|_{2}^{2} - \mathbb{E} \|\theta_{u} - \theta_{*}\|_{2}^{2} \right]$$
$$\leqslant \sum_{u=1}^{n} \frac{B^{2} \gamma_{u}}{2} + \frac{4D^{2}}{2\gamma_{n}} \leqslant \frac{2DB}{\sqrt{n}} \text{ with } \gamma_{n} = \frac{2D}{B\sqrt{n}}$$

• Using convexity: 
$$\mathbb{E}f\left(\frac{1}{n}\sum_{k=0}^{n-1}\theta_k\right) - f(\theta_*) \leqslant \frac{2DB}{\sqrt{n}}$$

## Stochastic subgradient descent - strong convexity - I

## • Assumptions

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

• Algorithm: 
$$\theta_n = \prod_D \left( \theta_{n-1} - \frac{2}{\mu(n+1)} f'_n(\theta_{n-1}) \right)$$

• Bound:

$$\mathbb{E}f\left(\frac{2}{n(n+1)}\sum_{k=1}^{n}k\theta_{k-1}\right) - f(\theta_*) \leqslant \frac{2B^2}{\mu(n+1)}$$

- "Same" proof than deterministic case (Lacoste-Julien et al., 2012)
- Minimax rate (Nemirovsky and Yudin, 1983; Agarwal et al., 2012)

## Stochastic subgradient descent - strong convexity - II

#### • Assumptions

- $f_n$  convex and *B*-Lipschitz-continuous
- $(f_n)$  i.i.d. functions such that  $\mathbb{E}f_n = f$
- $\theta_*$  global optimum of  $g=f+\frac{\mu}{2}\|\cdot\|_2^2$
- No compactness assumption no projections
- Algorithm:

$$\theta_n = \theta_{n-1} - \frac{2}{\mu(n+1)} g'_n(\theta_{n-1}) = \theta_{n-1} - \frac{2}{\mu(n+1)} \left[ f'_n(\theta_{n-1}) + \mu \theta_{n-1} \right]$$

• Bound: 
$$\mathbb{E}g\left(\frac{2}{n(n+1)}\sum_{k=1}^{n}k\theta_{k-1}\right) - g(\theta_*) \leqslant \frac{2B^2}{\mu(n+1)}$$

• Minimax convergence rate

# Outline

#### 1. Large-scale machine learning and optimization

- Traditional statistical analysis
- Classical methods for convex optimization

## 2. Non-smooth stochastic approximation

- Stochastic (sub)gradient and averaging
- Non-asymptotic results and lower bounds
- Strongly convex vs. non-strongly convex

### 3. Smooth stochastic approximation algorithms

- Asymptotic and non-asymptotic results
- 4. Beyond decaying step-sizes
- 5. Finite data sets

- Known global minimax rates of convergence for non-smooth problems (Nemirovsky and Yudin, 1983; Agarwal et al., 2012)
  - Strongly convex:  $O((\mu n)^{-1})$

Attained by averaged stochastic gradient descent with  $\gamma_n \propto (\mu n)^{-1}$ 

– Non-strongly convex:  $O(n^{-1/2})$ 

Attained by averaged stochastic gradient descent with  $\gamma_n \propto n^{-1/2}$ 

- Known global minimax rates of convergence for non-smooth problems (Nemirovsky and Yudin, 1983; Agarwal et al., 2012)
  - Strongly convex:  $O((\mu n)^{-1})$

Attained by averaged stochastic gradient descent with  $\gamma_n \propto (\mu n)^{-1}$ 

- Non-strongly convex:  $O(n^{-1/2})$ Attained by averaged stochastic gradient descent with  $\gamma_n \propto n^{-1/2}$
- Many contributions in optimization and online learning: Bottou and Le Cun (2005); Bottou and Bousquet (2008); Hazan et al. (2007); Shalev-Shwartz and Srebro (2008); Shalev-Shwartz et al. (2007, 2009); Xiao (2010); Duchi and Singer (2009); Nesterov and Vial (2008); Nemirovski et al. (2009)

- Known global minimax rates of convergence for non-smooth problems (Nemirovsky and Yudin, 1983; Agarwal et al., 2012)
  - Strongly convex:  $O((\mu n)^{-1})$

Attained by averaged stochastic gradient descent with  $\gamma_n \propto (\mu n)^{-1}$ 

– Non-strongly convex:  $O(n^{-1/2})$ 

Attained by averaged stochastic gradient descent with  $\gamma_n \propto n^{-1/2}$ 

- Asymptotic analysis of averaging (Polyak and Juditsky, 1992; Ruppert, 1988)
  - All step sizes  $\gamma_n = Cn^{-\alpha}$  with  $\alpha \in (1/2, 1)$  lead to  $O(n^{-1})$  for smooth strongly convex problems

- Known global minimax rates of convergence for non-smooth problems (Nemirovsky and Yudin, 1983; Agarwal et al., 2012)
  - Strongly convex:  $O((\mu n)^{-1})$

Attained by averaged stochastic gradient descent with  $\gamma_n \propto (\mu n)^{-1}$ 

- Non-strongly convex:  $O(n^{-1/2})$ Attained by averaged stochastic gradient descent with  $\gamma_n \propto n^{-1/2}$
- Asymptotic analysis of averaging (Polyak and Juditsky, 1992; Ruppert, 1988)
  - All step sizes  $\gamma_n = Cn^{-\alpha}$  with  $\alpha \in (1/2, 1)$  lead to  $O(n^{-1})$  for smooth strongly convex problems
- Non-asymptotic analysis for smooth problems?

## **Smoothness/convexity assumptions**

• Iteration: 
$$\theta_n = \theta_{n-1} - \gamma_n f'_n(\theta_{n-1})$$

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

- Smoothness of  $f_n$ : For each  $n \ge 1$ , the function  $f_n$  is a.s. convex, differentiable with *L*-Lipschitz-continuous gradient  $f'_n$ :
  - Smooth loss and bounded data
- Strong convexity of f: The function f is strongly convex with respect to the norm  $\|\cdot\|$ , with convexity constant  $\mu > 0$ :
  - Invertible population covariance matrix
  - or regularization by  $\frac{\mu}{2} \|\theta\|^2$

# Summary of new results (Bach and Moulines, 2011)

- Stochastic gradient descent with learning rate  $\gamma_n = C n^{-\alpha}$
- Strongly convex smooth objective functions
  - 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
  - Forgetting of initial conditions
  - Robustness to the choice of  ${\boldsymbol C}$

## Summary of new results (Bach and Moulines, 2011)

- Stochastic gradient descent with learning rate  $\gamma_n = C n^{-\alpha}$
- Strongly convex smooth objective functions
  - 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
  - Forgetting of initial conditions
  - Robustness to the choice of  ${\boldsymbol C}$
- Convergence rates for  $\mathbb{E} \| \theta_n \theta^* \|^2$  and  $\mathbb{E} \| \overline{\theta}_n \theta^* \|^2$

- no averaging:  $O\left(\frac{\sigma^2 \gamma_n}{\mu}\right) + O(e^{-\mu n \gamma_n}) \|\theta_0 - \theta^*\|^2$ - averaging:  $\frac{\operatorname{tr} H(\theta^*)^{-1}}{n} + \mu^{-1}O(n^{-2\alpha} + n^{-2+\alpha}) + O\left(\frac{\|\theta_0 - \theta^*\|^2}{\mu^2 n^2}\right)$ 

## **Classical proof sketch (no averaging)**

$$\begin{split} \|\theta_n - \theta_*\|_2^2 &= \|\theta_{n-1} - \gamma_n f'_n(\theta_{n-1}) - \theta_*\|_2^2 \\ &= \|\theta_{n-1} - \theta_*\|_2^2 - 2\gamma_n(\theta_{n-1} - \theta_*)^\top f'_n(\theta_{n-1}) + \gamma_n^2 \|f'_n(\theta_{n-1})\|_2^2 \\ &\leqslant \|\theta_{n-1} - \theta_*\|_2^2 - 2\gamma_n(\theta_{n-1} - \theta_*)^\top f'_n(\theta_{n-1}) \\ &+ 2\gamma_n^2 \|f'_n(\theta_*)\|_2^2 + 2\gamma_n^2 \|f'_n(\theta_{n-1}) - f'_n(\theta_*)\|_2^2 \\ &\leqslant \|\theta_{n-1} - \theta_*\|_2^2 - 2\gamma_n(\theta_{n-1} - \theta_*)^\top f'_n(\theta_{n-1}) \\ &+ 2\gamma_n^2 \|f'_n(\theta_*)\|_2^2 + 2\gamma_n^2 L[f'_n(\theta_{n-1}) - f'_n(\theta_*)]^\top (\theta_{n-1} - \theta_*) \\ &\mathbb{E}[\|\theta_n - \theta_*\|_2^2]\mathcal{F}_{n-1}] &\leqslant \|\theta_{n-1} - \theta_*\|_2^2 - 2\gamma_n(\theta_{n-1} - \theta_*)^\top f'(\theta_{n-1}) \\ &+ 2\gamma_n^2 \mathbb{E}\|f'_n(\theta_*)\|_2^2 + 2\gamma_n^2 L[f'(\theta_{n-1}) - 0]^\top (\theta_{n-1} - \theta_*) \\ &\leqslant \|\theta_{n-1} - \theta_*\|_2^2 - 2\gamma_n(1 - \gamma_n L)(\theta_{n-1} - \theta_*)^\top f'(\theta_{n-1}) + 2\gamma_n^2 \sigma^2 \\ &\leqslant \|\theta_{n-1} - \theta_*\|_2^2 - 2\gamma_n(1 - \gamma_n L)\frac{1}{2}\mu\|\theta_{n-1} - \theta_*\|_2^2 + 2\gamma_n^2 \sigma^2 \\ &= [1 - \mu\gamma_n(1 - \gamma_n L)]\|\theta_{n-1} - \theta_*\|_2^2 + 2\gamma_n^2 \sigma^2 \\ \mathbb{E}[\|\theta_{n-1} - \theta_*\|_2^2] &\leqslant [1 - \mu\gamma_n(1 - \gamma_n L)]\mathbb{E}[\|\theta_{n-1} - \theta_*\|_2^2] + 2\gamma_n^2 \sigma^2 \end{split}$$

## **Proof sketch (averaging)**

• From Polyak and Juditsky (1992):

$$\begin{aligned} \theta_{n} &= \theta_{n-1} - \gamma_{n} f_{n}'(\theta_{n-1}) \\ \Leftrightarrow & f_{n}'(\theta_{n-1}) = \frac{1}{\gamma_{n}}(\theta_{n-1} - \theta_{n}) \\ \Leftrightarrow & f_{n}'(\theta_{*}) + f_{n}''(\theta_{*})(\theta_{n-1} - \theta_{*}) = \frac{1}{\gamma_{n}}(\theta_{n-1} - \theta_{n}) + O(||\theta_{n-1} - \theta_{*}||^{2}) \\ \Leftrightarrow & f_{n}'(\theta_{*}) + f''(\theta_{*})(\theta_{n-1} - \theta_{*}) = \frac{1}{\gamma_{n}}(\theta_{n-1} - \theta_{n}) + O(||\theta_{n-1} - \theta_{*}||^{2}) \\ & + O(||\theta_{n-1} - \theta_{*}||)\varepsilon_{n} \\ \Leftrightarrow & \theta_{n-1} - \theta_{*} = -f''(\theta_{*})^{-1}f_{n}'(\theta_{*}) + \frac{1}{\gamma_{n}}f''(\theta_{*})^{-1}(\theta_{n-1} - \theta_{n}) \\ & + O(||\theta_{n-1} - \theta_{*}||^{2}) + O(||\theta_{n-1} - \theta_{*}||)\varepsilon_{n} \end{aligned}$$

• Averaging to cancel the term  $\frac{1}{\gamma_n}f''(\theta_*)^{-1}(\theta_{n-1}-\theta_n)$ 

#### Robustness to wrong constants for $\gamma_n = C n^{-\alpha}$

- $f(\theta) = \frac{1}{2} |\theta|^2$  with i.i.d. Gaussian noise (d = 1)
- Left:  $\alpha = 1/2$
- Right:  $\alpha = 1$



• See also http://leon.bottou.org/projects/sgd

# Summary of new results (Bach and Moulines, 2011)

- Stochastic gradient descent with learning rate  $\gamma_n = C n^{-\alpha}$
- Strongly convex smooth objective functions
  - 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

# Summary of new results (Bach and Moulines, 2011)

- Stochastic gradient descent with learning rate  $\gamma_n = C n^{-\alpha}$
- Strongly convex smooth objective functions
  - 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

### Non-strongly convex smooth objective functions

- Old:  $O(n^{-1/2})$  rate achieved with averaging for  $\alpha = 1/2$
- New:  $O(\max\{n^{1/2-3\alpha/2}, n^{-\alpha/2}, n^{\alpha-1}\})$  rate achieved without averaging for  $\alpha \in [1/3, 1]$

#### • Take-home message

– Use  $\alpha=1/2$  with averaging to be adaptive to strong convexity

#### **Robustness to lack of strong convexity**

- Left:  $f(\theta) = |\theta|^2$  between -1 and 1
- Right:  $f(\theta) = |\theta|^4$  between -1 and 1
- affine outside of [-1, 1], continuously differentiable.



## **Beyond stochastic gradient method**

#### • Adding a proximal step

- Goal:  $\min_{\theta \in \mathbb{R}^d} f(\theta) + \Omega(\theta) = \mathbb{E} f_n(\theta) + \Omega(\theta)$ 

– Replace recursion  $\theta_n = \theta_{n-1} - \gamma_n f'_n(\theta_n)$  by

$$\theta_n = \min_{\theta \in \mathbb{R}^d} \left\| \theta - \theta_{n-1} + \gamma_n f'_n(\theta_n) \right\|_2^2 + C\Omega(\theta)$$

- Xiao (2010); Hu et al. (2009)
- May be accelerated (Ghadimi and Lan, 2013)
- Related frameworks
  - Regularized dual averaging (Nesterov, 2009; Xiao, 2010)
  - Mirror descent (Nemirovski et al., 2009; Lan et al., 2012)

# Outline

#### 1. Large-scale machine learning and optimization

- Traditional statistical analysis
- Classical methods for convex optimization

## 2. Non-smooth stochastic approximation

- Stochastic (sub)gradient and averaging
- Non-asymptotic results and lower bounds
- Strongly convex vs. non-strongly convex

### 3. Smooth stochastic approximation algorithms

- Asymptotic and non-asymptotic results
- 4. Beyond decaying step-sizes
- 5. Finite data sets

- Known global minimax rates of convergence for non-smooth problems (Nemirovsky and Yudin, 1983; Agarwal et al., 2012)
  - Strongly convex:  $O((\mu n)^{-1})$

Attained by averaged stochastic gradient descent with  $\gamma_n \propto (\mu n)^{-1}$ 

- Non-strongly convex:  $O(n^{-1/2})$ Attained by averaged stochastic gradient descent with  $\gamma_n \propto n^{-1/2}$
- Asymptotic analysis of averaging (Polyak and Juditsky, 1992; Ruppert, 1988)
  - All step sizes  $\gamma_n = Cn^{-\alpha}$  with  $\alpha \in (1/2, 1)$  lead to  $O(n^{-1})$  for smooth strongly convex problems
- A single adaptive algorithm for smooth problems with convergence rate  $O(\min\{1/\mu n, 1/\sqrt{n}\})$  in all situations?

- Logistic regression:  $(\Phi(x_n), y_n) \in \mathbb{R}^d \times \{-1, 1\}$ 
  - Single data point:  $f_n(\theta) = \log(1 + \exp(-y_n \theta^\top \Phi(x_n)))$
  - Generalization error:  $f(\theta) = \mathbb{E}f_n(\theta)$

- Logistic regression:  $(\Phi(x_n), y_n) \in \mathbb{R}^d \times \{-1, 1\}$ 
  - Single data point:  $f_n(\theta) = \log(1 + \exp(-y_n \theta^\top \Phi(x_n)))$
  - Generalization error:  $f(\theta) = \mathbb{E}f_n(\theta)$
- Cannot be strongly convex  $\Rightarrow$  local strong convexity
  - unless restricted to  $|\theta^{\top}\Phi(x_n)| \leq M$  (and with constants  $e^M$ )
  - $\mu =$  lowest eigenvalue of the Hessian at the optimum  $f''( heta_*)$



- Logistic regression:  $(\Phi(x_n), y_n) \in \mathbb{R}^d \times \{-1, 1\}$ 
  - Single data point:  $f_n(\theta) = \log(1 + \exp(-y_n \theta^\top \Phi(x_n)))$
  - Generalization error:  $f(\theta) = \mathbb{E}f_n(\theta)$
- Cannot be strongly convex  $\Rightarrow$  local strong convexity
  - unless restricted to  $|\theta^{\top}\Phi(x_n)| \leq M$  (and with constants  $e^M$ )
  - $\mu$  = lowest eigenvalue of the Hessian at the optimum  $f''(\theta_*)$
- *n* steps of averaged SGD with constant step-size  $1/(2R^2\sqrt{n})$ - with R = radius of data (Bach, 2013):

$$\mathbb{E}f(\bar{\theta}_n) - f(\theta_*) \leqslant \min\left\{\frac{1}{\sqrt{n}}, \frac{R^2}{n\mu}\right\} \left(15 + 5R\|\theta_0 - \theta_*\|\right)^4$$

- Proof based on self-concordance (Nesterov and Nemirovski, 1994)

## **Self-concordance**

- Usual definition for convex  $\varphi : \mathbb{R} \to \mathbb{R}$ :  $|\varphi'''(t)| \leq 2\varphi''(t)^{3/2}$ 
  - Affine invariant
  - Extendable to all convex functions on  $\mathbb{R}^d$  by looking at rays
  - Used for the sharp proof of quadratic convergence of Newton method (Nesterov and Nemirovski, 1994)
- Generalized notion:  $|\varphi^{\prime\prime\prime}(t)| \leqslant \varphi^{\prime\prime}(t)$ 
  - Applicable to logistic regression (with extensions) -  $\varphi(t) = \log(1 + e^{-t})$

## **Self-concordance**

- Usual definition for convex  $\varphi : \mathbb{R} \to \mathbb{R}$ :  $|\varphi'''(t)| \leq 2\varphi''(t)^{3/2}$ 
  - Affine invariant
  - Extendable to all convex functions on  $\mathbb{R}^d$  by looking at rays
  - Used for the sharp proof of quadratic convergence of Newton method (Nesterov and Nemirovski, 1994)
- Generalized notion:  $|\varphi^{\prime\prime\prime}(t)| \leqslant \varphi^{\prime\prime}(t)$ 
  - Applicable to logistic regression (with extensions)
- Important properties
  - Allows global Taylor expansions
  - Relates expansions of derivatives of different orders

## Adaptive algorithm for logistic regression Proof sketch

- Step 1: use existing result  $f(\bar{\theta}_n) f(\theta_*) + \frac{R^2}{\sqrt{n}} \|\theta_0 \theta_*\|_2^2 = O(1/\sqrt{n})$
- Step 2:  $f'_n(\theta_{n-1}) = \frac{1}{\gamma}(\theta_{n-1} \theta_n) \Rightarrow \frac{1}{n} \sum_{k=1}^n f'_k(\theta_{k-1}) = \frac{1}{n\gamma}(\theta_0 \theta_n)$

• Step 3: 
$$\left\| f'\left(\frac{1}{n} \sum_{k=1}^{n} \theta_{k-1}\right) - \frac{1}{n} \sum_{k=1}^{n} f'(\theta_{k-1}) \right\|_{2}$$
$$= O\left(f(\bar{\theta}_{n}) - f(\theta_{*})\right) = O(1/\sqrt{n}) \text{ using self-concordance}$$

- Step 4a: if  $f \mu$ -strongly convex,  $f(\bar{\theta}_n) f(\theta_*) \leq \frac{1}{2\mu} \|f'(\bar{\theta}_n)\|_2^2$
- Step 4b: if f self-concordant, "locally true" with  $\mu = \lambda_{\min}(f''(\theta_*))$

- Logistic regression:  $(\Phi(x_n), y_n) \in \mathbb{R}^d \times \{-1, 1\}$ 
  - Single data point:  $f_n(\theta) = \log(1 + \exp(-y_n \theta^\top \Phi(x_n)))$
  - Generalization error:  $f(\theta) = \mathbb{E}f_n(\theta)$
- Cannot be strongly convex  $\Rightarrow$  local strong convexity
  - unless restricted to  $|\theta^{\top}\Phi(x_n)| \leq M$  (and with constants  $e^M$ )
  - $\mu$  = lowest eigenvalue of the Hessian at the optimum  $f''(\theta_*)$
- *n* steps of averaged SGD with constant step-size  $1/(2R^2\sqrt{n})$ - with R = radius of data (Bach, 2013):

$$\mathbb{E}f(\bar{\theta}_n) - f(\theta_*) \leqslant \min\left\{\frac{1}{\sqrt{n}}, \frac{R^2}{n\mu}\right\} \left(15 + 5R\|\theta_0 - \theta_*\|\right)^4$$

- Proof based on self-concordance (Nesterov and Nemirovski, 1994)
## Adaptive algorithm for logistic regression

- Logistic regression:  $(\Phi(x_n), y_n) \in \mathbb{R}^d \times \{-1, 1\}$ 
  - Single data point:  $f_n(\theta) = \log(1 + \exp(-y_n \theta^\top \Phi(x_n)))$
  - Generalization error:  $f(\theta) = \mathbb{E}f_n(\theta)$
- Cannot be strongly convex  $\Rightarrow$  local strong convexity
  - unless restricted to  $|\theta^{\top}\Phi(x_n)| \leq M$  (and with constants  $e^M$ )
  - $\mu =$  lowest eigenvalue of the Hessian at the optimum  $f''(\theta_*)$
- *n* steps of averaged SGD with constant step-size  $1/(2R^2\sqrt{n})$ - with R = radius of data (Bach, 2013):

$$\mathbb{E}f(\bar{\theta}_n) - f(\theta_*) \leqslant \min\left\{\frac{1}{\sqrt{n}}, \frac{R^2}{n\mu}\right\} \left(15 + 5R\|\theta_0 - \theta_*\|\right)^4$$

- A single adaptive algorithm for smooth problems with convergence rate O(1/n) in all situations?

### Least-mean-square algorithm

- Least-squares:  $f(\theta) = \frac{1}{2}\mathbb{E}[(y_n \langle \Phi(x_n), \theta \rangle)^2]$  with  $\theta \in \mathbb{R}^d$ 
  - SGD = least-mean-square algorithm (see, e.g., Macchi, 1995)
  - usually studied without averaging and decreasing step-sizes
  - with strong convexity assumption  $\mathbb{E}[\Phi(x_n) \otimes \Phi(x_n)] = H \succcurlyeq \mu \cdot \mathrm{Id}$

### Least-mean-square algorithm

- Least-squares:  $f(\theta) = \frac{1}{2}\mathbb{E}[(y_n \langle \Phi(x_n), \theta \rangle)^2]$  with  $\theta \in \mathbb{R}^d$ 
  - SGD = least-mean-square algorithm (see, e.g., Macchi, 1995)
  - usually studied without averaging and decreasing step-sizes
  - with strong convexity assumption  $\mathbb{E}[\Phi(x_n) \otimes \Phi(x_n)] = H \succcurlyeq \mu \cdot \mathrm{Id}$
- $\bullet$  New analysis for averaging and constant step-size  $\gamma = 1/(4R^2)$ 
  - Assume  $\|\Phi(x_n)\| \leq R$  and  $|y_n \langle \Phi(x_n), \theta_* \rangle| \leq \sigma$  almost surely

- No assumption regarding lowest eigenvalues of H

- Main result: 
$$\left| \mathbb{E}f(\bar{\theta}_{n-1}) - f(\theta_*) \leqslant \frac{4\sigma^2 d}{n} + \frac{4R^2 \|\theta_0 - \theta_*\|^2}{n} \right|$$

- Matches statistical lower bound (Tsybakov, 2003)
  - Non-asymptotic robust version of Györfi and Walk (1996)

### **Least-squares - Proof technique**

• LMS recursion:

$$\theta_n - \theta_* = \left[I - \gamma \Phi(x_n) \otimes \Phi(x_n)\right] (\theta_{n-1} - \theta_*) + \gamma \varepsilon_n \Phi(x_n)$$

• Simplified LMS recursion: with  $H = \mathbb{E}[\Phi(x_n) \otimes \Phi(x_n)]$ 

$$\theta_n - \theta_* = \left[I - \gamma \mathbf{H}\right](\theta_{n-1} - \theta_*) + \gamma \varepsilon_n \Phi(x_n)$$

- Direct proof technique of Polyak and Juditsky (1992), e.g.,

$$\theta_n - \theta_* = \left[I - \gamma \mathbf{H}\right]^n (\theta_0 - \theta_*) + \gamma \sum_{k=1}^n \left[I - \gamma \mathbf{H}\right]^{n-k} \varepsilon_k \Phi(x_k)$$

- Infinite expansion of Aguech, Moulines, and Priouret (2000) in powers of  $\gamma$ 

$$\theta_n = \theta_{n-1} - \gamma \big( \langle \Phi(x_n), \theta_{n-1} \rangle - y_n \big) \Phi(x_n)$$

- The sequence  $(\theta_n)_n$  is a homogeneous Markov chain
  - convergence to a stationary distribution  $\pi_{\gamma}$
  - with expectation  $\bar{\theta}_{\gamma} \stackrel{\text{def}}{=} \int \theta \pi_{\gamma}(\mathrm{d}\theta)$



$$\theta_n = \theta_{n-1} - \gamma \big( \langle \Phi(x_n), \theta_{n-1} \rangle - y_n \big) \Phi(x_n)$$

- The sequence  $(\theta_n)_n$  is a homogeneous Markov chain
  - convergence to a stationary distribution  $\pi_{\gamma}$
  - with expectation  $\bar{\theta}_{\gamma} \stackrel{\text{def}}{=} \int \theta \pi_{\gamma}(\mathrm{d}\theta)$
- For least-squares,  $\bar{\theta}_{\gamma} = \theta_*$



$$\theta_n = \theta_{n-1} - \gamma \big( \langle \Phi(x_n), \theta_{n-1} \rangle - y_n \big) \Phi(x_n)$$

- The sequence  $(\theta_n)_n$  is a homogeneous Markov chain
  - convergence to a stationary distribution  $\pi_{\gamma}$
  - with expectation  $\bar{\theta}_{\gamma} \stackrel{\mathrm{def}}{=} \int \theta \pi_{\gamma}(\mathrm{d}\theta)$
- For least-squares,  $\bar{\theta}_{\gamma} = \theta_{*}$



$$\theta_n = \theta_{n-1} - \gamma \big( \langle \Phi(x_n), \theta_{n-1} \rangle - y_n \big) \Phi(x_n)$$

- The sequence  $(\theta_n)_n$  is a homogeneous Markov chain
  - convergence to a stationary distribution  $\pi_\gamma$
  - with expectation  $\bar{\theta}_{\gamma} \stackrel{\mathrm{def}}{=} \int \theta \pi_{\gamma}(\mathrm{d}\theta)$
- For least-squares,  $\bar{\theta}_{\gamma} = \theta_{*}$ 
  - $\theta_n$  does not converge to  $\theta_*$  but oscillates around it
  - oscillations of order  $\sqrt{\gamma}$
- Ergodic theorem:
  - Averaged iterates converge to  $\bar{ heta}_{\gamma} = heta_*$  at rate O(1/n)

### **Simulations - synthetic examples**

 $\bullet$  Gaussian distributions - p=20



### **Simulations - benchmarks**

• alpha (p = 500,  $n = 500\ 000$ ), news ( $p = 1\ 300\ 000$ ,  $n = 20\ 000$ )



# **Optimal bounds for least-squares?**

- Least-squares: cannot beat  $\sigma^2 p/n$  (Tsybakov, 2003). Really? - What if  $p \gg n$ ?
- **Refined assumptions with adaptivity** (Dieuleveut and Bach, 2014)
  - Beyond strong convexity or lack thereof

## Finer assumptions (Dieuleveut and Bach, 2014)

#### • Covariance eigenvalues

- Pessimistic assumption: all eigenvalues  $\lambda_m$  less than a constant
- Actual decay as  $\lambda_m = o(m^{-\alpha})$  with tr  $H^{1/\alpha} = \sum_m \lambda_m^{1/\alpha}$  small



# Finer assumptions (Dieuleveut and Bach, 2014)

#### • Covariance eigenvalues

– Pessimistic assumption: all eigenvalues  $\lambda_m$  less than a constant

- Actual decay as  $\lambda_m = o(m^{-\alpha})$  with  $\operatorname{tr} H^{1/\alpha} = \sum_m \lambda_m^{1/\alpha}$  small  $\sigma^2 p = \sigma^2 (\gamma n)^{1/\alpha} \operatorname{tr} H^{1/\alpha}$ 

– New result: replace  $\frac{\sigma^2 p}{n}$  by  $\frac{\sigma^2 (\gamma n)^{1/\alpha} \operatorname{tr} H^{1/\alpha}}{n}$ 



# Finer assumptions (Dieuleveut and Bach, 2014)

#### • Covariance eigenvalues

- Pessimistic assumption: all eigenvalues  $\lambda_m$  less than a constant
- Actual decay as  $\lambda_m = o(m^{-\alpha})$  with  $\operatorname{tr} H^{1/\alpha} = \sum_m \lambda_m^{1/\alpha}$  small - New result: replace  $\frac{\sigma^2 p}{n}$  by  $\frac{\sigma^2 (\gamma n)^{1/\alpha} \operatorname{tr} H^{1/\alpha}}{n}$

#### • Optimal predictor

- Pessimistic assumption:  $\|\theta_0 \theta_*\|^2$  finite
- Finer assumption:  $||H^{1/2-r}(\theta_0 \theta_*)||_2$  small - Replace  $\frac{||\theta_0 - \theta_*||^2}{\gamma n}$  by  $\frac{4||H^{1/2-r}(\theta_0 - \theta_*)||_2}{\gamma^{2r}n^{2\min\{r,1\}}}$

# **Optimal bounds for least-squares?**

• Least-squares: cannot beat  $\sigma^2 p/n$  (Tsybakov, 2003). Really?

– What if  $p \gg n$ ?

- Refined assumptions with adaptivity (Dieuleveut and Bach, 2014)
  - Beyond strong convexity or lack thereof

$$f(\bar{\theta}_n) - f(\theta_*) \leqslant \frac{16\sigma^2 \operatorname{tr} H^{1/\alpha}}{n} (\gamma n)^{1/\alpha} + \frac{4\|H^{1/2 - r}(\theta_0 - \theta_*)\|_2}{\gamma^{2r} n^{2\min\{r, 1\}}}$$

- Previous results:  $\alpha=+\infty$  and r=1/2
- Valid for all  $\alpha$  and r
- Optimal step-size potentially decaying with n
- Extension to non-parametric estimation (kernels) with optimal rates

# Bias-variance decomposition (Défossez and Bach, 2015)

- $\bullet$  Simplification: dominating term when  $n \to \infty$  and  $\gamma \to 0$
- Variance (e.g., starting from the solution)

$$f(\bar{\theta}_n) - f(\theta_*) \sim \frac{1}{n} \mathbb{E} \Big[ \varepsilon^2 \Phi(x)^\top H^{-1} \Phi(x) \Big]$$

- NB: if noise  $\varepsilon$  is independent, then we obtain  $\frac{p\sigma^2}{n}$
- Exponentially decaying remainder terms (strongly convex problems)
- Bias (e.g., no noise)

$$f(\bar{\theta}_n) - f(\theta_*) \sim \frac{1}{n^2 \gamma^2} (\theta_0 - \theta_*)^\top H^{-1}(\theta_0 - \theta_*)$$

### **Bias-variance decomposition**



### **Bias-variance decomposition**



• Sampling from a different distribution with importance weights

$$\mathbb{E}_{\boldsymbol{p}(\boldsymbol{x})\boldsymbol{p}(\boldsymbol{y}|\boldsymbol{x})}|\boldsymbol{y} - \boldsymbol{\Phi}(\boldsymbol{x})^{\top}\boldsymbol{\theta}|^2 = \mathbb{E}_{\boldsymbol{q}(\boldsymbol{x})\boldsymbol{p}(\boldsymbol{y}|\boldsymbol{x})}\frac{d\boldsymbol{p}(\boldsymbol{x})}{d\boldsymbol{q}(\boldsymbol{x})}|\boldsymbol{y} - \boldsymbol{\Phi}(\boldsymbol{x})^{\top}\boldsymbol{\theta}|^2$$

• Sampling from a different distribution with importance weights

$$\mathbb{E}_{\boldsymbol{p}(\boldsymbol{x})\boldsymbol{p}(\boldsymbol{y}|\boldsymbol{x})}|\boldsymbol{y} - \boldsymbol{\Phi}(\boldsymbol{x})^{\top}\boldsymbol{\theta}|^2 = \mathbb{E}_{\boldsymbol{q}(\boldsymbol{x})\boldsymbol{p}(\boldsymbol{y}|\boldsymbol{x})}\frac{d\boldsymbol{p}(\boldsymbol{x})}{d\boldsymbol{q}(\boldsymbol{x})}|\boldsymbol{y} - \boldsymbol{\Phi}(\boldsymbol{x})^{\top}\boldsymbol{\theta}|^2$$

- Specific to least-squares =  $\mathbb{E}_{q(x)p(y|x)} \left| \sqrt{\frac{dp(x)}{dq(x)}} y \sqrt{\frac{dp(x)}{dq(x)}} \Phi(x)^{\top} \theta \right|^2$
- Reweighting of the data: same bounds apply!

• Sampling from a different distribution with importance weights

$$\mathbb{E}_{\boldsymbol{p}(\boldsymbol{x})\boldsymbol{p}(\boldsymbol{y}|\boldsymbol{x})}|\boldsymbol{y} - \boldsymbol{\Phi}(\boldsymbol{x})^{\top}\boldsymbol{\theta}|^2 = \mathbb{E}_{\boldsymbol{q}(\boldsymbol{x})\boldsymbol{p}(\boldsymbol{y}|\boldsymbol{x})}\frac{d\boldsymbol{p}(\boldsymbol{x})}{d\boldsymbol{q}(\boldsymbol{x})}|\boldsymbol{y} - \boldsymbol{\Phi}(\boldsymbol{x})^{\top}\boldsymbol{\theta}|^2$$

- Specific to least-squares =  $\mathbb{E}_{q(x)p(y|x)} \left| \sqrt{\frac{dp(x)}{dq(x)}} y \sqrt{\frac{dp(x)}{dq(x)}} \Phi(x)^{\top} \theta \right|^2$
- Reweighting of the data: same bounds apply!

• Optimal for variance: 
$$\frac{dq(x)}{dp(x)} \propto \sqrt{\Phi(x)^{\top} H^{-1} \Phi(x)}$$

- Same density as active learning (Kanamori and Shimodaira, 2003)
- Limited gains: different between first and second moments
- Caveat: need to know H

• Sampling from a different distribution with importance weights

$$\mathbb{E}_{\boldsymbol{p}(\boldsymbol{x})\boldsymbol{p}(\boldsymbol{y}|\boldsymbol{x})}|\boldsymbol{y} - \boldsymbol{\Phi}(\boldsymbol{x})^{\top}\boldsymbol{\theta}|^2 = \mathbb{E}_{\boldsymbol{q}(\boldsymbol{x})\boldsymbol{p}(\boldsymbol{y}|\boldsymbol{x})}\frac{d\boldsymbol{p}(\boldsymbol{x})}{d\boldsymbol{q}(\boldsymbol{x})}|\boldsymbol{y} - \boldsymbol{\Phi}(\boldsymbol{x})^{\top}\boldsymbol{\theta}|^2$$

- Specific to least-squares =  $\mathbb{E}_{q(x)p(y|x)} \left| \sqrt{\frac{dp(x)}{dq(x)}} y \sqrt{\frac{dp(x)}{dq(x)}} \Phi(x)^{\top} \theta \right|^2$
- Reweighting of the data: same bounds apply!

• Optimal for bias: 
$$\frac{dq(x)}{dp(x)} \propto \|\Phi(x)\|^2$$

- Simpy allows biggest possible step size  $\gamma < \frac{2}{\operatorname{tr} H}$
- Large gains in practice
- Corresponds to normalized least-mean-squares

## **Convergence on** *Sido* dataset



## Acceleration (Flammarion and Bach, 2015)

• Existing results (Bach and Moulines, 2013)

- Variance = 
$$\frac{\sigma^2 p}{n}$$
  
- Bias  $\leq \min\left\{\frac{R^2 \|\theta_0 - \theta_*\|^2}{n}, \frac{R^4 \langle \theta_0 - \theta_*, H^{-1}(\theta_0 - \theta_*) \rangle}{n^2}\right\}$ 

- Is it possible to get a bias term in  $\frac{R^2 \|\theta_0 \theta_*\|^2}{n^2}$ ?
  - Corresponds to acceleration (Nesterov, 1983)
  - Best (current) result:

$$\frac{\sigma^2 p}{n^{1-\alpha}} + \frac{R^2 \|\theta_0 - \theta_*\|^2}{n^{1+\alpha}}$$

# **Beyond least-squares - Markov chain interpretation**

- Recursion  $\theta_n = \theta_{n-1} \gamma f'_n(\theta_{n-1})$  also defines a Markov chain
  - Stationary distribution  $\pi_{\gamma}$  such that  $\int f'(\theta) \pi_{\gamma}(\mathrm{d}\theta) = 0$
  - When f' is not linear,  $f'(\int \theta \pi_{\gamma}(\mathrm{d}\theta)) \neq \int f'(\theta)\pi_{\gamma}(\mathrm{d}\theta) = 0$

### **Beyond least-squares - Markov chain interpretation**

- Recursion  $\theta_n = \theta_{n-1} \gamma f'_n(\theta_{n-1})$  also defines a Markov chain
  - Stationary distribution  $\pi_{\gamma}$  such that  $\int f'(\theta)\pi_{\gamma}(d\theta) = 0$ - When f' is not linear  $f'(\int \theta \pi_{\gamma}(d\theta)) \neq \int f'(\theta)\pi_{\gamma}(d\theta) = 0$
  - When f' is not linear,  $f'(\int \theta \pi_{\gamma}(\mathrm{d}\theta)) \neq \int f'(\theta) \pi_{\gamma}(\mathrm{d}\theta) = 0$
- $\theta_n$  oscillates around the wrong value  $\bar{\theta}_{\gamma} \neq \theta_*$



# **Beyond least-squares - Markov chain interpretation**

- Recursion  $\theta_n = \theta_{n-1} \gamma f'_n(\theta_{n-1})$  also defines a Markov chain
  - Stationary distribution  $\pi_{\gamma}$  such that  $\int f'(\theta)\pi_{\gamma}(d\theta) = 0$
  - When f' is not linear,  $f'(\int \theta \pi_{\gamma}(\mathrm{d}\theta)) \neq \int f'(\theta) \pi_{\gamma}(\mathrm{d}\theta) = 0$
- $\theta_n$  oscillates around the wrong value  $\bar{\theta}_{\gamma} \neq \theta_*$

– moreover, 
$$\| heta_* - heta_n\| = O_p(\sqrt{\gamma})$$

 Linear convergence up to the noise level for strongly-convex problems (Nedic and Bertsekas, 2000)

### • Ergodic theorem

- averaged iterates converge to  $\bar{\theta}_{\gamma} \neq \theta_*$  at rate O(1/n)
- moreover,  $\|\theta_* \bar{\theta}_{\gamma}\| = O(\gamma)$  (Bach, 2013)

### **Simulations - synthetic examples**

 $\bullet$  Gaussian distributions - p=20



## • Known facts

- 1. Averaged SGD with  $\gamma_n \propto n^{-1/2}$  leads to *robust* rate  $O(n^{-1/2})$  for all convex functions
- 2. Averaged SGD with  $\gamma_n$  constant leads to *robust* rate  $O(n^{-1})$  for all convex *quadratic* functions
- 3. Newton's method squares the error at each iteration for smooth functions
- 4. A single step of Newton's method is equivalent to minimizing the quadratic Taylor expansion

## • Known facts

- 1. Averaged SGD with  $\gamma_n \propto n^{-1/2}$  leads to *robust* rate  $O(n^{-1/2})$  for all convex functions
- 2. Averaged SGD with  $\gamma_n$  constant leads to *robust* rate  $O(n^{-1})$  for all convex *quadratic* functions  $\Rightarrow O(n^{-1})$
- 3. Newton's method squares the error at each iteration for smooth functions  $\Rightarrow O((n^{-1/2})^2)$
- 4. A single step of Newton's method is equivalent to minimizing the quadratic Taylor expansion

# • Online Newton step

- Rate:  $O((n^{-1/2})^2 + n^{-1}) = O(n^{-1})$
- Complexity: O(p) per iteration

• The Newton step for  $f = \mathbb{E}f_n(\theta) \stackrel{\text{def}}{=} \mathbb{E}[\ell(y_n, \langle \theta, \Phi(x_n) \rangle)]$  at  $\tilde{\theta}$  is equivalent to minimizing the quadratic approximation

$$g(\theta) = f(\tilde{\theta}) + \langle f'(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \frac{1}{2} \langle \theta - \tilde{\theta}, f''(\tilde{\theta})(\theta - \tilde{\theta}) \rangle$$
  
$$= f(\tilde{\theta}) + \langle \mathbb{E}f'_{n}(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \frac{1}{2} \langle \theta - \tilde{\theta}, \mathbb{E}f''_{n}(\tilde{\theta})(\theta - \tilde{\theta}) \rangle$$
  
$$= \mathbb{E}\Big[f(\tilde{\theta}) + \langle f'_{n}(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \frac{1}{2} \langle \theta - \tilde{\theta}, f''_{n}(\tilde{\theta})(\theta - \tilde{\theta}) \rangle\Big]$$

• The Newton step for  $f = \mathbb{E}f_n(\theta) \stackrel{\text{def}}{=} \mathbb{E}[\ell(y_n, \langle \theta, \Phi(x_n) \rangle)]$  at  $\tilde{\theta}$  is equivalent to minimizing the quadratic approximation

$$g(\theta) = f(\tilde{\theta}) + \langle f'(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \frac{1}{2} \langle \theta - \tilde{\theta}, f''(\tilde{\theta})(\theta - \tilde{\theta}) \rangle$$
  
$$= f(\tilde{\theta}) + \langle \mathbb{E}f'_{n}(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \frac{1}{2} \langle \theta - \tilde{\theta}, \mathbb{E}f''_{n}(\tilde{\theta})(\theta - \tilde{\theta}) \rangle$$
  
$$= \mathbb{E}\Big[f(\tilde{\theta}) + \langle f'_{n}(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \frac{1}{2} \langle \theta - \tilde{\theta}, f''_{n}(\tilde{\theta})(\theta - \tilde{\theta}) \rangle\Big]$$

• Complexity of least-mean-square recursion for g is O(p)

$$\theta_n = \theta_{n-1} - \gamma \left[ f'_n(\tilde{\theta}) + f''_n(\tilde{\theta})(\theta_{n-1} - \tilde{\theta}) \right]$$

 $-f_n''(\tilde{\theta}) = \ell''(y_n, \langle \tilde{\theta}, \Phi(x_n) \rangle) \Phi(x_n) \otimes \Phi(x_n)$  has rank one

- New online Newton step without computing/inverting Hessians

# **Choice of support point for online Newton step**

#### • Two-stage procedure

- (1) Run n/2 iterations of averaged SGD to obtain  $ilde{ heta}$
- (2) Run n/2 iterations of averaged constant step-size LMS
  - Reminiscent of one-step estimators (see, e.g., Van der Vaart, 2000)
  - Provable convergence rate of O(p/n) for logistic regression
  - Additional assumptions but no strong convexity

# **Logistic regression - Proof technique**

• Using generalized self-concordance of  $\varphi: u \mapsto \log(1 + e^{-u})$ :

 $|\varphi'''(u)| \leqslant \varphi''(u)$ 

- NB: difference with regular self-concordance:  $|\varphi'''(u)| \leq 2\varphi''(u)^{3/2}$
- Using novel high-probability convergence results for regular averaged stochastic gradient descent
- Requires assumption on the kurtosis in every direction, i.e.,

$$\mathbb{E}\langle \Phi(x_n), \eta \rangle^4 \leqslant \kappa \big[ \mathbb{E}\langle \Phi(x_n), \eta \rangle^2 \big]^2$$

# **Choice of support point for online Newton step**

#### • Two-stage procedure

- (1) Run n/2 iterations of averaged SGD to obtain  $ilde{ heta}$
- (2) Run n/2 iterations of averaged constant step-size LMS
  - Reminiscent of one-step estimators (see, e.g., Van der Vaart, 2000)
  - Provable convergence rate of O(p/n) for logistic regression
  - Additional assumptions but no strong convexity
- Update at each iteration using the current averaged iterate

- Recursion: 
$$\theta_n = \theta_{n-1} - \gamma \left[ f'_n(\bar{\theta}_{n-1}) + f''_n(\bar{\theta}_{n-1})(\theta_{n-1} - \bar{\theta}_{n-1}) \right]$$

- No provable convergence rate (yet) but best practical behavior
- Note (dis)similarity with regular SGD:  $\theta_n = \theta_{n-1} \gamma f'_n(\theta_{n-1})$

# Online Newton algorithm Current proof (Flammarion et al., 2014)

• Recursion

$$\begin{cases} \theta_n = \theta_{n-1} - \gamma \left[ f'_n(\bar{\theta}_{n-1}) + f''_n(\bar{\theta}_{n-1})(\theta_{n-1} - \bar{\theta}_{n-1}) \right] \\ \bar{\theta}_n = \bar{\theta}_{n-1} + \frac{1}{n}(\theta_n - \bar{\theta}_{n-1}) \end{cases}$$

- Instance of two-time-scale stochastic approximation (Borkar, 1997)
  - Given  $\bar{\theta}$ ,  $\theta_n = \theta_{n-1} \gamma \left[ f'_n(\bar{\theta}) + f''_n(\bar{\theta})(\theta_{n-1} \bar{\theta}) \right]$  defines a homogeneous Markov chain (fast dynamics) -  $\bar{\theta}_n$  is updated at rate 1/n (slow dynamics)
- **Difficulty**: preserving robustness to ill-conditioning
#### **Simulations - synthetic examples**

 $\bullet$  Gaussian distributions - p=20



#### **Simulations - benchmarks**

• alpha (p = 500,  $n = 500\ 000$ ), news ( $p = 1\ 300\ 000$ ,  $n = 20\ 000$ )



# Outline

#### 1. Large-scale machine learning and optimization

- Traditional statistical analysis
- Classical methods for convex optimization

### 2. Non-smooth stochastic approximation

- Stochastic (sub)gradient and averaging
- Non-asymptotic results and lower bounds
- Strongly convex vs. non-strongly convex

### 3. Smooth stochastic approximation algorithms

- Asymptotic and non-asymptotic results
- 4. Beyond decaying step-sizes
- 5. Finite data sets

# Going beyond a single pass over the data

- Stochastic approximation
  - Assumes infinite data stream
  - Observations are used only once
  - Directly minimizes testing cost  $\mathbb{E}_{(x,y)} \ell(y, \theta^{\top} \Phi(x))$

# Going beyond a single pass over the data

#### • Stochastic approximation

- Assumes infinite data stream
- Observations are used only once
- Directly minimizes testing cost  $\mathbb{E}_{(x,y)} \ell(y, \theta^{\top} \Phi(x))$
- Machine learning practice
  - Finite data set  $(x_1, y_1, \ldots, x_n, y_n)$
  - Multiple passes
  - Minimizes training cost  $\frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \theta^{\top} \Phi(x_i))$
  - Need to regularize (e.g., by the  $\ell_2$ -norm) to avoid overfitting

• **Goal**: minimize 
$$g(\theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\theta)$$

• Minimizing  $g(\theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\theta)$  with  $f_i(\theta) = \ell(y_i, \theta^\top \Phi(x_i)) + \mu \Omega(\theta)$ 

- Batch gradient descent:  $\theta_t = \theta_{t-1} \gamma_t g'(\theta_{t-1}) = \theta_{t-1} \frac{\gamma_t}{n} \sum_{i=1}^n f'_i(\theta_{t-1})$ 
  - Linear (e.g., exponential) convergence rate in  $O(e^{-\alpha t})$
  - Iteration complexity is linear in *n* (with line search)

- Minimizing  $g(\theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\theta)$  with  $f_i(\theta) = \ell(y_i, \theta^\top \Phi(x_i)) + \mu \Omega(\theta)$
- Batch gradient descent:  $\theta_t = \theta_{t-1} \gamma_t g'(\theta_{t-1}) = \theta_{t-1} \frac{\gamma_t}{n} \sum_{i=1}^n f'_i(\theta_{t-1})$



• Minimizing  $g(\theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\theta)$  with  $f_i(\theta) = \ell(y_i, \theta^\top \Phi(x_i)) + \mu \Omega(\theta)$ 

- Batch gradient descent:  $\theta_t = \theta_{t-1} \gamma_t g'(\theta_{t-1}) = \theta_{t-1} \frac{\gamma_t}{n} \sum_{i=1}^n f'_i(\theta_{t-1})$ 
  - Linear (e.g., exponential) convergence rate in  $O(e^{-\alpha t})$
  - Iteration complexity is linear in *n* (with line search)

- Stochastic gradient descent:  $\theta_t = \theta_{t-1} \gamma_t f'_{i(t)}(\theta_{t-1})$ 
  - Sampling with replacement: i(t) random element of  $\{1, \ldots, n\}$
  - Convergence rate in O(1/t)
  - Iteration complexity is independent of n (step size selection?)

- Minimizing  $g(\theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\theta)$  with  $f_i(\theta) = \ell(y_i, \theta^\top \Phi(x_i)) + \mu \Omega(\theta)$
- Batch gradient descent:  $\theta_t = \theta_{t-1} \gamma_t g'(\theta_{t-1}) = \theta_{t-1} \frac{\gamma_t}{n} \sum_{i=1}^n f'_i(\theta_{t-1})$



• Stochastic gradient descent:  $\theta_t = \theta_{t-1} - \gamma_t f'_{i(t)}(\theta_{t-1})$ 



• Goal = best of both worlds: Linear rate with O(1) iteration cost Robustness to step size



• Goal = best of both worlds: Linear rate with O(1) iteration cost Robustness to step size



# **Accelerating gradient methods - Related work**

- Nesterov acceleration
  - Nesterov (1983, 2004)
  - Better linear rate but still O(n) iteration cost
- Hybrid methods, incremental average gradient, increasing batch size
  - Bertsekas (1997); Blatt et al. (2008); Friedlander and Schmidt (2011)
  - Linear rate, but iterations make full passes through the data.

# **Accelerating gradient methods - Related work**

- Momentum, gradient/iterate averaging, stochastic version of accelerated batch gradient methods
  - Polyak and Juditsky (1992); Tseng (1998); Sunehag et al. (2009);
    Ghadimi and Lan (2010); Xiao (2010)
  - Can improve constants, but still have sublinear O(1/t) rate
- Constant step-size stochastic gradient (SG), accelerated SG
  - Kesten (1958); Delyon and Juditsky (1993); Solodov (1998); Nedic and Bertsekas (2000)
  - Linear convergence, but only up to a fixed tolerance.
- Stochastic methods in the dual
  - Shalev-Shwartz and Zhang (2012)
  - Similar linear rate but limited choice for the  $f_i$ 's

# Stochastic average gradient (Le Roux, Schmidt, and Bach, 2012)

- Stochastic average gradient (SAG) iteration
  - Keep in memory the gradients of all functions  $f_i$ ,  $i = 1, \ldots, n$
  - Random selection  $i(t) \in \{1, \ldots, n\}$  with replacement

- Iteration: 
$$\theta_t = \theta_{t-1} - \frac{\gamma_t}{n} \sum_{i=1}^n y_i^t$$
 with  $y_i^t = \begin{cases} f'_i(\theta_{t-1}) & \text{if } i = i(t) \\ y_i^{t-1} & \text{otherwise} \end{cases}$ 

# Stochastic average gradient (Le Roux, Schmidt, and Bach, 2012)

- Stochastic average gradient (SAG) iteration
  - Keep in memory the gradients of all functions  $f_i$ ,  $i = 1, \ldots, n$
  - Random selection  $i(t) \in \{1, \ldots, n\}$  with replacement

- Iteration: 
$$\theta_t = \theta_{t-1} - \frac{\gamma_t}{n} \sum_{i=1}^n y_i^t$$
 with  $y_i^t = \begin{cases} f'_i(\theta_{t-1}) & \text{if } i = i(t) \\ y_i^{t-1} & \text{otherwise} \end{cases}$ 

- Stochastic version of incremental average gradient (Blatt et al., 2008)
- Extra memory requirement
  - Supervised machine learning
    - If  $f_i(\theta) = \ell_i(y_i, \Phi(x_i)^\top \theta)$ , then  $f'_i(\theta) = \ell'_i(y_i, \Phi(x_i)^\top \theta) \Phi(x_i)$
    - Only need to store  $\boldsymbol{n}$  real numbers

# **Stochastic average gradient - Convergence analysis**

#### • Assumptions

- Each  $f_i$  is L-smooth,  $i = 1, \ldots, n$
- $-g = \frac{1}{n} \sum_{i=1}^{n} f_i$  is  $\mu$ -strongly convex (with potentially  $\mu = 0$ )
- constant step size  $\gamma_t = 1/(16L)$
- initialization with one pass of averaged SGD

# **Stochastic average gradient - Convergence analysis**

#### • Assumptions

- Each  $f_i$  is L-smooth,  $i = 1, \ldots, n$
- $-g = \frac{1}{n} \sum_{i=1}^{n} f_i$  is  $\mu$ -strongly convex (with potentially  $\mu = 0$ )
- constant step size  $\gamma_t = 1/(16L)$
- initialization with one pass of averaged SGD
- Strongly convex case (Le Roux et al., 2012, 2013)

$$\mathbb{E}\big[g(\theta_t) - g(\theta_*)\big] \leqslant \Big(\frac{8\sigma^2}{n\mu} + \frac{4L\|\theta_0 - \theta_*\|^2}{n}\Big) \, \exp\Big(-t\min\Big\{\frac{1}{8n}, \frac{\mu}{16L}\Big\}\Big)$$

– Linear (exponential) convergence rate with O(1) iteration cost

- After one pass, reduction of cost by  $\exp\left(-\min\left\{\frac{1}{8},\frac{n\mu}{16L}\right\}\right)$ 

# **Stochastic average gradient - Convergence analysis**

#### • Assumptions

- Each  $f_i$  is L-smooth,  $i = 1, \ldots, n$
- $-g = \frac{1}{n} \sum_{i=1}^{n} f_i$  is  $\mu$ -strongly convex (with potentially  $\mu = 0$ )
- constant step size  $\gamma_t = 1/(16L)$
- initialization with one pass of averaged SGD
- Non-strongly convex case (Le Roux et al., 2013)

$$\mathbb{E}\left[g(\theta_t) - g(\theta_*)\right] \leqslant 48 \frac{\sigma^2 + L \|\theta_0 - \theta_*\|^2}{\sqrt{n}} \frac{n}{t}$$

- Improvement over regular batch and stochastic gradient
- Adaptivity to potentially hidden strong convexity

# **Convergence analysis - Proof sketch**

- Main step: find "good" Lyapunov function  $J(\theta_t, y_1^t, \dots, y_n^t)$ 
  - such that  $\mathbb{E}[J(\theta_t, y_1^t, \dots, y_n^t) | \mathcal{F}_{t-1}] < J(\theta_{t-1}, y_1^{t-1}, \dots, y_n^{t-1})$ - no natural candidates

### • Computer-aided proof

- Parameterize function  $J(\theta_t, y_1^t, \dots, y_n^t) = g(\theta_t) g(\theta_*) + quadratic$
- Solve semidefinite program to obtain candidates (that depend on  $n,\mu,L)$
- Check validity with symbolic computations

# Rate of convergence comparison

 $\bullet$  Assume that  $L=100\text{, }\mu=.01\text{, and }n=80000$ 

- Full gradient method has rate

$$\left(1 - \frac{\mu}{L}\right) = 0.9999$$

- Accelerated gradient method has rate

$$\left(1 - \sqrt{\frac{\mu}{L}}\right) = 0.9900$$

– Running  $\boldsymbol{n}$  iterations of SAG for the same cost has rate

$$\left(1 - \frac{1}{8n}\right)^n = 0.8825$$

- Fastest possible first-order method has rate

$$\left(\frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}\right)^2 = 0.9608$$

- Beating two lower bounds (with additional assumptions)
  - -(1) stochastic gradient and (2) full gradient

# **Stochastic average gradient Implementation details and extensions**

- The algorithm can use sparsity in the features to reduce the storage and iteration cost
- Grouping functions together can further reduce the memory requirement
- We have obtained good performance when L is not known with a heuristic line-search
- Algorithm allows non-uniform sampling
- Possibility of making proximal, coordinate-wise, and Newton-like variants





# **Extensions and related work**

- Exponential convergence rate for strongly convex problems
- Need to store gradients
  - SVRG (Johnson and Zhang, 2013)
- Adaptivity to non-strong convexity
  - SAGA (Defazio, Bach, and Lacoste-Julien, 2014)
- Simple proof
  - SVRG, SAGA
- Lower bounds
  - Agarwal and Bottou (2014)

### Dual stochastic coordinate ascent - I

• General learning formulation:

$$\begin{split} \min_{\theta \in \mathbb{R}^d} & \frac{1}{n} \sum_{i=1}^n \ell_i(\theta^\top \Phi(x_i)) + \frac{\mu}{2} \|\theta\|_2^2 \\ = & \min_{\theta \in \mathbb{R}^d, u \in \mathbb{R}^n} & \frac{1}{n} \sum_{i=1}^n \ell_i(u_i) + \frac{\mu}{2} \|\theta\|_2^2 \text{ such that } \forall i, u_i = \theta^\top \Phi(x_i) \\ = & \min_{\theta \in \mathbb{R}^d, u \in \mathbb{R}^n} & \max_{\alpha \in \mathbb{R}^n} & \frac{1}{n} \sum_{i=1}^n \ell_i(u_i) + \frac{\mu}{2} \|\theta\|_2^2 + \sum_{i=1}^n \alpha_i(u_i - \theta^\top \Phi(x_i)) \\ = & \max_{\alpha \in \mathbb{R}^n} \min_{\theta \in \mathbb{R}^d, u \in \mathbb{R}^n} & \frac{1}{n} \sum_{i=1}^n \ell_i(u_i) + \frac{\mu}{2} \|\theta\|_2^2 + \sum_{i=1}^n \alpha_i(u_i - \theta^\top \Phi(x_i)) \\ = & \max_{\alpha \in \mathbb{R}^n} \min_{\theta \in \mathbb{R}^d, u \in \mathbb{R}^n} & \frac{1}{n} \sum_{i=1}^n \ell_i(u_i) + \frac{\mu}{2} \|\theta\|_2^2 + \sum_{i=1}^n \alpha_i(u_i - \theta^\top \Phi(x_i)) \end{split}$$

### Dual stochastic coordinate ascent - II

• General learning formulation:

$$\begin{split} \min_{\theta \in \mathbb{R}^d} & \frac{1}{n} \sum_{i=1}^n \ell_i(\theta^\top \Phi(x_i)) + \frac{\mu}{2} \|\theta\|_2^2 \\ = & \max_{\alpha \in \mathbb{R}^n} \min_{\theta \in \mathbb{R}^d, u \in \mathbb{R}^n} & \frac{1}{n} \sum_{i=1}^n \ell_i(u_i) + \frac{\mu}{2} \|\theta\|_2^2 + \sum_{i=1}^n \alpha_i(u_i - \theta^\top \Phi(x_i)) \\ = & \max_{\alpha \in \mathbb{R}^n} & \sum_{i=1}^n \max_{u_i \in \mathbb{R}} \left\{ \frac{1}{n} \ell_i(u_i) + \alpha_i u_i \right\} - \frac{1}{2\mu} \left\| \sum_{i=1}^n \alpha_i \Phi(x_i) \right\|_2^2 \\ = & \max_{\alpha \in \mathbb{R}^n} & - \sum_{i=1}^n \psi_i(\alpha_i) - \frac{1}{2\mu} \left\| \sum_{i=1}^n \alpha_i \Phi(x_i) \right\|_2^2 \end{split}$$

- Minimizers obtained as  $\theta = \frac{1}{\mu} \sum_{i=1}^{n} \alpha_i \Phi(x_i)$ -  $\psi_i$  convex (up to affine transform = Fenchel-Legendre dual of  $\ell_i$ )

## Dual stochastic coordinate ascent - III

• General learning formulation:

$$\min_{\theta \in \mathbb{R}^d} \left\| \frac{1}{n} \sum_{i=1}^n \ell_i(\theta^\top \Phi(x_i)) + \frac{\mu}{2} \|\theta\|_2^2 = \max_{\alpha \in \mathbb{R}^n} \left\| -\sum_{i=1}^n \psi_i(\alpha_i) - \frac{1}{2\mu} \right\| \sum_{i=1}^n \alpha_i \Phi(x_i) \right\|_2^2$$

#### • From primal to dual

- $\ell_i$  smooth  $\Leftrightarrow \psi_i$  strongly convex
- $\ell_i$  strongly convex  $\Leftrightarrow \psi_i$  smooth
- Applying coordinate descent in the dual
  - Nesterov (2012); Shalev-Shwartz and Zhang (2012)
  - Linear convergence rate with simple iterations

## Dual stochastic coordinate ascent - IV

• Dual formulation: 
$$\max_{\alpha \in \mathbb{R}^n} - \sum_{i=1}^n \psi_i(\alpha_i) - \frac{1}{2\mu} \left\| \sum_{i=1}^n \alpha_i \Phi(x_i) \right\|_2^2$$

- Stochastic coordinate descent: at iteration t
  - Choose a coordinate i at random
  - Optimite w.r.t.  $\alpha_i: \max_{\alpha_i \in \mathbb{R}} -\psi_i(\alpha_i) \frac{1}{2\mu} \left\| \alpha_i \Phi(x_i) + \sum_{j \neq i} \alpha_i \Phi(x_i) \right\|_2^2$
  - Can be done by a single access to  $\Phi(x_i)$  and updating  $\sum_{i=1}^{j \neq i} \alpha_i \Phi(x_i)$

#### • Convergence proof

- See Nesterov (2012); Shalev-Shwartz and Zhang (2012)
- Similar linear convergence than SAG

# Summary and future work

- Constant-step-size averaged stochastic gradient descent
  - Reaches convergence rate O(1/n) in all regimes
  - Improves on the  $O(1/\sqrt{n})$  lower-bound of non-smooth problems
  - Efficient online Newton step for non-quadratic problems
  - Robustness to step-size selection
- Going beyond a single pass through the data

# Summary and future work

- Constant-step-size averaged stochastic gradient descent
  - Reaches convergence rate O(1/n) in all regimes
  - Improves on the  $O(1/\sqrt{n})$  lower-bound of non-smooth problems
  - Efficient online Newton step for non-quadratic problems
  - Robustness to step-size selection
- Going beyond a single pass through the data
- Extensions and future work
  - Pre-conditioning
  - Proximal extensions fo non-differentiable terms
  - kernels and non-parametric estimation
  - line-search
  - parallelization

# Outline

#### 1. Large-scale machine learning and optimization

- Traditional statistical analysis
- Classical methods for convex optimization

### 2. Non-smooth stochastic approximation

- Stochastic (sub)gradient and averaging
- Non-asymptotic results and lower bounds
- Strongly convex vs. non-strongly convex

### 3. Smooth stochastic approximation algorithms

- Asymptotic and non-asymptotic results
- 4. Beyond decaying step-sizes
- 5. Finite data sets

# Conclusions

# Machine learning and convex optimization

## • Statistics with or without optimization?

- Significance of mixing algorithms with analysis
- Benefits of mixing algorithms with analysis

## • Open problems

- Non-parametric stochastic approximation
- Going beyond a single pass over the data (testing performance)
- Characterization of implicit regularization of online methods
- Further links between convex optimization and online learning/bandits

# References

- A. Agarwal, P. L. Bartlett, P. Ravikumar, and M. J. Wainwright. Information-theoretic lower bounds on the oracle complexity of stochastic convex optimization. *Information Theory, IEEE Transactions* on, 58(5):3235–3249, 2012.
- Alekh Agarwal and Leon Bottou. A lower bound for the optimization of finite sums. *arXiv preprint arXiv:1410.0723*, 2014.
- R. Aguech, E. Moulines, and P. Priouret. On a perturbation approach for the analysis of stochastic tracking algorithms. *SIAM J. Control and Optimization*, 39(3):872–899, 2000.
- F. Bach. Adaptivity of averaged stochastic gradient descent to local strong convexity for logistic regression. Technical Report 00804431, HAL, 2013.
- F. Bach and E. Moulines. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Adv. NIPS*, 2011.
- F. Bach and E. Moulines. Non-strongly-convex smooth stochastic approximation with convergence rate o(1/n). Technical Report 00831977, HAL, 2013.
- F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Structured sparsity through convex optimization, 2012a.
- Francis Bach, Rodolphe Jenatton, Julien Mairal, and Guillaume Obozinski. Optimization with sparsityinducing penalties. *Foundations and Trends*® *in Machine Learning*, 4(1):1–106, 2012b.
- A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.

- Albert Benveniste, Michel Métivier, and Pierre Priouret. *Adaptive algorithms and stochastic approximations*. Springer Publishing Company, Incorporated, 2012.
- D. P. Bertsekas. A new class of incremental gradient methods for least squares problems. *SIAM Journal on Optimization*, 7(4):913–926, 1997.
- D. Blatt, A. O. Hero, and H. Gauchman. A convergent incremental gradient method with a constant step size. *SIAM Journal on Optimization*, 18(1):29–51, 2008.
- V. S. Borkar. Stochastic approximation with two time scales. *Systems & Control Letters*, 29(5): 291–294, 1997.
- L. Bottou and O. Bousquet. The tradeoffs of large scale learning. In Adv. NIPS, 2008.
- L. Bottou and Y. Le Cun. On-line learning for very large data sets. *Applied Stochastic Models in Business and Industry*, 21(2):137–151, 2005.
- S. Boucheron and P. Massart. A high-dimensional wilks phenomenon. *Probability theory and related fields*, 150(3-4):405–433, 2011.
- S. Boucheron, O. Bousquet, G. Lugosi, et al. Theory of classification: A survey of some recent advances. *ESAIM Probability and statistics*, 9:323–375, 2005.
- N. Cesa-Bianchi, A. Conconi, and C. Gentile. On the generalization ability of on-line learning algorithms. *Information Theory, IEEE Transactions on*, 50(9):2050–2057, 2004.
- Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems*, pages 1646–1654, 2014.
- A. Défossez and F. Bach. Constant step size least-mean-square: Bias-variance trade-offs and optimal

sampling distributions. 2015.

- B. Delyon and A. Juditsky. Accelerated stochastic approximation. *SIAM Journal on Optimization*, 3: 868–881, 1993.
- A. Dieuleveut and F. Bach. Non-parametric Stochastic Approximation with Large Step sizes. Technical report, ArXiv, 2014.
- J. Duchi and Y. Singer. Efficient online and batch learning using forward backward splitting. *Journal of Machine Learning Research*, 10:2899–2934, 2009. ISSN 1532-4435.
- N. Flammarion and F. Bach. From averaging to acceleration, there is only a step-size. *arXiv preprint arXiv:1504.01577*, 2015.
- M. P. Friedlander and M. Schmidt. Hybrid deterministic-stochastic methods for data fitting. *arXiv:1104.2373*, 2011.
- S. Ghadimi and G. Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization. *Optimization Online*, July, 2010.
- Saeed Ghadimi and Guanghui Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization, ii: shrinking procedures and optimal algorithms. *SIAM Journal on Optimization*, 23(4):2061–2089, 2013.
- L. Györfi and H. Walk. On the averaged stochastic approximation for linear regression. *SIAM Journal* on Control and Optimization, 34(1):31–61, 1996.
- E. Hazan, A. Agarwal, and S. Kale. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 69(2):169–192, 2007.

Chonghai Hu, James T Kwok, and Weike Pan. Accelerated gradient methods for stochastic optimization

and online learning. In NIPS, volume 22, pages 781-789, 2009.

- Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In Advances in Neural Information Processing Systems, pages 315–323, 2013.
- Takafumi Kanamori and Hidetoshi Shimodaira. Active learning algorithm using the maximum weighted log-likelihood estimator. *Journal of statistical planning and inference*, 116(1):149–162, 2003.
- H. Kesten. Accelerated stochastic approximation. Ann. Math. Stat., 29(1):41-59, 1958.
- H. J. Kushner and G. G. Yin. *Stochastic approximation and recursive algorithms and applications*. Springer-Verlag, second edition, 2003.
- S. Lacoste-Julien, M. Schmidt, and F. Bach. A simpler approach to obtaining an o (1/t) convergence rate for projected stochastic subgradient descent. Technical Report 1212.2002, ArXiv, 2012.
- Guanghui Lan, Arkadi Nemirovski, and Alexander Shapiro. Validation analysis of mirror descent stochastic approximation method. *Mathematical programming*, 134(2):425–458, 2012.
- N. Le Roux, M. Schmidt, and F. Bach. A stochastic gradient method with an exponential convergence rate for strongly-convex optimization with finite training sets. In *Adv. NIPS*, 2012.
- N. Le Roux, M. Schmidt, and F. Bach. A stochastic gradient method with an exponential convergence rate for strongly-convex optimization with finite training sets. Technical Report 00674995, HAL, 2013.
- O. Macchi. Adaptive processing: The least mean squares approach with applications in transmission. Wiley West Sussex, 1995.
- A. Nedic and D. Bertsekas. Convergence rate of incremental subgradient algorithms. *Stochastic Optimization: Algorithms and Applications*, pages 263–304, 2000.

- A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.
- A. S. Nemirovsky and D. B. Yudin. Problem complexity and method efficiency in optimization. Wiley & Sons, 1983.
- Y. Nesterov. A method for solving a convex programming problem with rate of convergence  $O(1/k^2)$ . Soviet Math. Doklady, 269(3):543–547, 1983.
- Y. Nesterov. *Introductory lectures on convex optimization: a basic course*. Kluwer Academic Publishers, 2004.
- Y. Nesterov. Gradient methods for minimizing composite objective function. *Center for Operations Research and Econometrics (CORE), Catholic University of Louvain, Tech. Rep*, 76, 2007.
- Y. Nesterov. Primal-dual subgradient methods for convex problems. *Mathematical programming*, 120 (1):221–259, 2009.
- Y. Nesterov and A. Nemirovski. *Interior-point polynomial algorithms in convex programming*. SIAM studies in Applied Mathematics, 1994.
- Y. Nesterov and J. P. Vial. Confidence level solutions for stochastic programming. *Automatica*, 44(6): 1559–1568, 2008. ISSN 0005-1098.
- Yu Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- B. T. Polyak and A. B. Juditsky. Acceleration of stochastic approximation by averaging. *SIAM Journal* on Control and Optimization, 30(4):838–855, 1992.
- H. Robbins and S. Monro. A stochastic approximation method. Ann. Math. Statistics, 22:400-407,
1951. ISSN 0003-4851.

- D. Ruppert. Efficient estimations from a slowly convergent Robbins-Monro process. Technical Report 781, Cornell University Operations Research and Industrial Engineering, 1988.
- M. Schmidt, N. Le Roux, and F. Bach. Optimization with approximate gradients. Technical report, HAL, 2011.
- B. Schölkopf and A. J. Smola. Learning with Kernels. MIT Press, 2001.
- S. Shalev-Shwartz and N. Srebro. SVM optimization: inverse dependence on training set size. In *Proc. ICML*, 2008.
- S. Shalev-Shwartz and T. Zhang. Stochastic dual coordinate ascent methods for regularized loss minimization. Technical Report 1209.1873, Arxiv, 2012.
- S. Shalev-Shwartz, Y. Singer, and N. Srebro. Pegasos: Primal estimated sub-gradient solver for svm. In *Proc. ICML*, 2007.
- S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan. Stochastic convex optimization. In *proc. COLT*, 2009.
- J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- Naum Zuselevich Shor, Krzysztof C. Kiwiel, and Andrzej Ruszcay?ski. *Minimization methods for non-differentiable functions*. Springer-Verlag New York, Inc., 1985.
- M.V. Solodov. Incremental gradient algorithms with stepsizes bounded away from zero. *Computational Optimization and Applications*, 11(1):23–35, 1998.
- K. Sridharan, N. Srebro, and S. Shalev-Shwartz. Fast rates for regularized objectives. 2008.

- P. Sunehag, J. Trumpf, SVN Vishwanathan, and N. Schraudolph. Variable metric stochastic approximation theory. *International Conference on Artificial Intelligence and Statistics*, 2009.
- P. Tseng. An incremental gradient(-projection) method with momentum term and adaptive stepsize rule. *SIAM Journal on Optimization*, 8(2):506–531, 1998.
- A. B. Tsybakov. Optimal rates of aggregation. 2003.
- A. W. Van der Vaart. Asymptotic statistics, volume 3. Cambridge Univ. press, 2000.
- L. Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *Journal of Machine Learning Research*, 9:2543–2596, 2010. ISSN 1532-4435.