

Stochastic gradient descent and robustness to ill-conditioning

Francis Bach

INRIA - Ecole Normale Supérieure, Paris, France



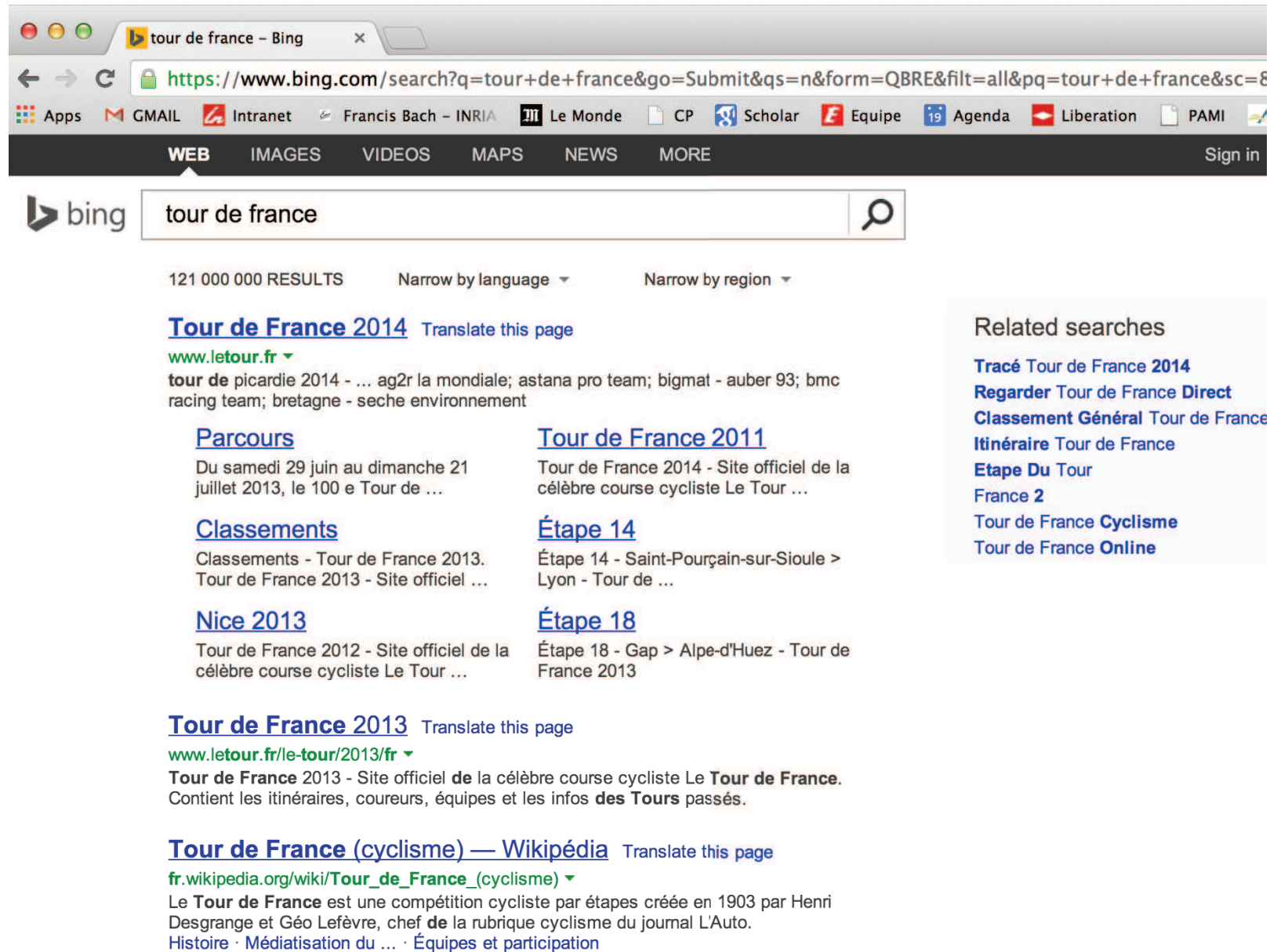
Joint work with Aymeric Dieuleveut, Nicolas Flammarion,
Eric Moulines - Ecole des Ponts, June 2016

“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 p

Search engines - Advertising



The image shows a screenshot of a web browser displaying a Bing search results page for the query "tour de france". The browser's address bar shows the URL: <https://www.bing.com/search?q=tour+de+france&go=Submit&qsn=n&form=QBRE&filt=all&pq=tour+de+france&sc=8>. The search bar contains the text "tour de france". Below the search bar, the results show 121,000,000 results. The first result is for "Tour de France 2014" from www.letour.fr, with a snippet: "tour de picardie 2014 - ... ag2r la mondiale; astana pro team; bigmat - auber 93; bmc racing team; bretagne - seche environnement". Other results include "Parcours", "Classements", "Nice 2013", "Tour de France 2011", "Étape 14", "Étape 18", "Tour de France 2013", and "Tour de France (cyclisme) — Wikipédia". A "Related searches" sidebar on the right lists: "Tracé Tour de France 2014", "Regarder Tour de France Direct", "Classement Général Tour de France", "Itinéraire Tour de France", "Étape Du Tour", "France 2", "Tour de France Cyclisme", and "Tour de France Online".

tour de france - Bing

<https://www.bing.com/search?q=tour+de+france&go=Submit&qsn=n&form=QBRE&filt=all&pq=tour+de+france&sc=8>

Apps GMAIL Intranet Francis Bach - INRIA Le Monde CP Scholar Equipe 19 Agenda Liberation PAMI

WEB IMAGES VIDEOS MAPS NEWS MORE Sign in

bing tour de france

121 000 000 RESULTS Narrow by language Narrow by region

Tour de France 2014 [Translate this page](#)
www.letour.fr
tour de picardie 2014 - ... ag2r la mondiale; astana pro team; bigmat - auber 93; bmc racing team; bretagne - seche environnement

Parcours
Du samedi 29 juin au dimanche 21 juillet 2013, le 100 e Tour de ...

Classements
Classements - Tour de France 2013. Tour de France 2013 - Site officiel ...

Nice 2013
Tour de France 2012 - Site officiel de la célèbre course cycliste Le Tour ...

Tour de France 2011
Tour de France 2014 - Site officiel de la célèbre course cycliste Le Tour ...

Étape 14
Étape 14 - Saint-Pourçain-sur-Sioule > Lyon - Tour de ...

Étape 18
Étape 18 - Gap > Alpe-d'Huez - Tour de France 2013

Tour de France 2013 [Translate this page](#)
www.letour.fr/le-tour/2013/fr
Tour de France 2013 - Site officiel de la célèbre course cycliste Le Tour de France. Contient les itinéraires, coureurs, équipes et les infos des Tours passés.

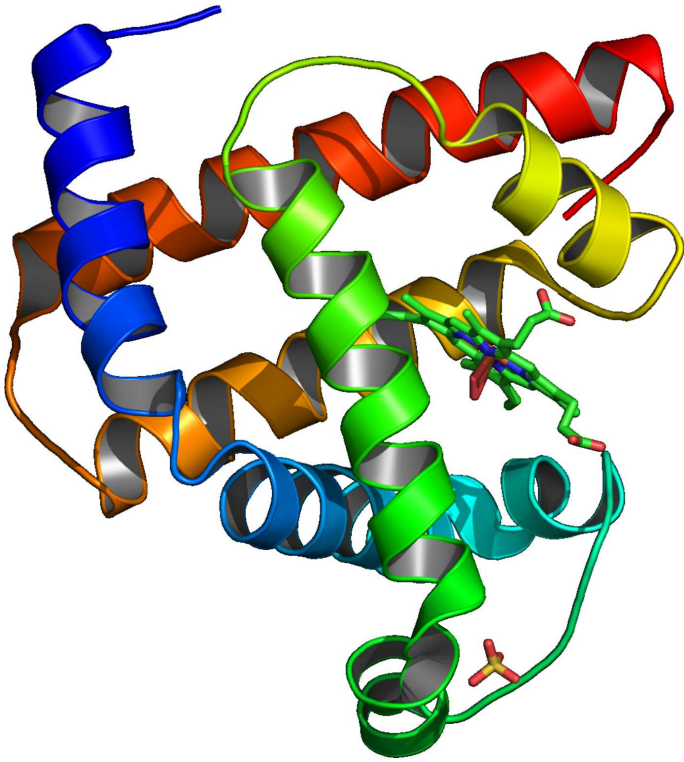
Tour de France (cyclisme) — Wikipédia [Translate this page](#)
[fr.wikipedia.org/wiki/Tour_de_France_\(cyclisme\)](http://fr.wikipedia.org/wiki/Tour_de_France_(cyclisme))
Le Tour de France est une compétition cycliste par étapes créée en 1903 par Henri Desgrange et Géo Lefèvre, chef de la rubrique cyclisme du journal L'Auto.
[Histoire](#) · [Médiatisation du ...](#) · [Équipes et participation](#)

Related searches
[Tracé Tour de France 2014](#)
[Regarder Tour de France Direct](#)
[Classement Général Tour de France](#)
[Itinéraire Tour de France](#)
[Étape Du Tour](#)
[France 2](#)
[Tour de France Cyclisme](#)
[Tour de France Online](#)

Visual object recognition



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 p , large n**
 - p : 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 p , large n**
 - p : dimension of each observation (input)
 - n : number of observations
- **Examples:** computer vision, bioinformatics, advertising
- **Ideal running-time complexity:** $O(pn)$

Context

Machine learning for “big data”

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

Supervised machine learning

- **Data:** n observations $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$, $i = 1, \dots, n$, **i.i.d.**
- Prediction as a **linear function** $\langle \theta, \Phi(x) \rangle$ of features $\Phi(x) \in \mathbb{R}^p$
 - Explicit features adapted to inputs (can be learned as well)
 - Using Hilbert spaces for non-linear / non-parametric estimation

Supervised machine learning

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

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

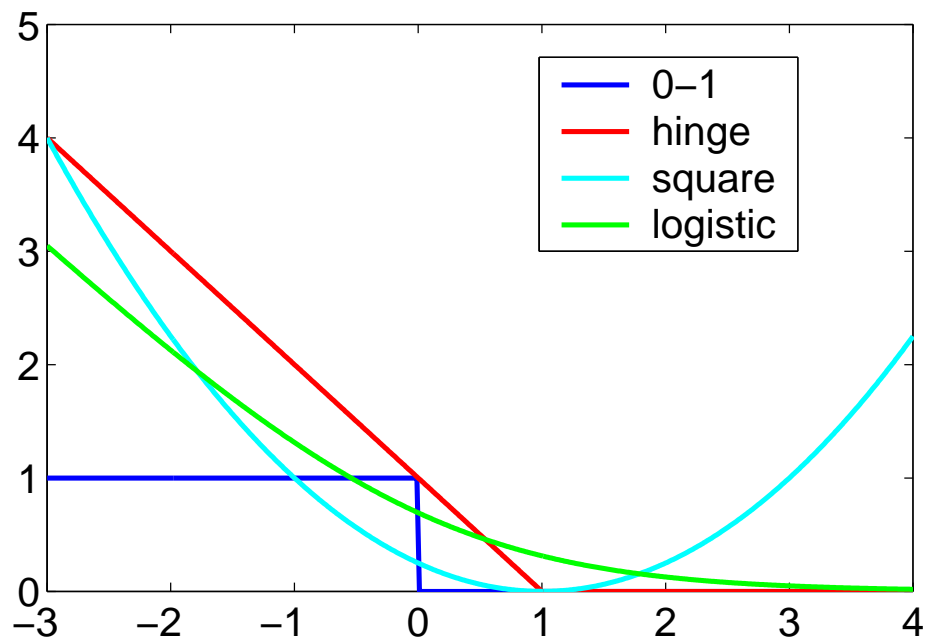
convex data fitting term + regularizer

Usual losses

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

Usual losses

- **Regression:** $y \in \mathbb{R}$, prediction $\hat{y} = \langle \theta, \Phi(x) \rangle$
 - quadratic loss $\frac{1}{2}(y - \hat{y})^2 = \frac{1}{2}(y - \langle \theta, \Phi(x) \rangle)^2$
- **Classification :** $y \in \{-1, 1\}$, prediction $\hat{y} = \text{sign}(\langle \theta, \Phi(x) \rangle)$
 - loss of the form $\ell(y \langle \theta, \Phi(x) \rangle)$
 - “True” **0-1** loss: $\ell(y \langle \theta, \Phi(x) \rangle) = 1_{y \langle \theta, \Phi(x) \rangle < 0}$
 - Usual **convex** losses:



Supervised machine learning

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

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

convex data fitting term + regularizer

Supervised machine learning

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

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

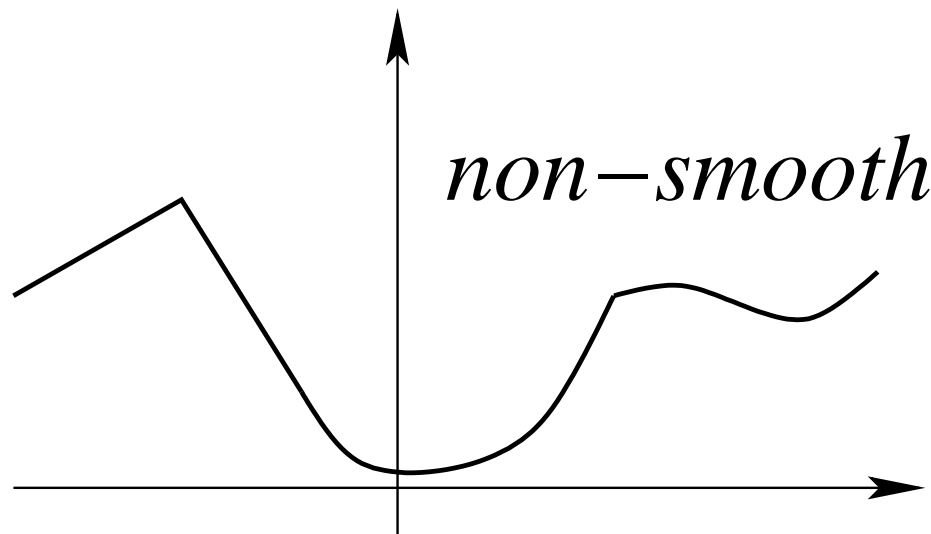
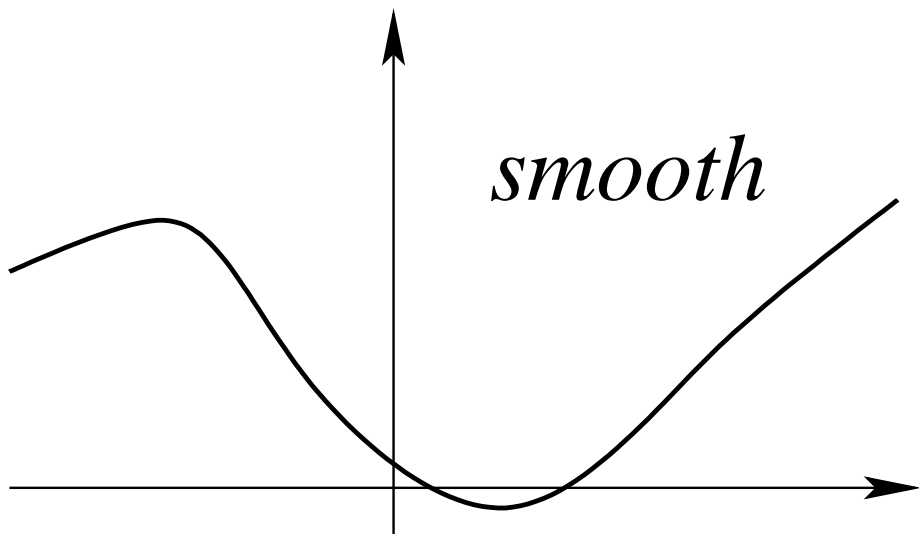
convex data fitting term + regularizer

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

Smoothness and strong convexity

- A function $g : \mathbb{R}^p \rightarrow \mathbb{R}$ is L -smooth if and only if it is twice differentiable and

$$\forall \theta \in \mathbb{R}^p, \text{ eigenvalues}[g''(\theta)] \leq L$$



Smoothness and strong convexity

- A function $g : \mathbb{R}^p \rightarrow \mathbb{R}$ is **L -smooth** if and only if it is twice differentiable and

$$\forall \theta \in \mathbb{R}^p, \text{ eigenvalues}[g''(\theta)] \leq L$$

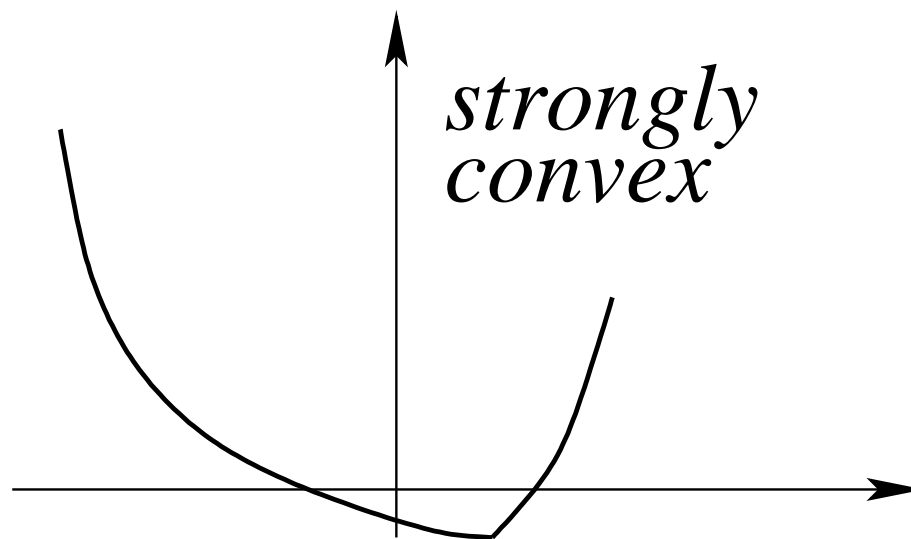
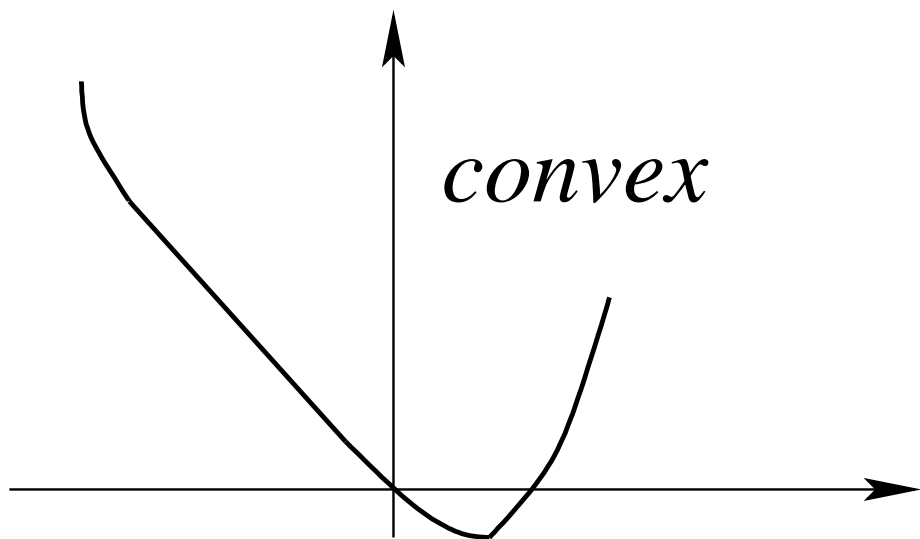
- **Machine learning**

- with $g(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \langle \theta, \Phi(x_i) \rangle)$
- Hessian \approx covariance matrix $\frac{1}{n} \sum_{i=1}^n \Phi(x_i) \otimes \Phi(x_i)$
- **Bounded data**: $\|\Phi(x)\| \leq R \Rightarrow L = O(R^2)$

Smoothness and strong convexity

- A twice differentiable function $g : \mathbb{R}^p \rightarrow \mathbb{R}$ is μ -strongly convex if and only if

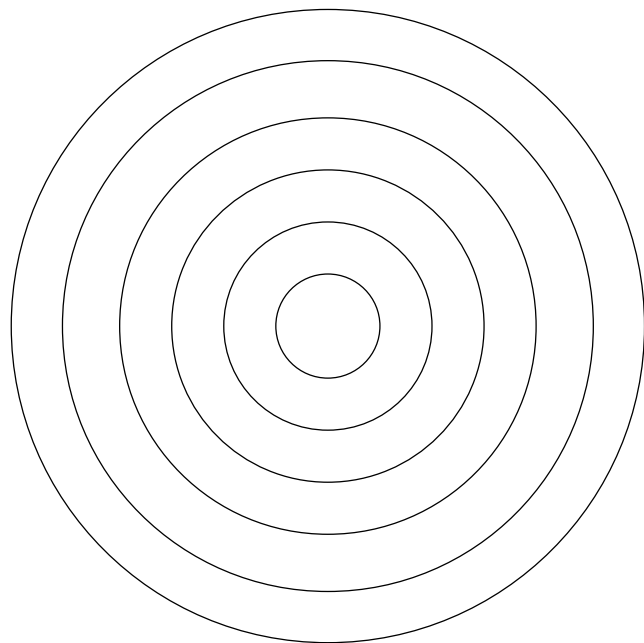
$$\forall \theta \in \mathbb{R}^p, \text{ eigenvalues}[g''(\theta)] \geq \mu$$



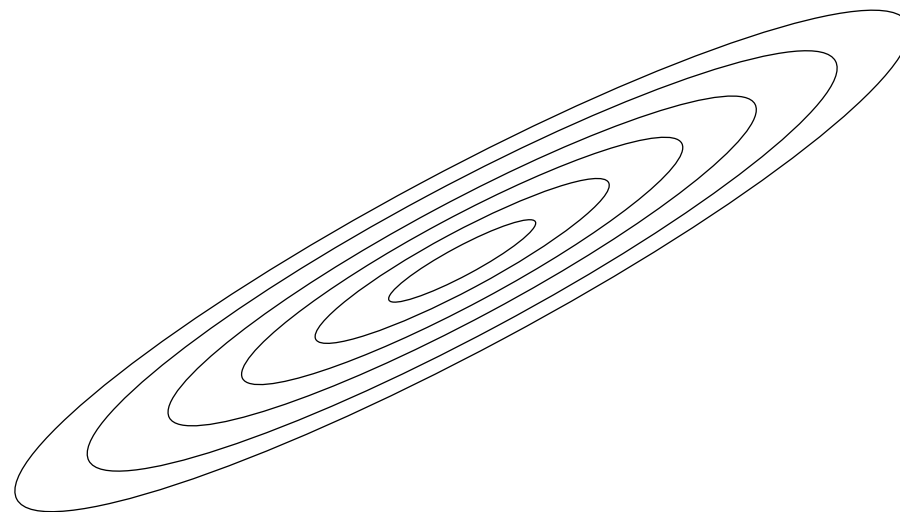
Smoothness and strong convexity

- A twice differentiable function $g : \mathbb{R}^p \rightarrow \mathbb{R}$ is μ -strongly convex if and only if

$$\forall \theta \in \mathbb{R}^p, \text{ eigenvalues}[g''(\theta)] \geq \mu$$



(large μ/L)



(small μ/L)

Smoothness and strong convexity

- A twice differentiable function $g : \mathbb{R}^p \rightarrow \mathbb{R}$ is μ -strongly convex if and only if

$$\forall \theta \in \mathbb{R}^p, \text{ eigenvalues}[g''(\theta)] \geq \mu$$

- **Machine learning**

- with $g(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \langle \theta, \Phi(x_i) \rangle)$
- Hessian \approx covariance matrix $\frac{1}{n} \sum_{i=1}^n \Phi(x_i) \otimes \Phi(x_i)$
- **Data with invertible covariance matrix** (low correlation/dimension)

Smoothness and strong convexity

- A twice differentiable function $g : \mathbb{R}^p \rightarrow \mathbb{R}$ is μ -strongly convex if and only if

$$\forall \theta \in \mathbb{R}^p, \text{ eigenvalues}[g''(\theta)] \geq \mu$$

- **Machine learning**

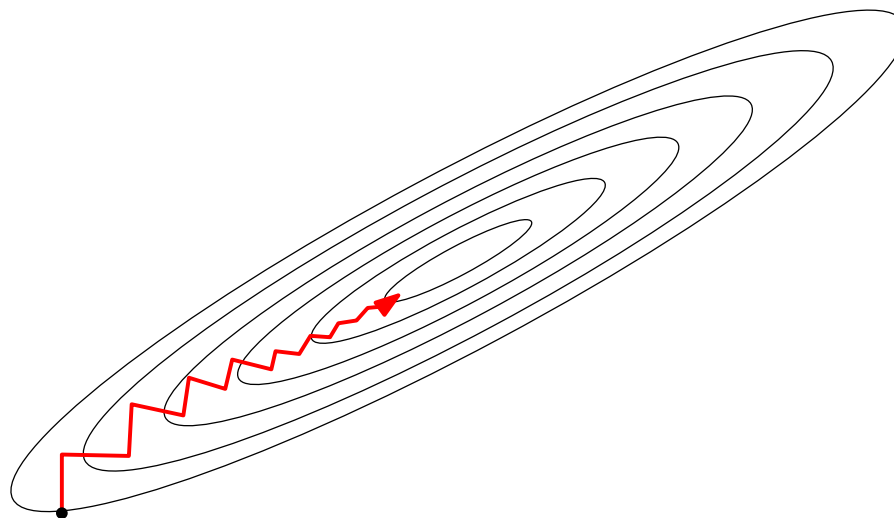
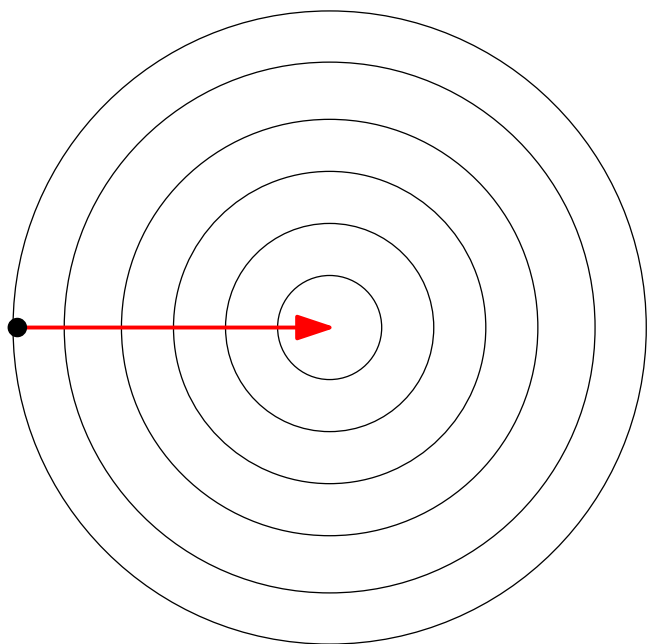
- with $g(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \langle \theta, \Phi(x_i) \rangle)$
- Hessian \approx covariance matrix $\frac{1}{n} \sum_{i=1}^n \Phi(x_i) \otimes \Phi(x_i)$
- **Data with invertible covariance matrix** (low correlation/dimension)

- **Adding regularization by $\frac{\mu}{2} \|\theta\|^2$**

- **creates additional bias unless μ is small**

Iterative methods for minimizing smooth functions

- **Assumption:** g convex and smooth on \mathbb{R}^p
- **Gradient descent:** $\theta_t = \theta_{t-1} - \gamma_t g'(\theta_{t-1})$
 - $O(1/t)$ convergence rate for convex functions
 - $O(e^{-(\mu/L)t})$ convergence rate for strongly convex functions



Iterative methods for minimizing smooth functions

- **Assumption:** g convex and smooth on \mathbb{R}^p
- **Gradient descent:** $\theta_t = \theta_{t-1} - \gamma_t g'(\theta_{t-1})$
 - $O(1/t)$ convergence rate for convex functions
 - $O(e^{-(\mu/L)t})$ convergence rate for strongly convex functions
- **Newton method:** $\theta_t = \theta_{t-1} - g''(\theta_{t-1})^{-1}g'(\theta_{t-1})$
 - $O(e^{-\rho 2^t})$ convergence rate

Iterative methods for minimizing smooth functions

- **Assumption:** g convex and smooth on \mathbb{R}^p
- **Gradient descent:** $\theta_t = \theta_{t-1} - \gamma_t g'(\theta_{t-1})$
 - $O(1/t)$ convergence rate for convex functions
 - $O(e^{-(\mu/L)t})$ convergence rate for strongly convex functions
- **Newton method:** $\theta_t = \theta_{t-1} - g''(\theta_{t-1})^{-1}g'(\theta_{t-1})$
 - $O(e^{-\rho 2^t})$ convergence rate
- **Key insights from Bottou and Bousquet (2008)**
 1. In machine learning, no need to optimize below statistical error
 2. In machine learning, cost functions are averages

\Rightarrow **Stochastic approximation**

Stochastic approximation

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

Stochastic approximation

- **Goal:** Minimizing a function f defined on \mathbb{R}^p
 - given only unbiased estimates $f'_n(\theta_n)$ of its gradients $f'(\theta_n)$ at certain points $\theta_n \in \mathbb{R}^p$
- **Machine learning - statistics**
 - $f(\theta) = \mathbb{E} f_n(\theta) = \mathbb{E} \ell(y_n, \langle \theta, \Phi(x_n) \rangle) =$ **generalization error**
 - **Loss for a single pair of observations:** $f_n(\theta) = \ell(y_n, \langle \theta, \Phi(x_n) \rangle)$
 - Expected gradient:

$$f'(\theta) = \mathbb{E} f'_n(\theta) = \mathbb{E} \{ \ell'(y_n, \langle \theta, \Phi(x_n) \rangle) \Phi(x_n) \}$$

- Beyond convex optimization: see, e.g., Benveniste et al. (2012)

Convex stochastic approximation

- **Key assumption:** smoothness and/or strong convexity
- **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: $\bar{\theta}_n = \frac{1}{n+1} \sum_{k=0}^n \theta_k$

– Which learning rate sequence γ_n ? Classical setting:

$$\gamma_n = Cn^{-\alpha}$$

Convex stochastic approximation

- **Key assumption:** smoothness and/or strong convexity
- **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: $\bar{\theta}_n = \frac{1}{n+1} \sum_{k=0}^n \theta_k$

– Which learning rate sequence γ_n ? Classical setting:

$$\gamma_n = Cn^{-\alpha}$$

- **Running-time** = $O(np)$
 - Single pass through the data
 - One line of code among many

Convex stochastic approximation

Existing analysis

- 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 $\gamma_n \propto n^{-1/2}$

Convex stochastic approximation

Existing analysis

- **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

Convex stochastic approximation

Existing analysis

- 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 algorithm for smooth problems with global 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 $\theta \in \mathbb{R}^p$
 - SGD = least-mean-square algorithm (see, e.g., Macchi, 1995)
 - usually studied without averaging and decreasing step-sizes
 - with strong convexity assumption $\mathbb{E}[\Phi(x_n) \otimes \Phi(x_n)] = H \succcurlyeq \mu \cdot \text{Id}$

Least-mean-square algorithm

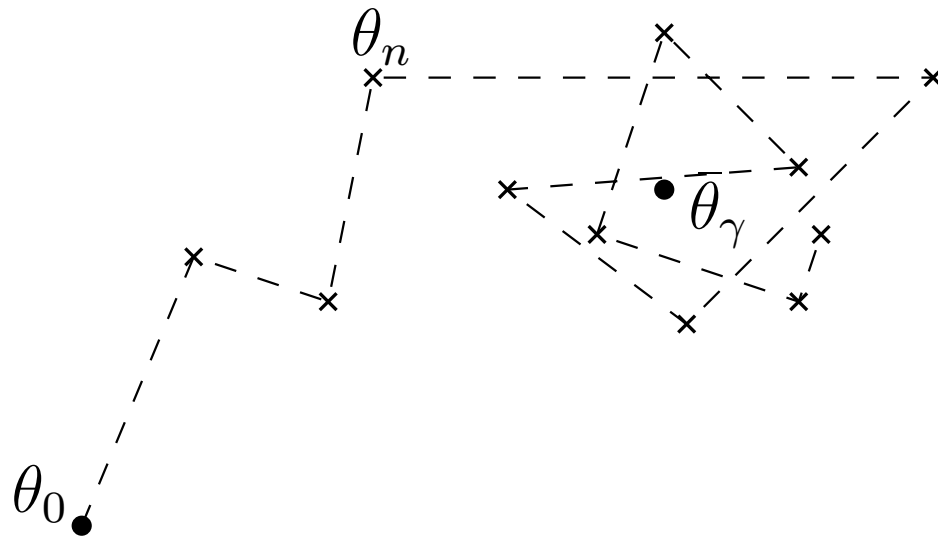
- **Least-squares:** $f(\theta) = \frac{1}{2}\mathbb{E}[(y_n - \langle \Phi(x_n), \theta \rangle)^2]$ with $\theta \in \mathbb{R}^p$
 - SGD = least-mean-square algorithm (see, e.g., Macchi, 1995)
 - usually studied without averaging and decreasing step-sizes
 - with strong convexity assumption $\mathbb{E}[\Phi(x_n) \otimes \Phi(x_n)] = H \succcurlyeq \mu \cdot \text{Id}$
- **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:
$$\mathbb{E}f(\bar{\theta}_n) - f(\theta_*) \leq \frac{4\sigma^2 p}{n} + \frac{4R^2 \|\theta_0 - \theta_*\|^2}{n}$$
- **Matches statistical lower bound** (Tsybakov, 2003)
 - Non-asymptotic robust version of Györfi and Walk (1996)

Markov chain interpretation of constant step sizes

- LMS recursion for $f_n(\theta) = \frac{1}{2}(y_n - \langle \Phi(x_n), \theta \rangle)^2$

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

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



Markov chain interpretation of constant step sizes

- LMS recursion for $f_n(\theta) = \frac{1}{2}(y_n - \langle \Phi(x_n), \theta \rangle)^2$

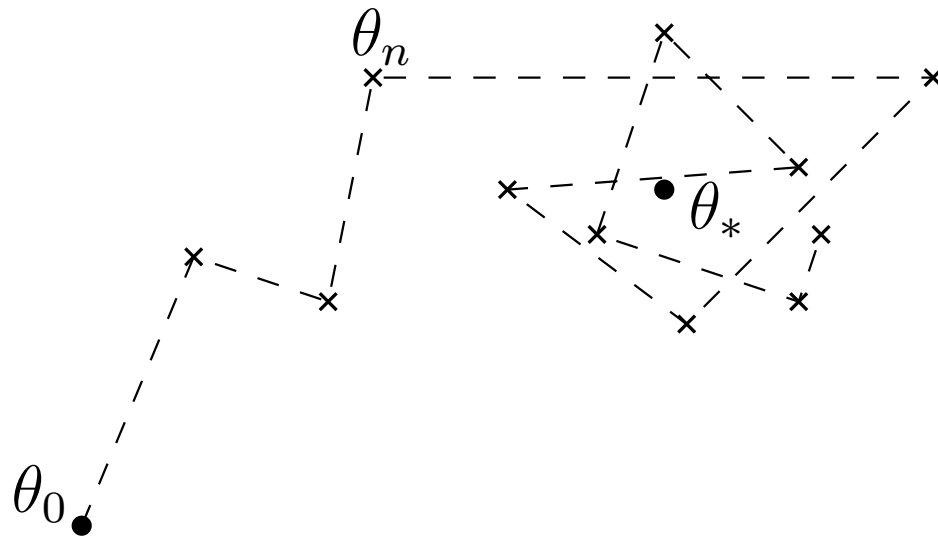
$$\theta_n = \theta_{n-1} - \gamma(\langle \Phi(x_n), \theta_{n-1} \rangle - y_n)\Phi(x_n)$$

- The sequence $(\theta_n)_n$ is a **homogeneous Markov chain**

– convergence to a stationary distribution π_γ

– with expectation $\bar{\theta}_\gamma \stackrel{\text{def}}{=} \int \theta \pi_\gamma(d\theta)$

- **For least-squares, $\bar{\theta}_\gamma = \theta_*$**



Markov chain interpretation of constant step sizes

- LMS recursion for $f_n(\theta) = \frac{1}{2}(y_n - \langle \Phi(x_n), \theta \rangle)^2$

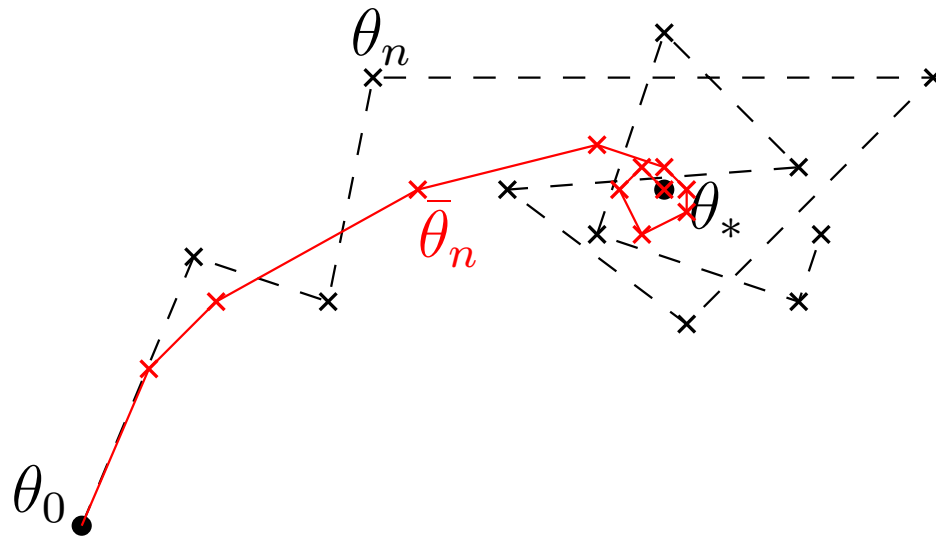
$$\theta_n = \theta_{n-1} - \gamma(\langle \Phi(x_n), \theta_{n-1} \rangle - y_n)\Phi(x_n)$$

- The sequence $(\theta_n)_n$ is a **homogeneous Markov chain**

– convergence to a stationary distribution π_γ

– with expectation $\bar{\theta}_\gamma \stackrel{\text{def}}{=} \int \theta \pi_\gamma(d\theta)$

- **For least-squares, $\bar{\theta}_\gamma = \theta_*$**



Markov chain interpretation of constant step sizes

- LMS recursion for $f_n(\theta) = \frac{1}{2}(y_n - \langle \Phi(x_n), \theta \rangle)^2$

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

- The sequence $(\theta_n)_n$ is a **homogeneous Markov chain**

- convergence to a stationary distribution π_γ

- with expectation $\bar{\theta}_\gamma \stackrel{\text{def}}{=} \int \theta \pi_\gamma(d\theta)$

- **For least-squares, $\bar{\theta}_\gamma = \theta_*$**

- θ_n does not converge to θ_* but oscillates around it

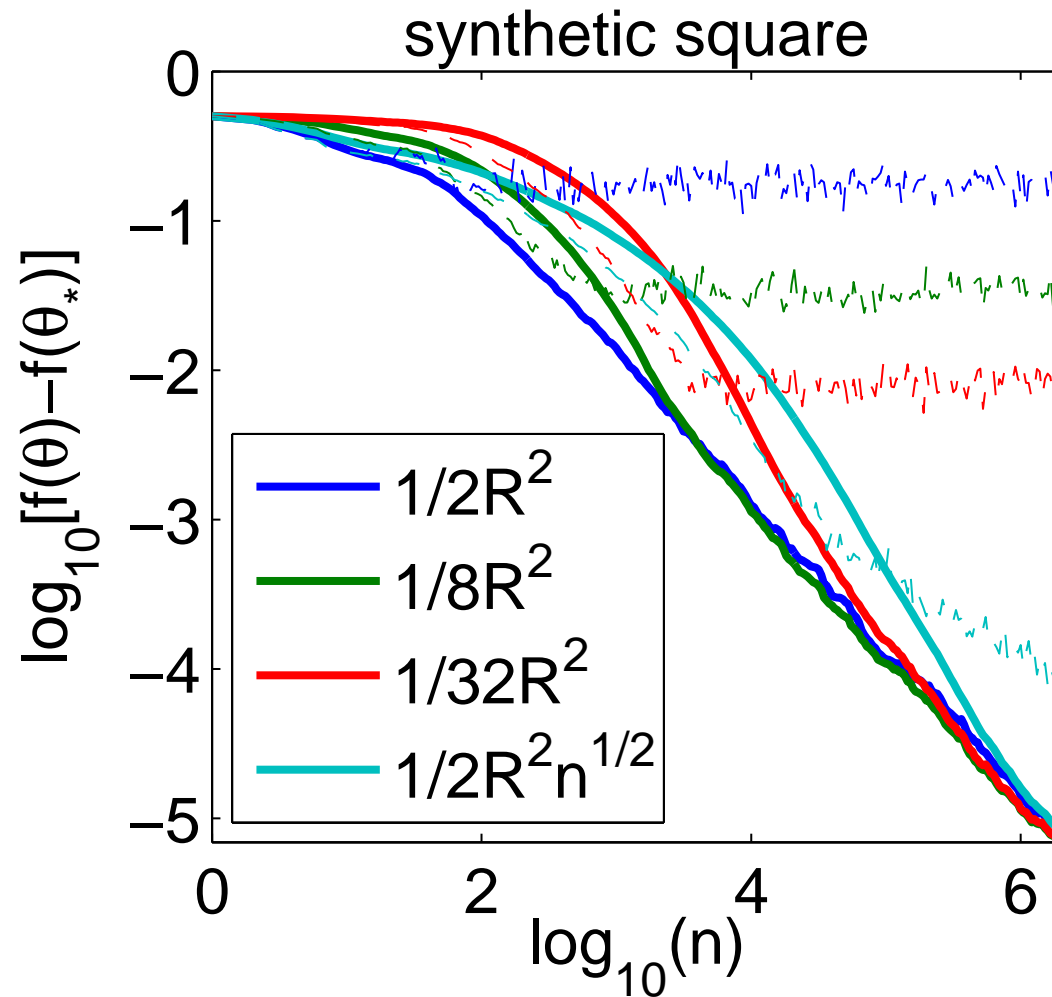
- oscillations of order $\sqrt{\gamma}$

- **Ergodic theorem:**

- Averaged iterates converge to $\bar{\theta}_\gamma = \theta_*$ at rate $O(1/n)$

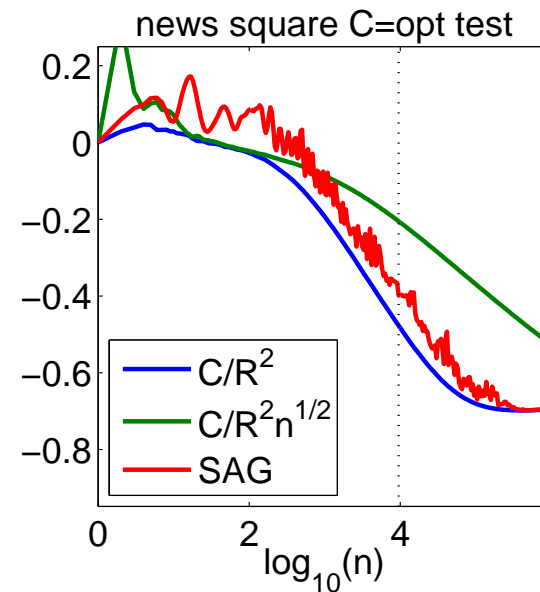
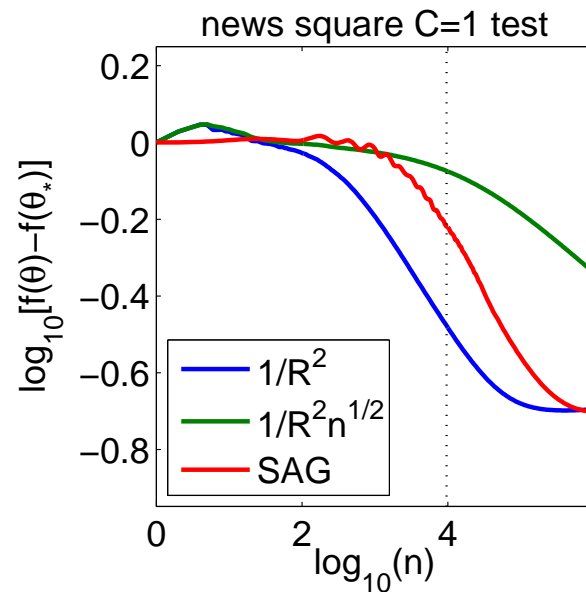
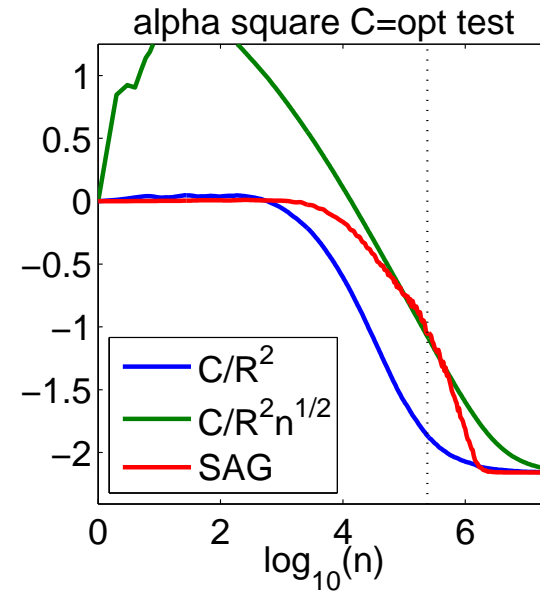
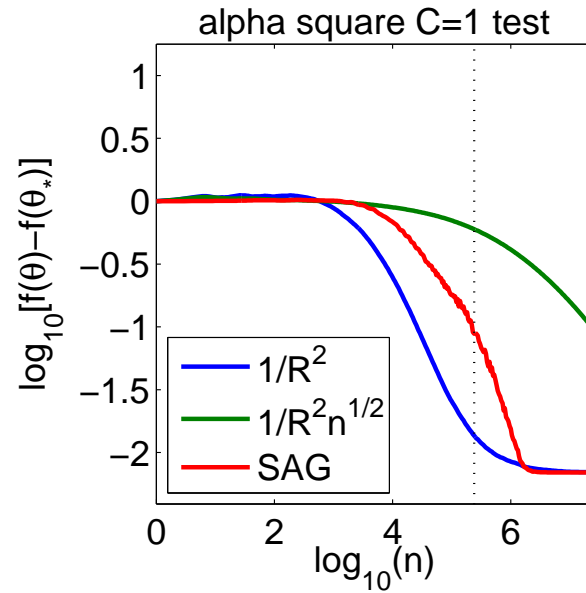
Simulations - synthetic examples

- Gaussian distributions - $p = 20$



Simulations - benchmarks

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



Isn't least-squares regression a "regression"?

Isn't least-squares regression a “regression”?

- **Least-squares regression**
 - Simpler to analyze and understand
 - **Explicit relationship to bias/variance trade-offs** (next slides)
- **Many important loss functions are not quadratic**
 - **Beyond least-squares with online Newton steps**
 - Complexity of $O(p)$ per iteration with rate $O(p/n)$
 - See Bach and Moulines (2013) for details

Optimal bounds for least-squares?

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

Optimal bounds for least-squares?

- **Least-squares:** cannot beat $\sigma^2 p/n$ (Tsybakov, 2003). Really?
- **Refined analysis** (Défossez and Bach, 2015)

$$\mathbb{E}f(\bar{\theta}_n) - f(\theta_*) \leq \frac{\sigma^2 p}{n} + \frac{\|H^{-1/2}(\theta_0 - \theta_*)\|_2^2}{\gamma^2 n^2}$$

- In practice: bias may be larger than variance, $\sigma^2 p/n$ pessimistic

Optimal bounds for least-squares?

- **Least-squares:** cannot beat $\sigma^2 p/n$ (Tsybakov, 2003). Really?
- **Refined analysis** (Défossez and Bach, 2015)

$$\mathbb{E}f(\bar{\theta}_n) - f(\theta_*) \leq \frac{\sigma^2 p}{n} + \frac{\|H^{-1/2}(\theta_0 - \theta_*)\|_2^2}{\gamma^2 n^2}$$

– In practice: bias may be larger than variance, $\sigma^2 p/n$ pessimistic

- **Refined assumptions with adaptivity** (Dieuleveut and Bach, 2014)

$$\mathbb{E}f(\bar{\theta}_n) - f(\theta_*) \leq \frac{\sigma^2 \gamma^{1/\alpha} \text{tr } H^{1/\alpha}}{n^{1-1/\alpha}} + \frac{\|H^{1/2-r}(\theta_0 - \theta_*)\|_2^2}{\gamma^{2r} n^{2 \min\{r, 1\}}}$$

- SGD is adaptive to the covariance matrix eigenvalue decay
- Leads to optimal rates for non-parametric regression

Achieving optimal bias and variance terms

- Current results with averaged SGD

- **Variance** (starting from optimal θ_*) = $\frac{\sigma^2 p}{n}$

- **Bias** (no noise) = $\min \left\{ \frac{R^2 \|\theta_0 - \theta_*\|^2}{n}, \frac{R^4 \langle \theta_0 - \theta_*, H^{-1}(\theta_0 - \theta_*) \rangle}{n^2} \right\}$

Achieving optimal bias and variance terms

- **Current results with averaged SGD** (ill-conditioned problems)

- **Variance** (starting from optimal θ_*) = $\frac{\sigma^2 p}{n}$

- **Bias** (no noise) = $\frac{R^2 \|\theta_0 - \theta_*\|^2}{n}$

Achieving optimal bias and variance terms

- **Current results with averaged SGD** (ill-conditioned problems)

- **Variance** (starting from optimal θ_*) = $\frac{\sigma^2 p}{n}$

- **Bias** (no noise) = $\frac{R^2 \|\theta_0 - \theta_*\|^2}{n}$

| | Bias | Variance |
|---|---|------------------------|
| Averaged gradient descent (Bach and Moulines, 2013) | $\frac{R^2 \ \theta_0 - \theta_*\ ^2}{n}$ | $\frac{\sigma^2 p}{n}$ |

Achieving optimal bias and variance terms

| | Bias | Variance |
|---|---|------------------------|
| Averaged gradient descent (Bach and Moulines, 2013) | $\frac{R^2 \ \theta_0 - \theta_*\ ^2}{n}$ | $\frac{\sigma^2 p}{n}$ |

Achieving optimal bias and variance terms

| | Bias | Variance |
|---|---|------------------------|
| Averaged gradient descent (Bach and Moulines, 2013) | $\frac{R^2 \ \theta_0 - \theta_*\ ^2}{n}$ | $\frac{\sigma^2 p}{n}$ |
| Accelerated gradient descent (Nesterov, 1983) | $\frac{R^2 \ \theta_0 - \theta_*\ ^2}{n^2}$ | $\sigma^2 p$ |

- **Acceleration is notoriously non-robust to noise** (d'Aspremont, 2008; Schmidt et al., 2011)
 - For non-structured noise, see Lan (2012)

Achieving optimal bias and variance terms

| | Bias | Variance |
|--|--|-----------------------------------|
| Averaged gradient descent (Bach and Moulines, 2013) | $\frac{R^2 \ \theta_0 - \theta_*\ ^2}{n}$ | $\frac{\sigma^2 p}{n}$ |
| Accelerated gradient descent (Nesterov, 1983) | $\frac{R^2 \ \theta_0 - \theta_*\ ^2}{n^2}$ | $\sigma^2 p$ |
| “Between” averaging and acceleration (Flammarion and Bach, 2015) | $\frac{R^2 \ \theta_0 - \theta_*\ ^2}{n^{1+\alpha}}$ | $\frac{\sigma^2 p}{n^{1-\alpha}}$ |

Achieving optimal bias and variance terms

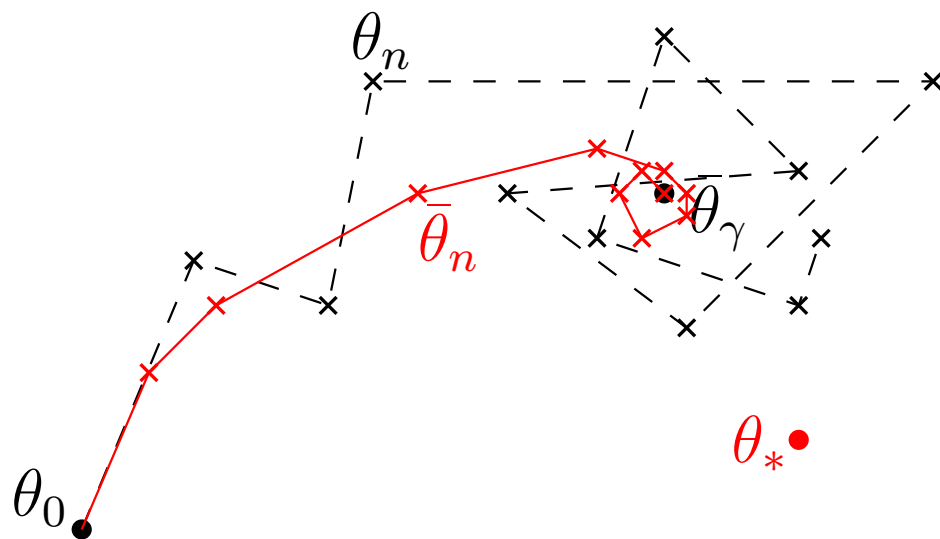
| | Bias | Variance |
|---|--|-----------------------------------|
| Averaged gradient descent (Bach and Moulines, 2013) | $\frac{R^2 \ \theta_0 - \theta_*\ ^2}{n}$ | $\frac{\sigma^2 p}{n}$ |
| Accelerated gradient descent (Nesterov, 1983) | $\frac{R^2 \ \theta_0 - \theta_*\ ^2}{n^2}$ | $\sigma^2 p$ |
| “Between” averaging and acceleration (Flammarion and Bach, 2015) | $\frac{R^2 \ \theta_0 - \theta_*\ ^2}{n^{1+\alpha}}$ | $\frac{\sigma^2 p}{n^{1-\alpha}}$ |
| Averaging and acceleration (Dieuleveut, Flammarion, and Bach, 2016) | $\frac{R^2 \ \theta_0 - \theta_*\ ^2}{n^2}$ | $\frac{\sigma^2 p}{n}$ |

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 π_γ such that $\int f'(\theta)\pi_\gamma(d\theta) = 0$
 - When f' is not linear, $f'(\int \theta\pi_\gamma(d\theta)) \neq \int f'(\theta)\pi_\gamma(d\theta) = 0$

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 π_γ such that $\int f'(\theta)\pi_\gamma(d\theta) = 0$
 - When f' is not linear, $f'(\int \theta\pi_\gamma(d\theta)) \neq \int f'(\theta)\pi_\gamma(d\theta) = 0$
- θ_n oscillates around the wrong value $\bar{\theta}_\gamma \neq \theta_*$

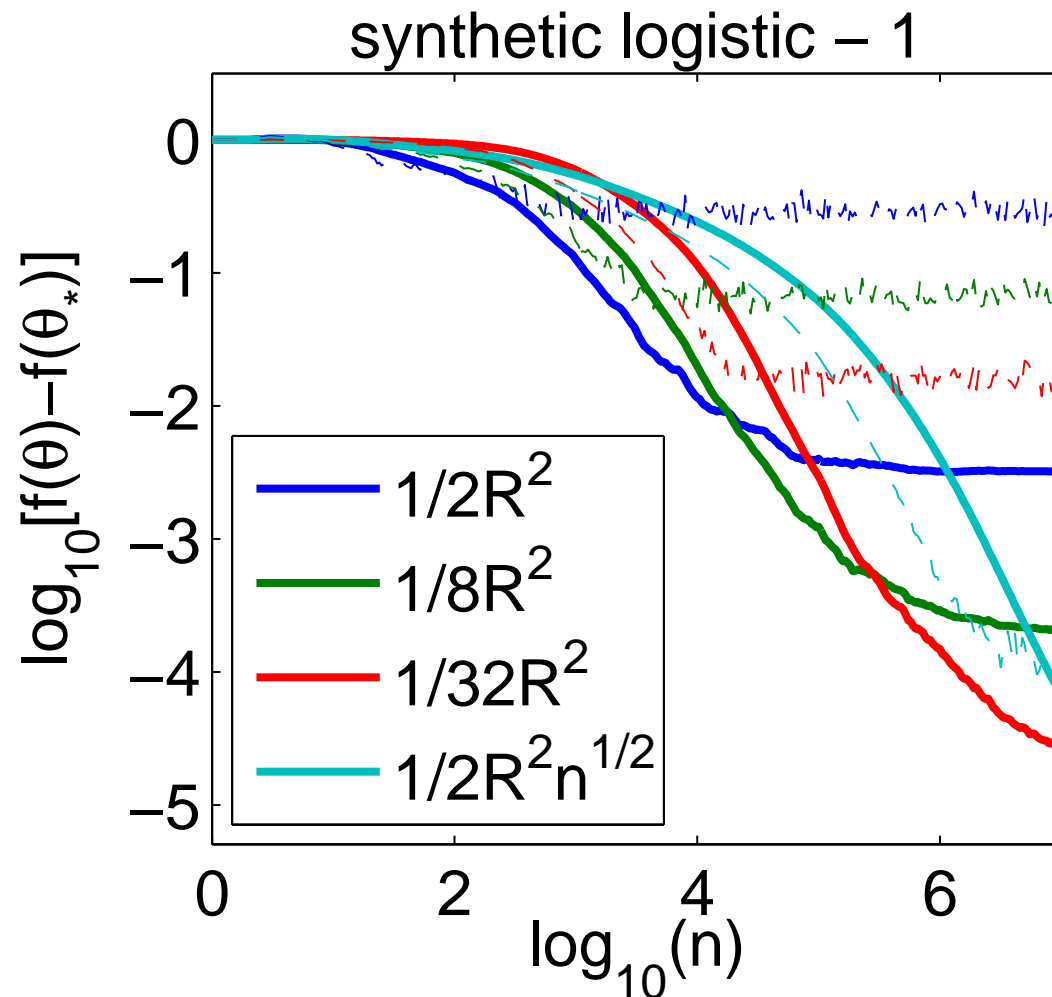


Beyond least-squares - Markov chain interpretation

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

- Gaussian distributions - $p = 20$



Restoring convergence through online Newton steps

- **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 γ_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

Restoring convergence through online Newton steps

- **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 γ_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

Restoring convergence through online Newton steps

- 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

$$\begin{aligned} g(\theta) &= f(\tilde{\theta}) + \langle f'(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \frac{1}{2} \langle \theta - \tilde{\theta}, f''(\tilde{\theta})(\theta - \tilde{\theta}) \rangle \\ &= f(\tilde{\theta}) + \langle \mathbb{E}f'_n(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \frac{1}{2} \langle \theta - \tilde{\theta}, \mathbb{E}f''_n(\tilde{\theta})(\theta - \tilde{\theta}) \rangle \\ &= \mathbb{E} \left[f(\tilde{\theta}) + \langle f'_n(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \frac{1}{2} \langle \theta - \tilde{\theta}, f''_n(\tilde{\theta})(\theta - \tilde{\theta}) \rangle \right] \end{aligned}$$

Restoring convergence through online Newton steps

- 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

$$\begin{aligned}g(\theta) &= f(\tilde{\theta}) + \langle f'(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \frac{1}{2} \langle \theta - \tilde{\theta}, f''(\tilde{\theta})(\theta - \tilde{\theta}) \rangle \\ &= f(\tilde{\theta}) + \langle \mathbb{E}f'_n(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \frac{1}{2} \langle \theta - \tilde{\theta}, \mathbb{E}f''_n(\tilde{\theta})(\theta - \tilde{\theta}) \rangle \\ &= \mathbb{E} \left[f(\tilde{\theta}) + \langle f'_n(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \frac{1}{2} \langle \theta - \tilde{\theta}, f''_n(\tilde{\theta})(\theta - \tilde{\theta}) \rangle \right]\end{aligned}$$

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

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

- $f''_n(\tilde{\theta}) = \ell''(y_n, \langle \tilde{\theta}, \Phi(x_n) \rangle) \Phi(x_n) \otimes \Phi(x_n)$ 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 $\tilde{\theta}$
- (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**

Choice of support point for online Newton step

- **Two-stage procedure**

- (1) Run $n/2$ iterations of averaged SGD to obtain $\tilde{\theta}$

- (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**

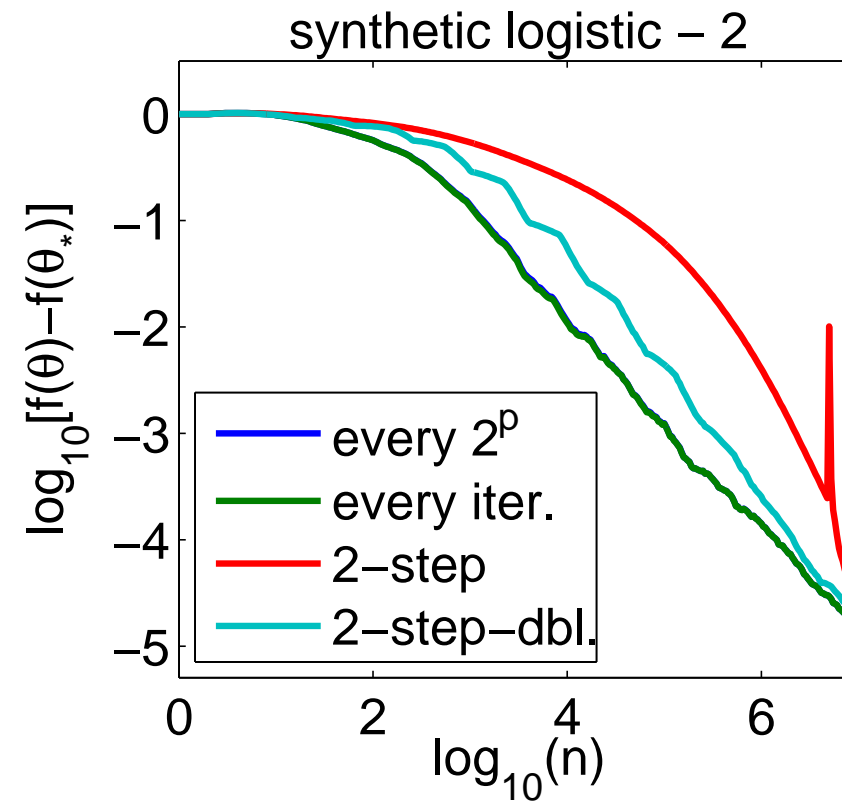
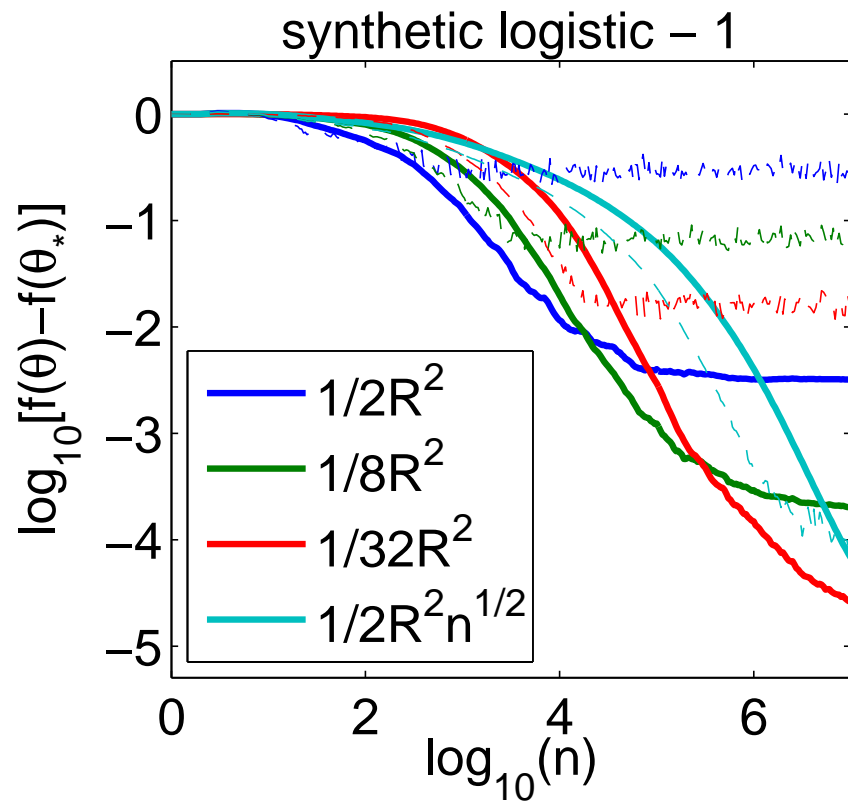
- Recursion:
$$\theta_n = \theta_{n-1} - \gamma [f'_n(\bar{\theta}_{n-1}) + f''_n(\bar{\theta}_{n-1})(\theta_{n-1} - \bar{\theta}_{n-1})]$$

- No provable convergence rate (yet) but best practical behavior

- Note (dis)similarity with regular SGD: $\theta_n = \theta_{n-1} - \gamma f'_n(\theta_{n-1})$

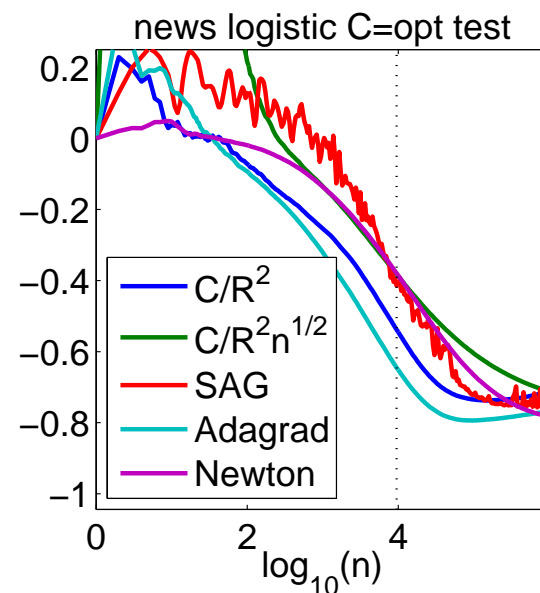
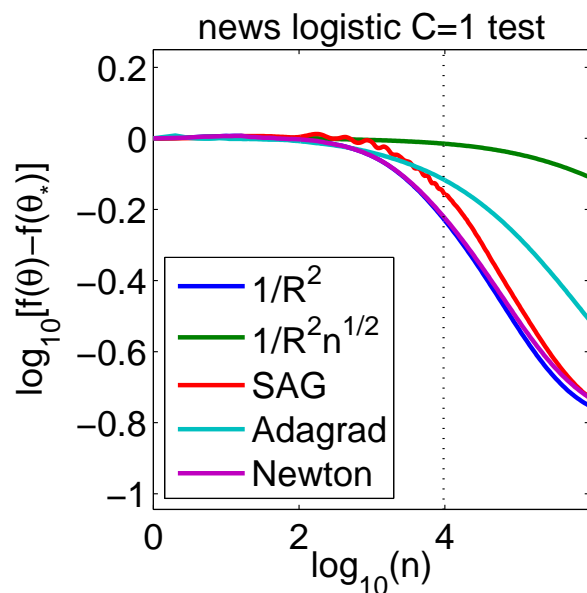
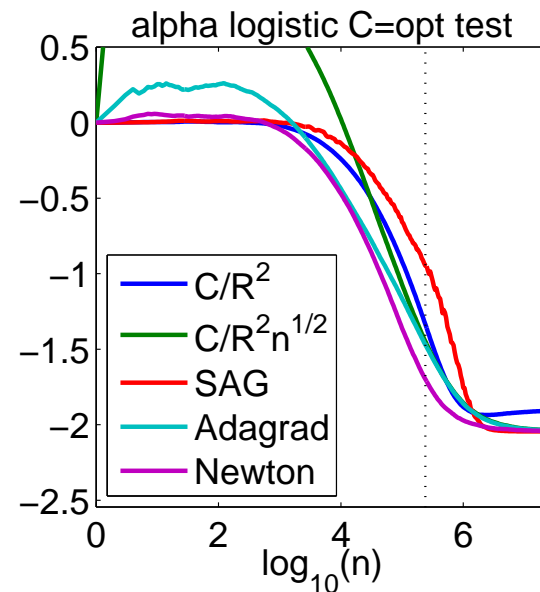
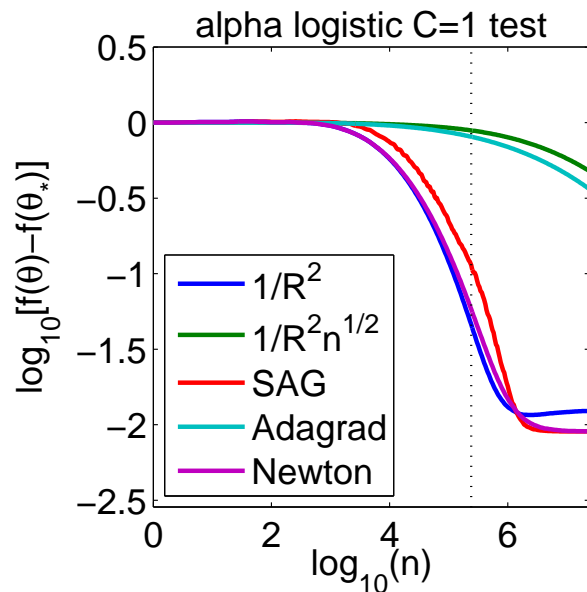
Simulations - synthetic examples

- Gaussian distributions - $p = 20$



Simulations - benchmarks

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



Conclusions

- **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 and adaptivity

Conclusions

- **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 and adaptivity
- **Extensions and future work**
 - Going beyond a single pass (Le Roux, Schmidt, and Bach, 2012; Defazio, Bach, and Lacoste-Julien, 2014)
 - Proximal extensions for non-differentiable terms
 - Kernels and nonparametric estimation (Dieuleveut and Bach, 2014)
 - **Parallelization**
 - **Non-convex problems**

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.
- 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-strongly-convex smooth stochastic approximation with convergence rate $o(1/n)$. Technical Report 00831977, HAL, 2013.
- Albert Benveniste, Michel Métivier, and Pierre Priouret. *Adaptive algorithms and stochastic approximations*. Springer Publishing Company, Incorporated, 2012.
- L. Bottou and O. Bousquet. The tradeoffs of large scale learning. In *Adv. NIPS*, 2008.
- A. d’Aspremont. Smooth optimization with approximate gradient. *SIAM J. Optim.*, 19(3):1171–1183, 2008.
- 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.
- A. Dieuleveut and F. Bach. Non-parametric Stochastic Approximation with Large Step sizes. Technical report, ArXiv, 2014.

- A. Dieuleveut, N. Flammarion, and F. Bach. Harder, better, faster, stronger convergence rates for least-squares regression. Technical Report 1602.05419, arXiv, 2016.
- N. Flammarion and F. Bach. From averaging to acceleration, there is only a step-size. *arXiv preprint arXiv:1504.01577*, 2015.
- L. Györfi and H. Walk. On the averaged stochastic approximation for linear regression. *SIAM Journal on Control and Optimization*, 34(1):31–61, 1996.
- G. Lan. An optimal method for stochastic composite optimization. *Math. Program.*, 133(1-2, Ser. A): 365–397, 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.
- O. Macchi. *Adaptive processing: The least mean squares approach with applications in transmission*. Wiley West Sussex, 1995.
- 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.
- 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. Convergence rates for inexact proximal-gradient method. In *Adv. NIPS*, 2011.

A. B. Tsybakov. Optimal rates of aggregation. 2003.

A. W. Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge Univ. press, 2000.