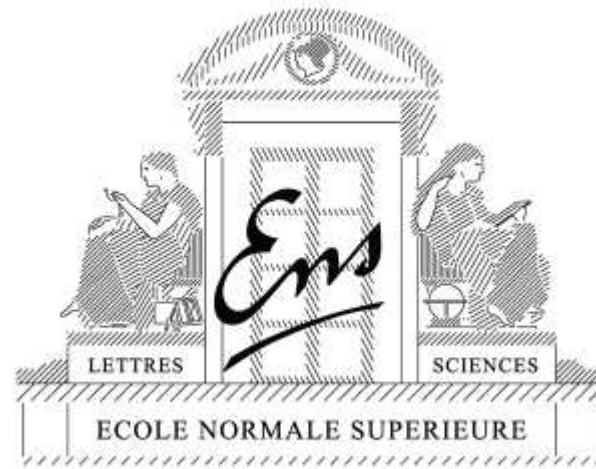


Beyond stochastic gradient descent for large-scale machine learning

Francis Bach

INRIA - Ecole Normale Supérieure, Paris, France



Joint work with Eric Moulines - September 2014

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 $\Phi(x) \in \mathbb{R}^p$
- **(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 $\hat{\theta}$

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

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

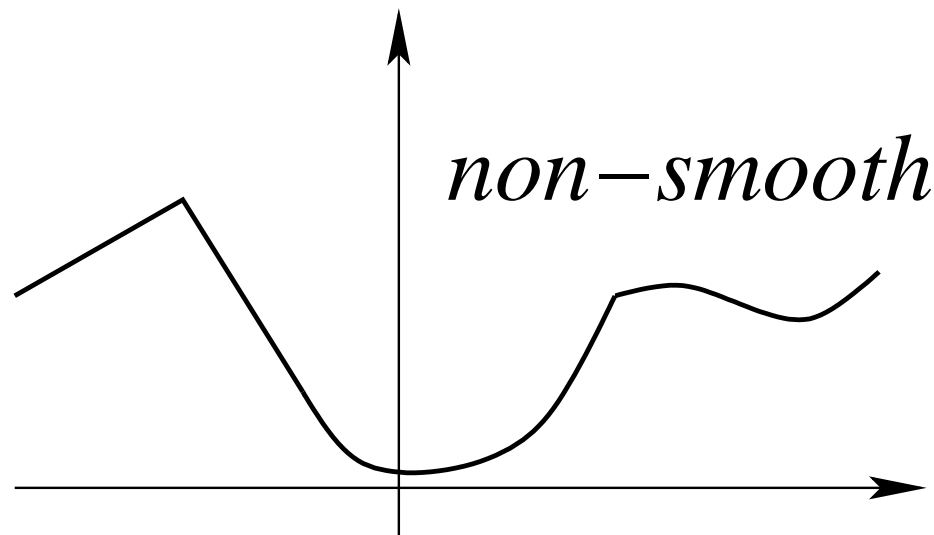
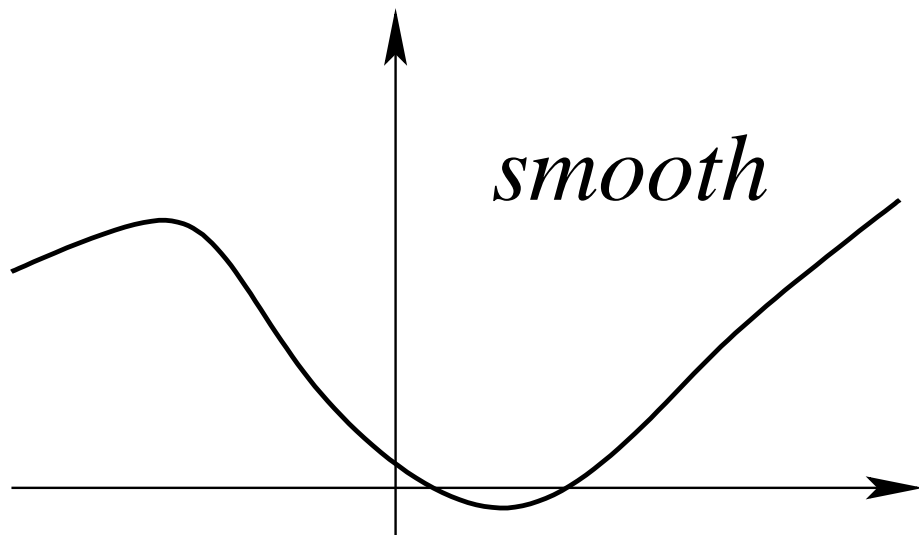
convex data fitting term + regularizer

- 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 $\hat{\theta}$
 - **May be tackled simultaneously**

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

$$\forall \theta \in \mathbb{R}^p, g''(\theta) \preceq L \cdot \text{Id}$$



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

$$\forall \theta \in \mathbb{R}^p, g''(\theta) \preceq L \cdot \text{Id}$$

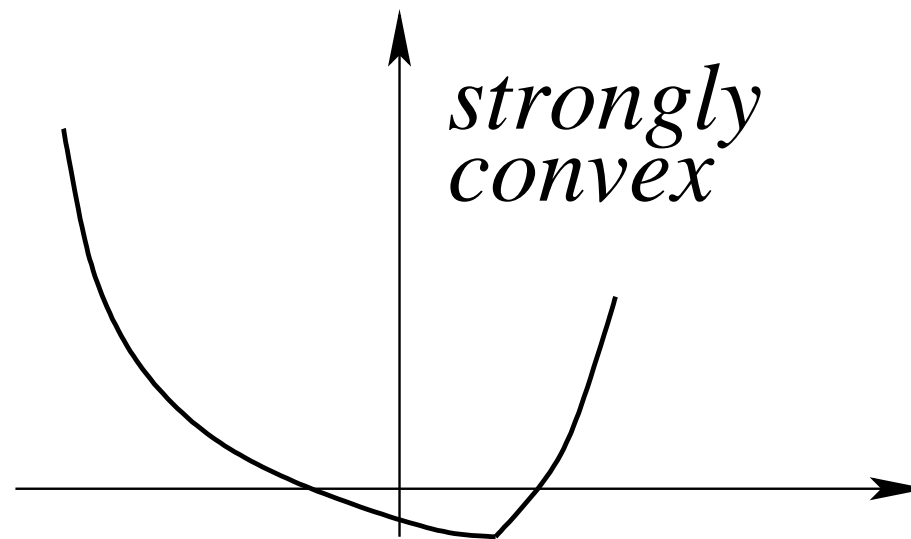
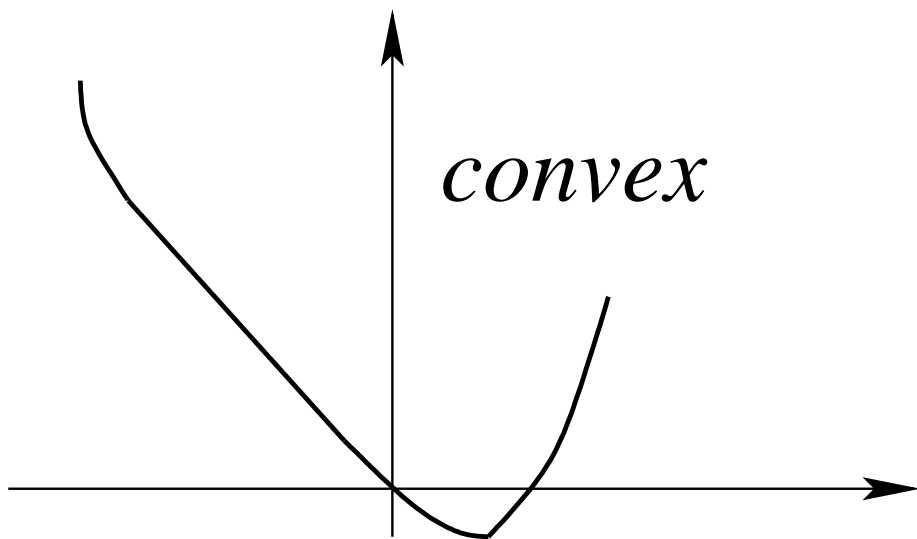
- **Machine learning**

- 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**

Smoothness and strong convexity

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

$$\forall \theta \in \mathbb{R}^p, g''(\theta) \succcurlyeq \mu \cdot \text{Id}$$



Smoothness and strong convexity

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

$$\forall \theta \in \mathbb{R}^p, g''(\theta) \succcurlyeq \mu \cdot \text{Id}$$

- **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

$$\forall \theta \in \mathbb{R}^p, g''(\theta) \succcurlyeq \mu \cdot \text{Id}$$

- **Machine learning**

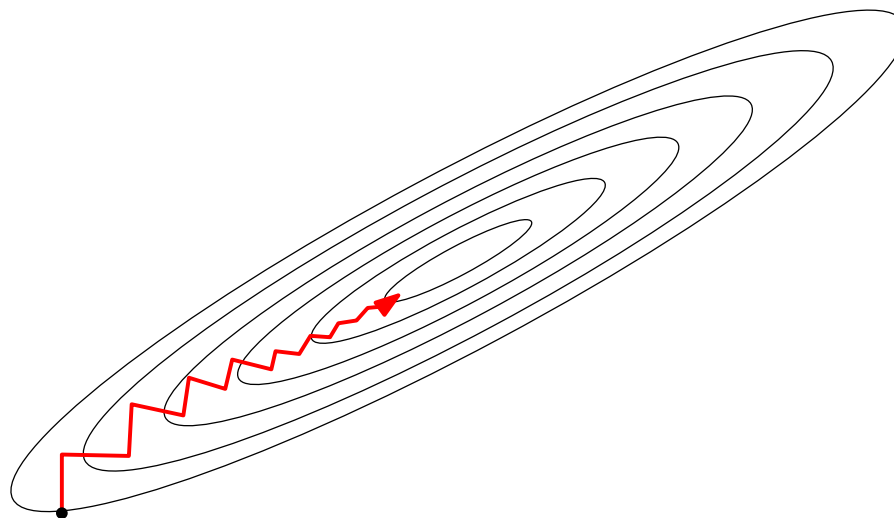
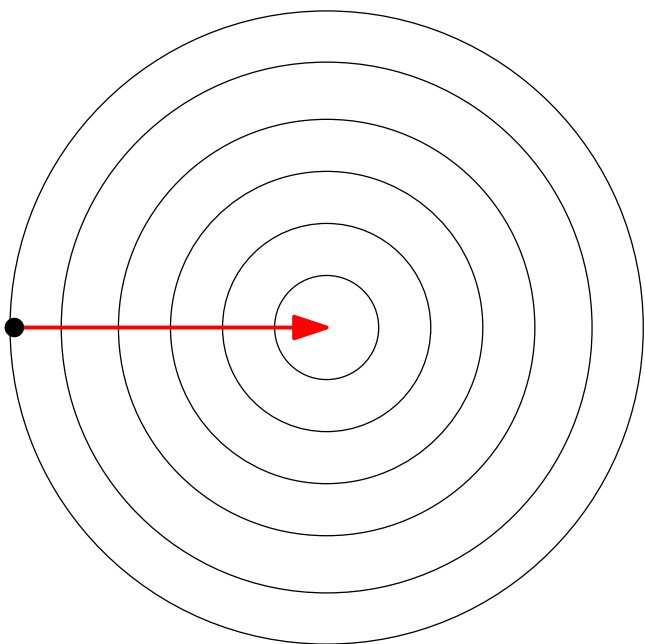
- with $g(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \langle \theta, \Phi(x_i) \rangle)$
- Hessian \approx covariance matrix $\frac{1}{n} \sum_{i=1}^n \Phi(x_i) \otimes \Phi(x_i)$
- **Data with invertible covariance matrix** (low correlation/dimension)

- **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
 - $O(e^{-\rho t})$ 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
 - $O(e^{-\rho t})$ 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
 - $O(e^{-\rho t})$ 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 work

- 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 work

- 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 work

- **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 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

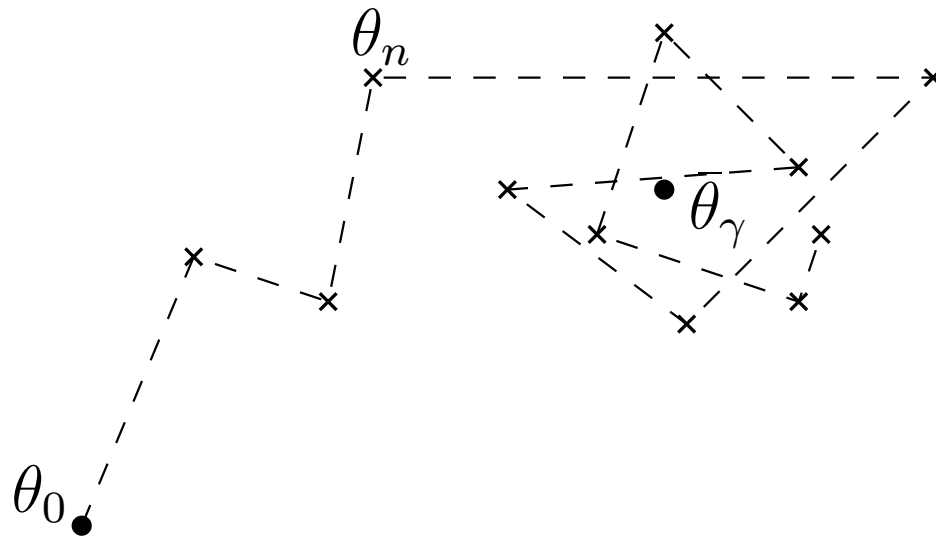
- **Least-squares:** $f(\theta) = \frac{1}{2}\mathbb{E}[(y_n - \langle \Phi(x_n), \theta \rangle)^2]$ with $\theta \in \mathbb{R}^p$
 - SGD = least-mean-square algorithm (see, e.g., Macchi, 1995)
 - usually studied without averaging and decreasing step-sizes
 - with strong convexity assumption $\mathbb{E}[\Phi(x_n)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

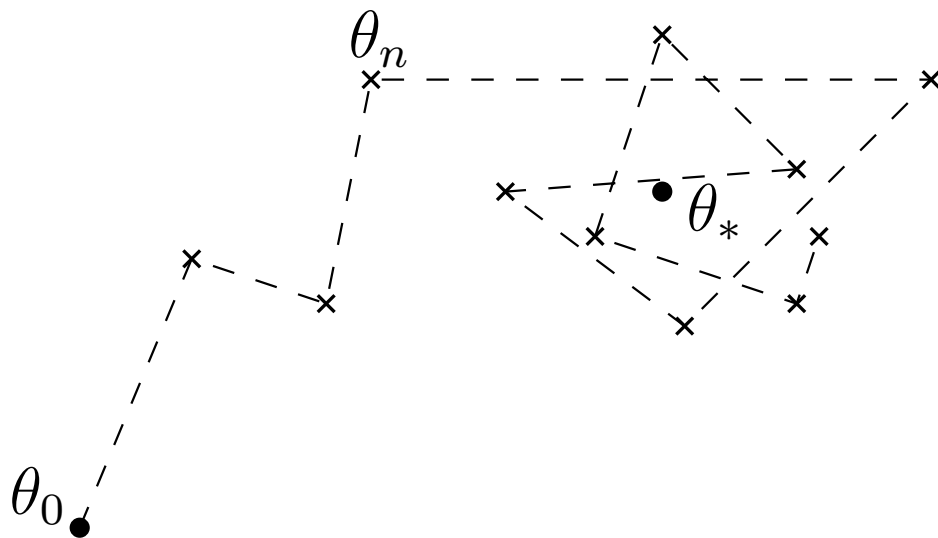
$$\theta_n = \theta_{n-1} - \gamma(\langle \Phi(x_n), \theta_{n-1} \rangle - y_n)\Phi(x_n)$$

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

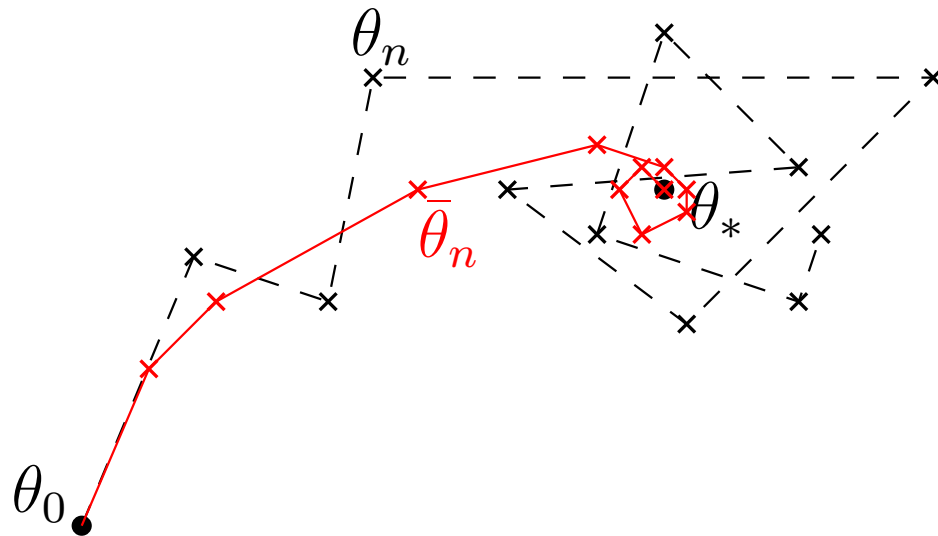
$$\theta_n = \theta_{n-1} - \gamma(\langle \Phi(x_n), \theta_{n-1} \rangle - y_n)\Phi(x_n)$$

- 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

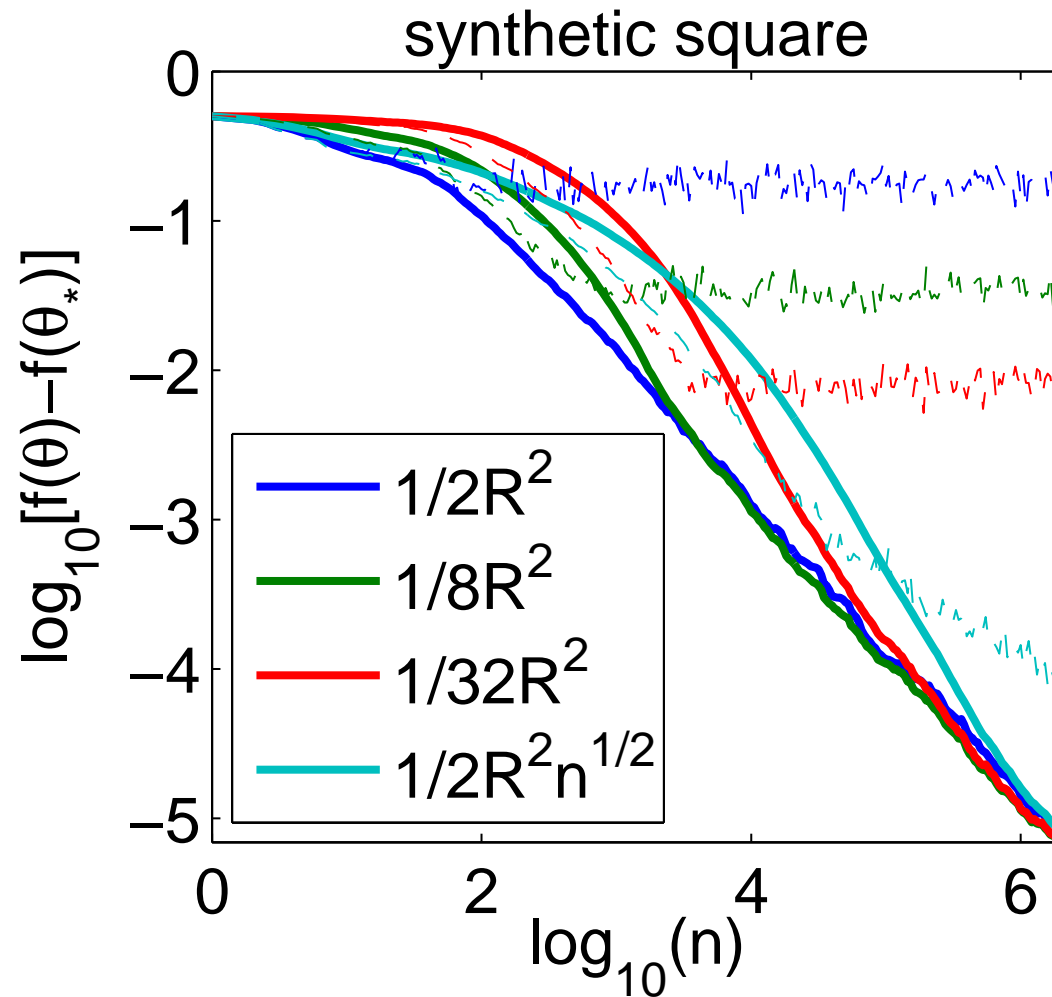
- oscillations of order $\sqrt{\gamma}$

- **Ergodic theorem:**

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

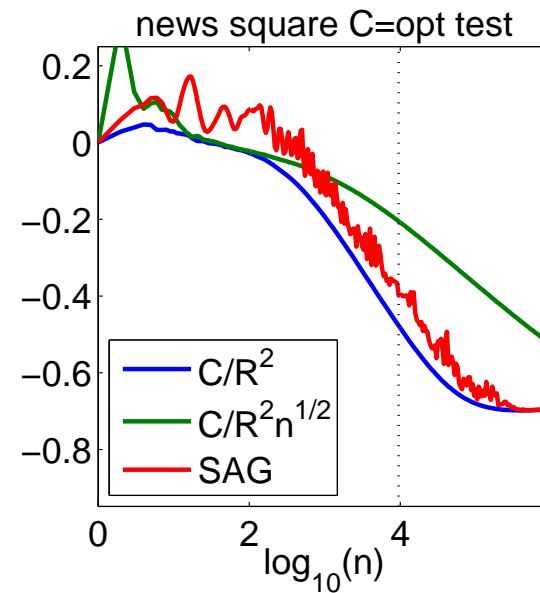
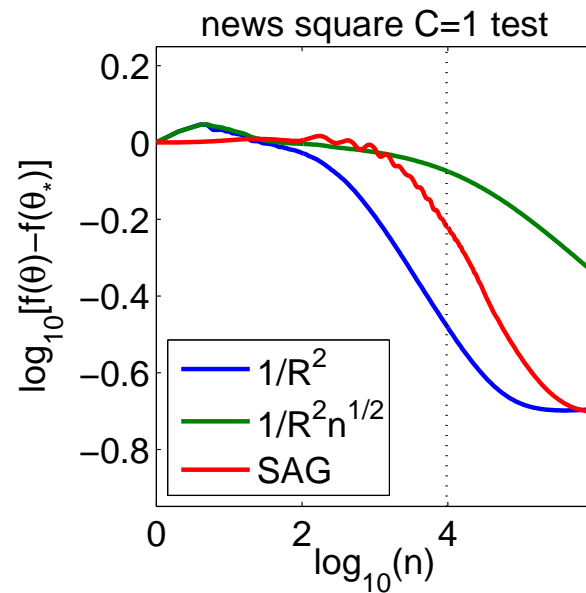
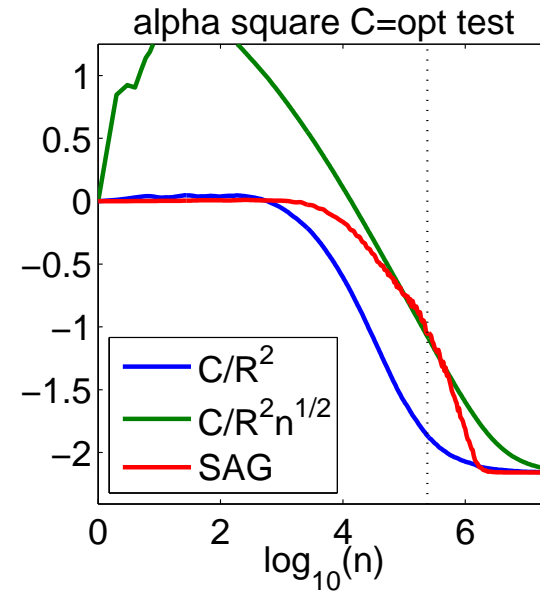
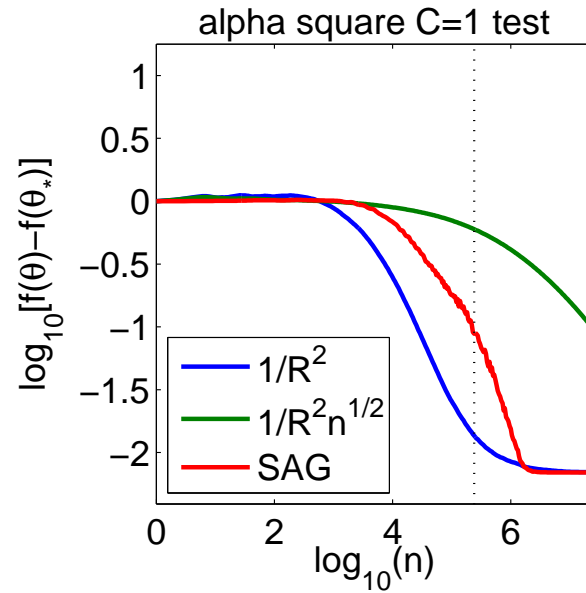
Simulations - synthetic examples

- Gaussian distributions - $p = 20$



Simulations - benchmarks

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

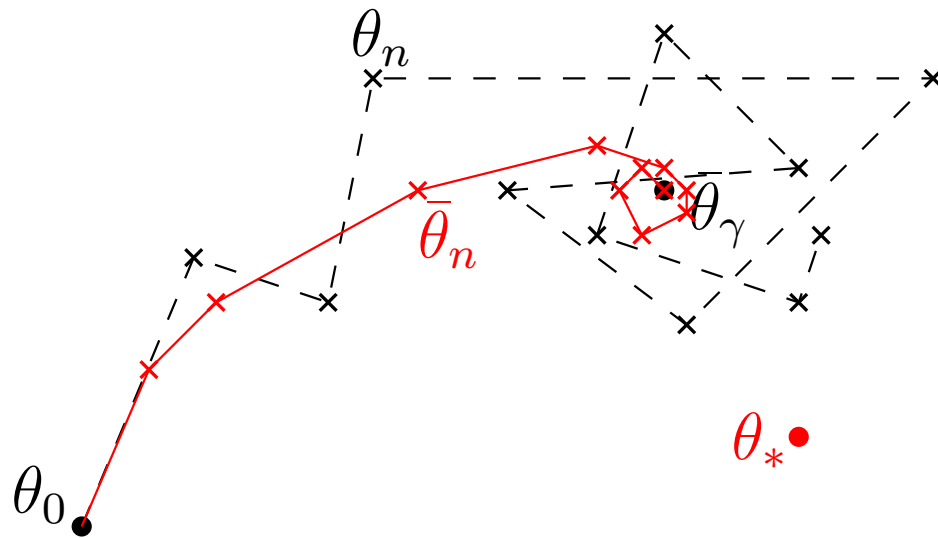


Beyond least-squares - Markov chain interpretation

- 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 $\bar{\theta}_\gamma \neq \theta_*$

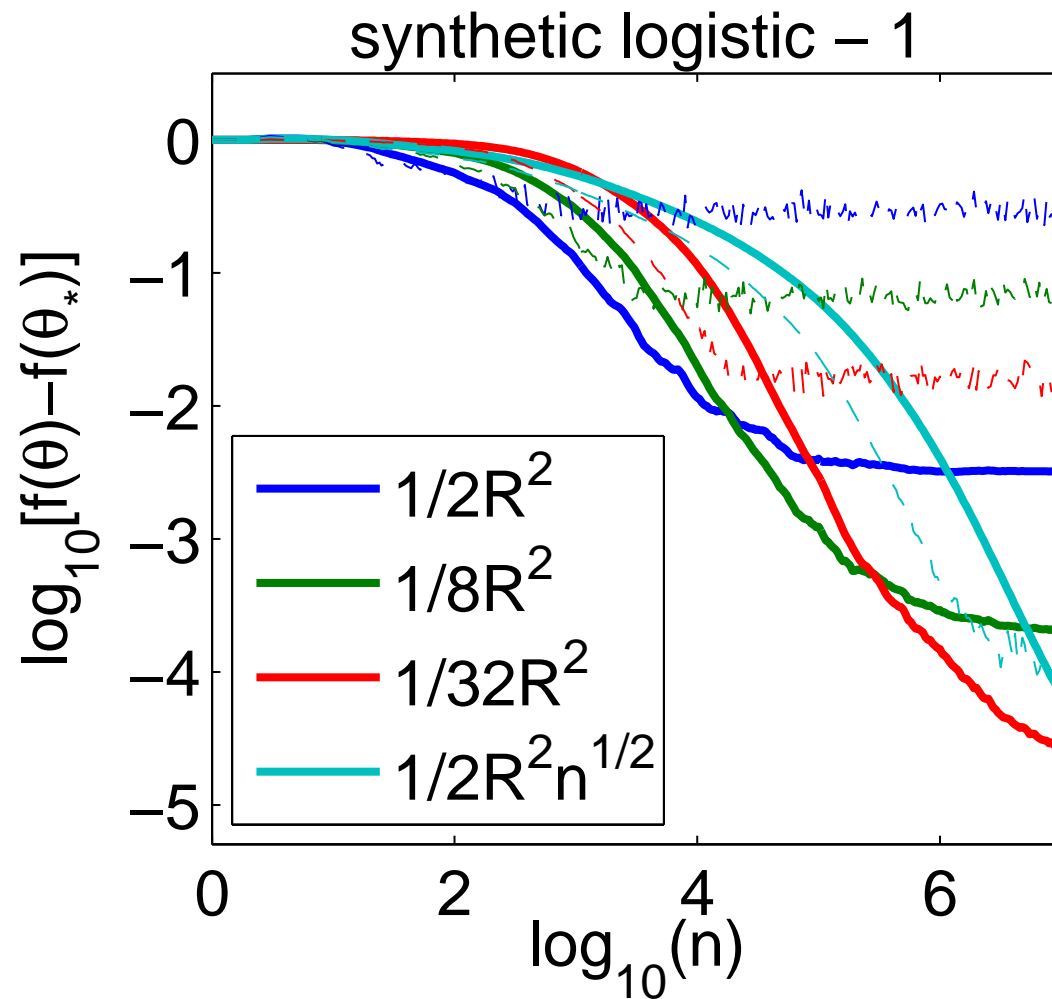


Beyond least-squares - Markov chain interpretation

- 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**

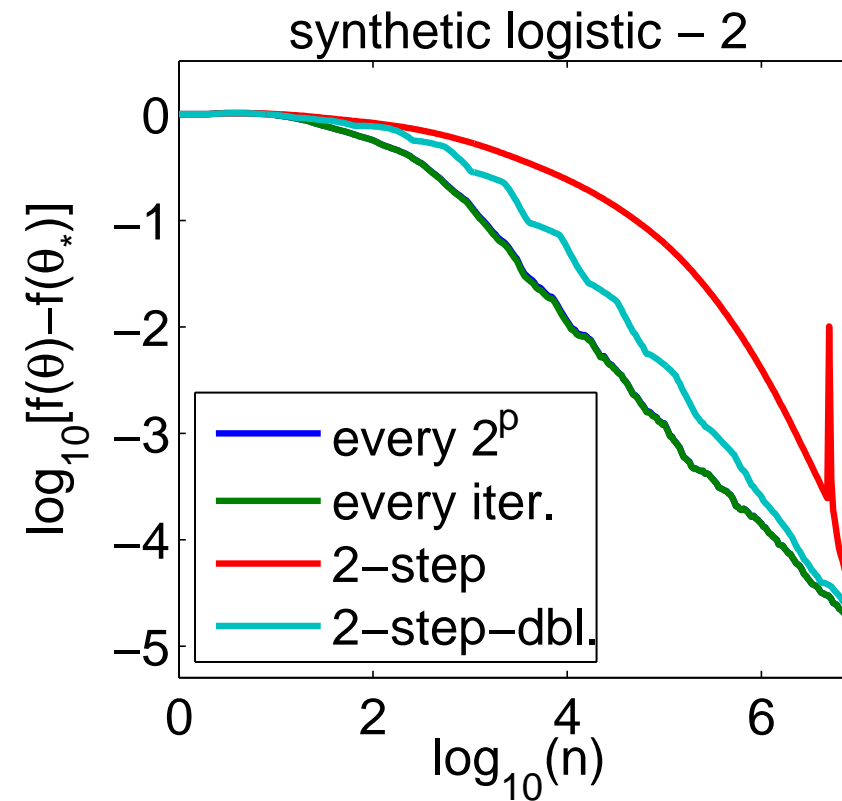
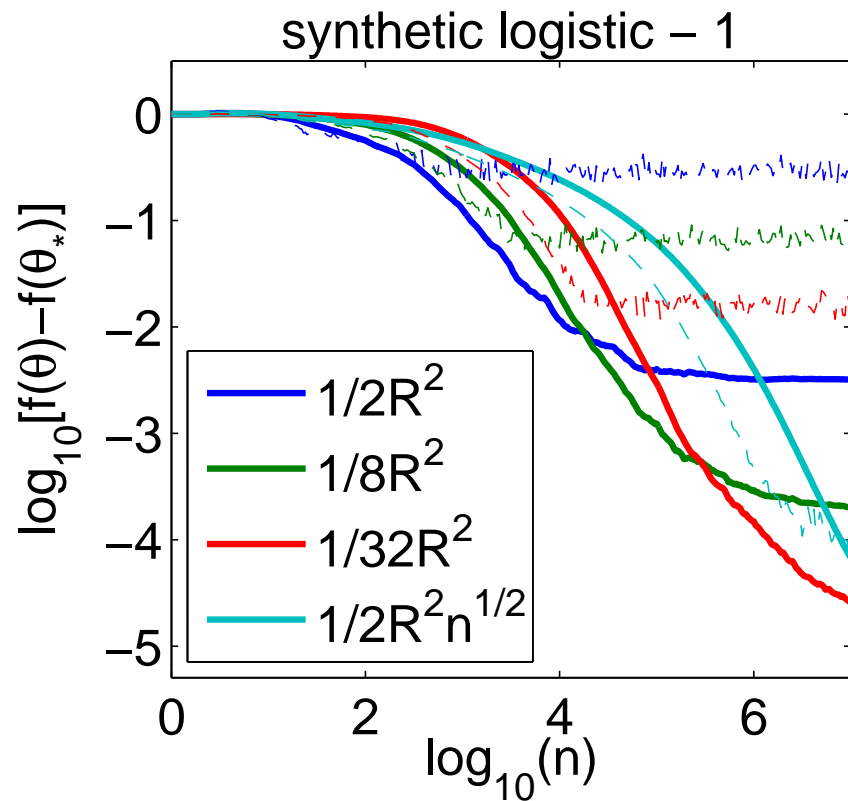
- Recursion:
$$\theta_n = \theta_{n-1} - \gamma [f'_n(\bar{\theta}_{n-1}) + f''_n(\bar{\theta}_{n-1})(\theta_{n-1} - \bar{\theta}_{n-1})]$$

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

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

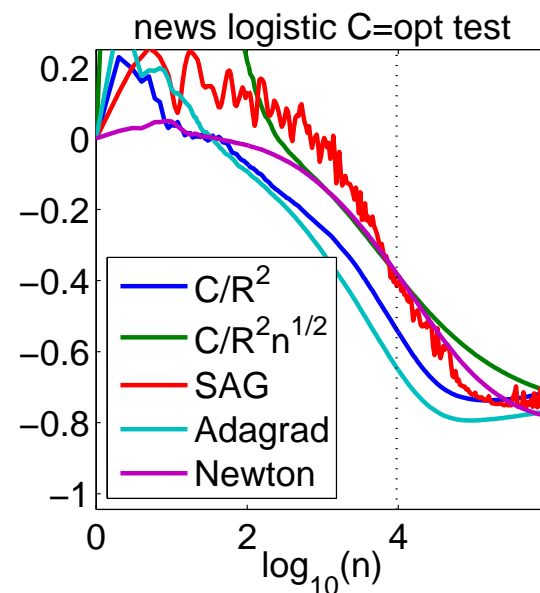
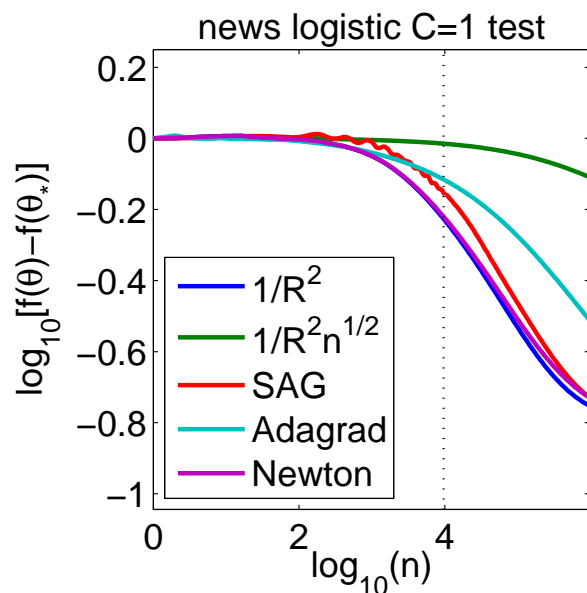
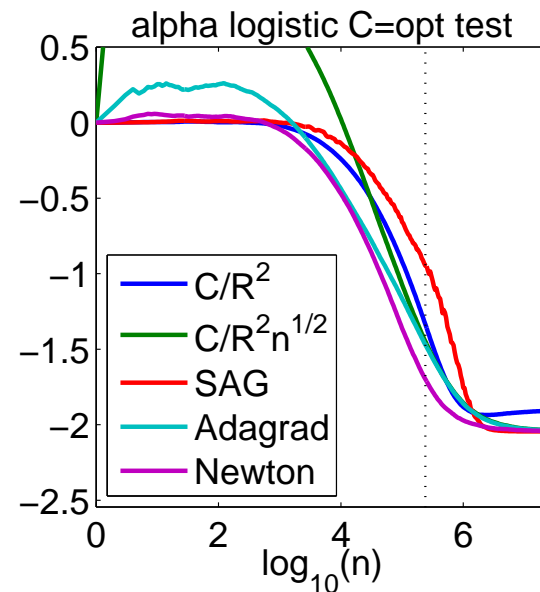
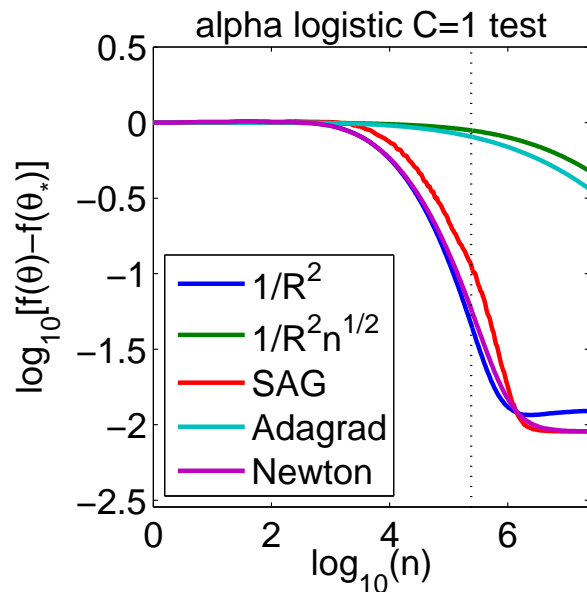
Simulations - synthetic examples

- 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

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
- **Extensions and future work**
 - Going beyond a single pass
 - Pre-conditioning
 - Proximal extensions fo non-differentiable terms
 - kernels and non-parametric estimation
 - line-search
 - 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.
- 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.
- L. Györfi and H. Walk. On the averaged stochastic approximation for linear regression. *SIAM Journal on Control and Optimization*, 34(1):31–61, 1996.
- O. Macchi. *Adaptive processing: The least mean squares approach with applications in transmission*. Wiley West Sussex, 1995.
- Julien Mairal. Optimization with first-order surrogate functions. *arXiv preprint arXiv:1305.3120*, 2013.
- A. S. Nemirovsky and D. B. Yudin. *Problem complexity and method efficiency in optimization*. Wiley & Sons, 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.
- S. Shalev-Shwartz and T. Zhang. Stochastic dual coordinate ascent methods for regularized loss minimization. Technical Report 1209.1873, Arxiv, 2012.
- A. B. Tsybakov. Optimal rates of aggregation. In *Proc. COLT*, 2003.
- A. W. Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge Univ. press, 2000.