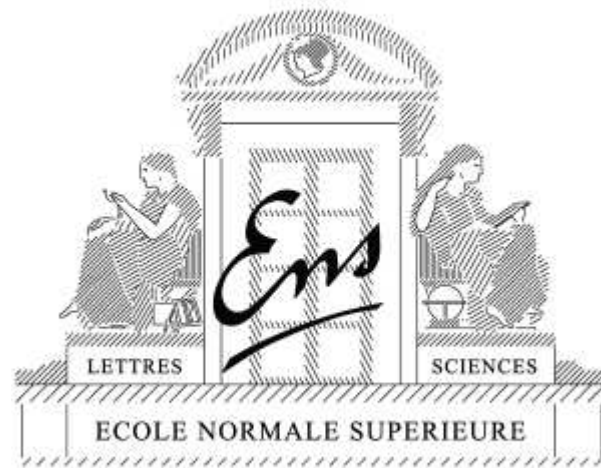# Efficient and robust stochastic approximation through an online Newton method

**Francis Bach**

*INRIA - Ecole Normale Supérieure, Paris, France*

# Context
## Large-scale supervised machine learning

- **Large $p$, large $n$, large $k$**

  - $p$ : dimension of each observation (input)
  - $n$ : number of observations
  - $k$ : number of tasks (dimension of outputs)

- **Examples**: computer vision, bioinformatics, etc.

- **Ideal running-time complexity**: $O(pn + kn)$

# Context
## Large-scale supervised machine learning

- **Large $p$, large $n$, large $k$**

  - $p$ : dimension of each observation (input)
  - $n$ : number of observations
  - $k$ : number of tasks (dimension of outputs)

- **Examples**: computer vision, bioinformatics, etc.

- **Ideal running-time complexity**: $O(pn + kn)$

- **Going back to simple methods**

  - Stochastic gradient methods (Robbins and Monro, 1951)
  - Mixing statistics and optimization

# Outline

- **Introduction**: Stochastic gradient and averaging
  - Strongly convex $O\left(\frac{1}{\mu n}\right)$ vs. non-strongly convex $O\left(\frac{1}{\sqrt{n}}\right)$

# Outline

- **Introduction**: Stochastic gradient and averaging

  - Strongly convex $O\left(\frac{1}{\mu n}\right)$ vs. non-strongly convex $O\left(\frac{1}{\sqrt{n}}\right)$

- **Adaptivity of averaging** (Bach, 2013)

  - Averaged stochastic gradient with step-sizes $\propto 1/\sqrt{n}$
  - Local strong convexity: rate of $O\left(\min\left\{\frac{1}{\mu n}, \frac{1}{\sqrt{n}}\right\}\right)$

# Outline

- **Introduction**: Stochastic gradient and averaging

  - Strongly convex $O\big(\frac{1}{\mu n}\big)$ vs. non-strongly convex $O\big(\frac{1}{\sqrt{n}}\big)$

- **Adaptivity of averaging** (Bach, 2013)

  - Averaged stochastic gradient with step-sizes $\propto 1/\sqrt{n}$
  - Local strong convexity: rate of $O\big(\min\big\{\frac{1}{\mu n}, \frac{1}{\sqrt{n}}\big\}\big)$

- **Least-squares regression** (Bach and Moulines, 2013)

  - Constant step-size averaged stochastic gradient descent
  - Convergence rate of $O(1/n)$ in all situations

# Outline

- **Introduction**: Stochastic gradient and averaging
  - Strongly convex $O\left(\frac{1}{\mu n}\right)$ vs. non-strongly convex $O\left(\frac{1}{\sqrt{n}}\right)$

- **Adaptivity of averaging** (Bach, 2013)
  - Averaged stochastic gradient with step-sizes $\propto 1/\sqrt{n}$
  - Local strong convexity: rate of $O\left(\min\left\{\frac{1}{\mu n}, \frac{1}{\sqrt{n}}\right\}\right)$

- **Least-squares regression** (Bach and Moulines, 2013)
  - Constant step-size averaged stochastic gradient descent
  - Convergence rate of $O(1/n)$ in all situations

- **Logistic regression** (Bach and Moulines, 2013)
  - Online Newton steps with linear time complexity
  - Convergence rate of $O(1/n)$ in all situations

# 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

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

convex data fitting term $+$ regularizer

# Supervised machine learning

- **Data**: $n$ observations $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$, $i = 1, \ldots, n$, **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

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

  convex data fitting term $+$  regularizer

- Empirical risk: $\hat{f}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \langle \theta, \Phi(x_i) \rangle)$  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 $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$, $i = 1, \ldots, n$, **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

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

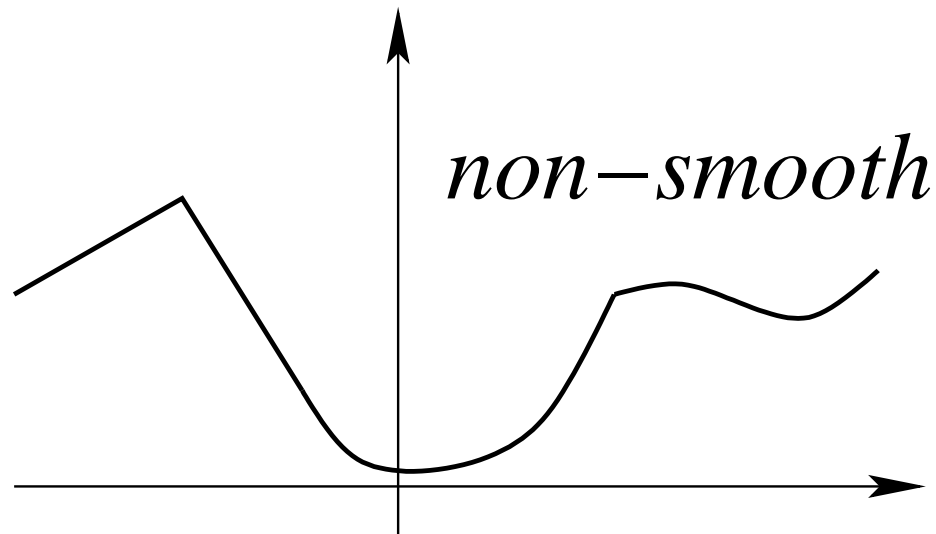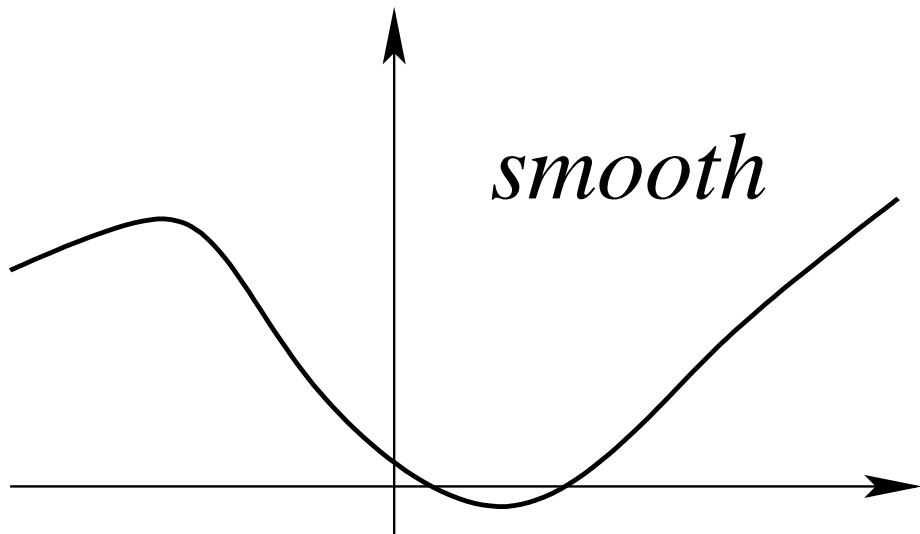<span style="color:blue">convex data fitting term</span> + <span style="color:blue">regularizer</span>

- Empirical risk: $\hat{f}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \langle \theta, \Phi(x_i) \rangle)$    <span style="color:red">training cost</span>

- Expected risk: $f(\theta) = \mathbb{E}_{(x,y)} \ell(y, \langle \theta, \Phi(x) \rangle)$    <span style="color:red">testing cost</span>

- **Two fundamental questions**: <span style="color:red">(1)</span> computing $\hat{\theta}$ and <span style="color:red">(2)</span> analyzing $\hat{\theta}$
  - **May be tackled simultaneously**

# Smoothness and strong convexity

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

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

*smooth*

*non−smooth*

# Smoothness and strong convexity

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

$$\forall \theta \in \mathbb{R}^p, \ g''(\theta) \preccurlyeq L \cdot \mathrm{Id}$$
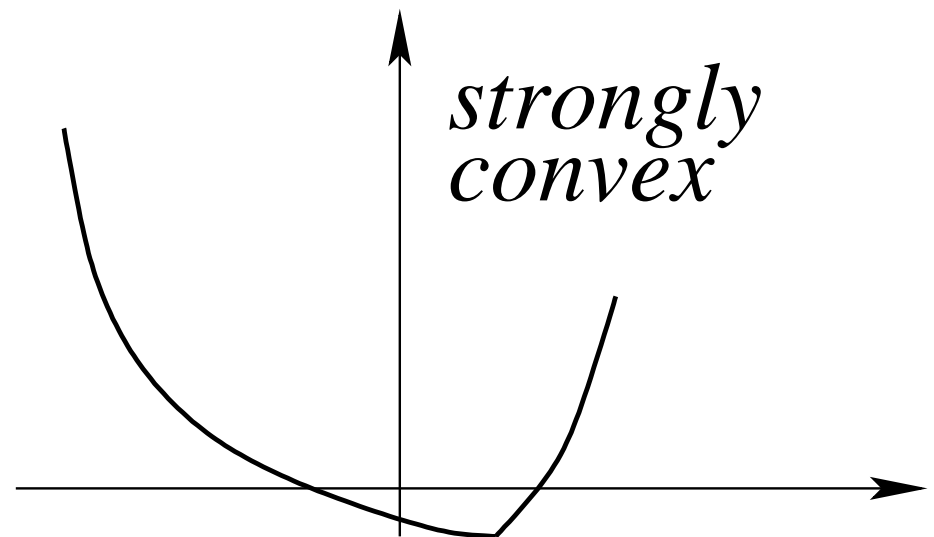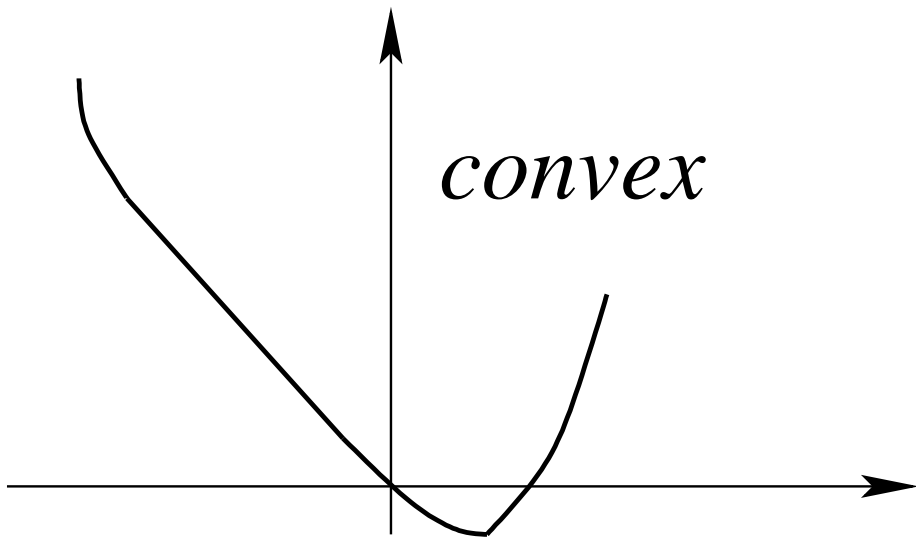
- **Machine learning**

  - 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 function $g : \mathbb{R}^p \to \mathbb{R}$ is $\mu$-strongly convex if and only if

$$\forall \theta_1, \theta_2 \in \mathbb{R}^p, \ g(\theta_1) \geqslant g(\theta_2) + \langle g'(\theta_2), \theta_1 - \theta_2 \rangle + \tfrac{\mu}{2} \|\theta_1 - \theta_2\|^2$$

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



*convex*

*strongly convex*

# Smoothness and strong convexity

- A function $g : \mathbb{R}^p \to \mathbb{R}$ is $\mu$-strongly convex if and only if

$$\forall \theta_1, \theta_2 \in \mathbb{R}^p, \ g(\theta_1) \geqslant g(\theta_2) + \langle g'(\theta_2), \theta_1 - \theta_2 \rangle + \tfrac{\mu}{2} \|\theta_1 - \theta_2\|^2$$

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

- **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 function $g : \mathbb{R}^p \to \mathbb{R}$ is $\mu$-strongly convex if and only if

$$\forall \theta_1, \theta_2 \in \mathbb{R}^p, \ g(\theta_1) \geqslant g(\theta_2) + \langle g'(\theta_2), \theta_1 - \theta_2 \rangle + \tfrac{\mu}{2} \|\theta_1 - \theta_2\|^2$$

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

- **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 $\mu$ 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

- **Newton method**: $\theta_t = \theta_{t-1} - g''(\theta_{t-1})^{-1} g'(\theta_{t-1})$

  – $O\big(e^{-\rho 2^t}\big)$ 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**: $\theta_t = \theta_{t-1} - g''(\theta_{t-1})^{-1} g'(\theta_{t-1})$

  - $O\big(e^{-\rho 2^t}\big)$ 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

$$\Rightarrow \textbf{Stochastic approximation}$$

# 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**

  - (much) broader applicability beyond convex optimization

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

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

# 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**

  - **loss for a single pair of observations**: $\boxed{f_n(\theta) = \ell(y_n, \langle \theta, \Phi(x_n) \rangle)}$

  - $f(\theta) = \mathbb{E}f_n(\theta) = \mathbb{E}\,\ell(y_n, \langle \theta, \Phi(x_n) \rangle) = $ **generalization error**
  - Expected gradient: $f'(\theta) = \mathbb{E}f_n'(\theta) = \mathbb{E}\left\{ \ell'(y_n, \langle \theta, \Phi(x_n) \rangle)\,\Phi(x_n) \right\}$
  - Non-asymptotic results

# Convex stochastic approximation

- **Key assumption**: smoothness and/or strongly convexity

- **Key algorithm:** stochastic gradient descent (a.k.a. Robbins-Monro)

$$\boxed{\theta_n = \theta_{n-1} - \gamma_n \, f_n'(\theta_{n-1})}$$

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

- Which learning rate sequence $\gamma_n$? Classical setting: $\boxed{\gamma_n = C n^{-\alpha}}$

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

- **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)

# 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 adaptive algorithm for smooth problems with convergence rate $O(\min\{1/\mu n, 1/\sqrt{n}\})$ in all situations?**
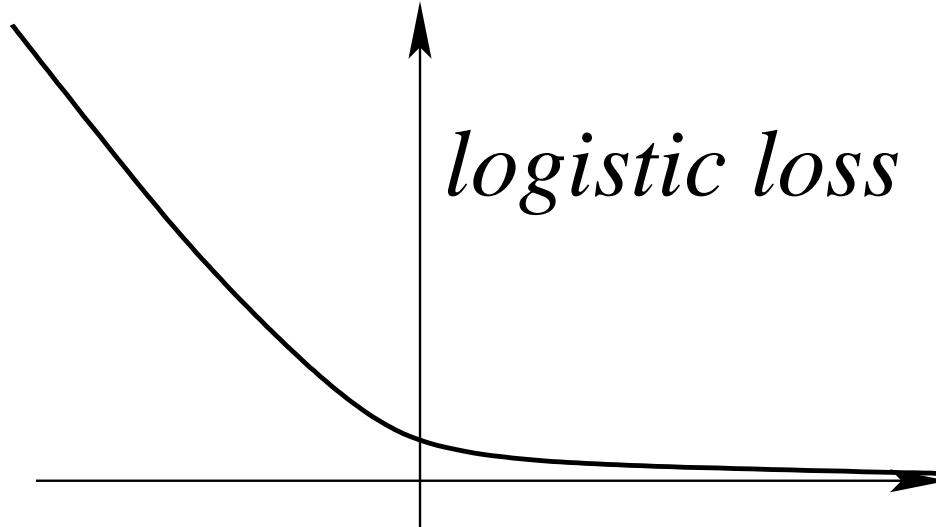
# Adaptive algorithm for logistic regression

- **Logistic regression**: $(\Phi(x_n), y_n) \in \mathbb{R}^p \times \{-1, 1\}$

  - Single data point: $f_n(\theta) = \log(1 + \exp(-y_n\langle\theta, \Phi(x_n)\rangle))$
  - Generalization error: $f(\theta) = \mathbb{E}f_n(\theta)$

# Adaptive algorithm for logistic regression

- **Logistic regression**: $(\Phi(x_n), y_n) \in \mathbb{R}^p \times \{-1, 1\}$

  – Single data point: $f_n(\theta) = \log(1 + \exp(-y_n \langle \theta, \Phi(x_n) \rangle))$
  – Generalization error: $f(\theta) = \mathbb{E} f_n(\theta)$

- **Cannot be strongly convex** $\Rightarrow$ <span style="color:red">local</span> strong convexity

  – unless restricted to $|\langle \theta, \Phi(x_n) \rangle| \leqslant M$ (and with constants $e^M$)
  – $\mu =$ lowest eigenvalue of the Hessian at the optimum $f''(\theta_*)$

*logistic loss*

# Adaptive algorithm for logistic regression

- **Logistic regression**: $(\Phi(x_n), y_n) \in \mathbb{R}^p \times \{-1, 1\}$

  – Single data point: $f_n(\theta) = \log(1 + \exp(-y_n\langle\theta, \Phi(x_n)\rangle))$
  – Generalization error: $f(\theta) = \mathbb{E}f_n(\theta)$

- **Cannot be strongly convex** $\Rightarrow$ <span style="color:red">local</span> strong convexity

  – unless restricted to $|\langle\theta, \Phi(x_n)\rangle| \leqslant M$ (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** $1/\big(2R^2\sqrt{n}\big)$

  – with $R =$ radius of data (Bach, 2013):

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

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

# Adaptive algorithm for logistic regression

- **Logistic regression**: $(\Phi(x_n), y_n) \in \mathbb{R}^p \times \{-1, 1\}$

  - Single data point: $f_n(\theta) = \log(1 + \exp(-y_n \langle \theta, \Phi(x_n) \rangle))$
  - Generalization error: $f(\theta) = \mathbb{E} f_n(\theta)$

- **Cannot be strongly convex** $\Rightarrow$ local strong convexity

  - unless restricted to $|\langle \theta, \Phi(x_n) \rangle| \leqslant M$ (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 (LMS) algorithm

- **Least-squares**: $f(\theta) = \frac{1}{2}\mathbb{E}\big[(y_n - \langle \Phi(x_n), \theta \rangle)^2\big]$ 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 $\mathbb{E}\big[\Phi(x_n) \otimes \Phi(x_n)\big] = H \succcurlyeq \mu \cdot \mathrm{Id}$

# Least-mean-square (LMS) algorithm

- **Least-squares**: $f(\theta) = \frac{1}{2}\mathbb{E}\big[(y_n - \langle \Phi(x_n), \theta \rangle)^2\big]$ 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 $\mathbb{E}\big[\Phi(x_n) \otimes \Phi(x_n)\big] = H \succcurlyeq \mu \cdot \mathrm{Id}$

- **New analysis for averaging and constant step-size** $\gamma = 1/(4R^2)$

  - Assume $\|\Phi(x_n)\| \leqslant R$ and $|y_n - \langle \Phi(x_n), \theta_* \rangle| \leqslant \sigma$ almost surely
  - No assumption regarding lowest eigenvalues of $H$

  - Main result: $\boxed{\mathbb{E}f(\bar{\theta}_n) - f(\theta_*) \leqslant \dfrac{4\sigma^2 p}{n} + \dfrac{2R^2\|\theta_0 - \theta_*\|^2}{n}}$

- **Matches statistical lower bound** (Tsybakov, 2003)

  - Non-asymptotic robust version of Györfi and Walk (1996)

# Least-mean-square (LMS) algorithm

- **Least-squares**: $f(\theta) = \frac{1}{2}\mathbb{E}\big[(y_n - \langle \Phi(x_n), \theta \rangle)^2\big]$ 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 $\mathbb{E}\big[\Phi(x_n) \otimes \Phi(x_n)\big] = H \succcurlyeq \mu \cdot \mathrm{Id}$

- **New analysis for averaging and constant step-size** $\gamma = 1/(4R^2)$

  - Assume $\|\Phi(x_n)\| \leqslant R$ and $|y_n - \langle \Phi(x_n), \theta_* \rangle| \leqslant \sigma$ almost surely
  - No assumption regarding lowest eigenvalues of $H$

  - Main result: $\boxed{\mathbb{E}f(\bar{\theta}_n) - f(\theta_*) \leqslant \dfrac{4\sigma^2 p}{n} + \dfrac{2R^2\|\theta_0 - \theta_*\|^2}{n}}$

- **Improvement of bias term** (Flammarion and Bach, 2014):
$$\min\left\{ \frac{R^2\|\theta_0 - \theta_*\|^2}{n}, \frac{R^4\langle \theta_0 - \theta_*, H^{-1}(\theta_0 - \theta_*)\rangle}{n^2} \right\}$$

# Least-mean-square (LMS) algorithm

- **Least-squares**: $f(\theta) = \frac{1}{2}\mathbb{E}\big[(y_n - \langle\Phi(x_n), \theta\rangle)^2\big]$ 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 $\mathbb{E}\big[\Phi(x_n) \otimes \Phi(x_n)\big] = H \succcurlyeq \mu \cdot \mathrm{Id}$

- **New analysis for averaging and constant step-size** $\gamma = 1/(4R^2)$

  – Assume $\|\Phi(x_n)\| \leqslant R$ and $|y_n - \langle\Phi(x_n), \theta_*\rangle| \leqslant \sigma$ almost surely
  – <span style="color:red">No assumption regarding lowest eigenvalues of $H$</span>

  – Main result: $\boxed{\mathbb{E}f(\bar{\theta}_n) - f(\theta_*) \leqslant \dfrac{4\sigma^2 p}{n} + \dfrac{2R^2\|\theta_0 - \theta_*\|^2}{n}}$

- **Extension to Hilbert spaces** (Dieuleveult and Bach, 2014):

  – Achieves minimax statistical rates given decay of spectrum of $H$

# Least-squares - Proof technique

- LMS recursion with $\varepsilon_n = y_n - \langle \Phi(x_n), \theta_* \rangle$ :

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

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

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

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

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

  – Exact computations

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

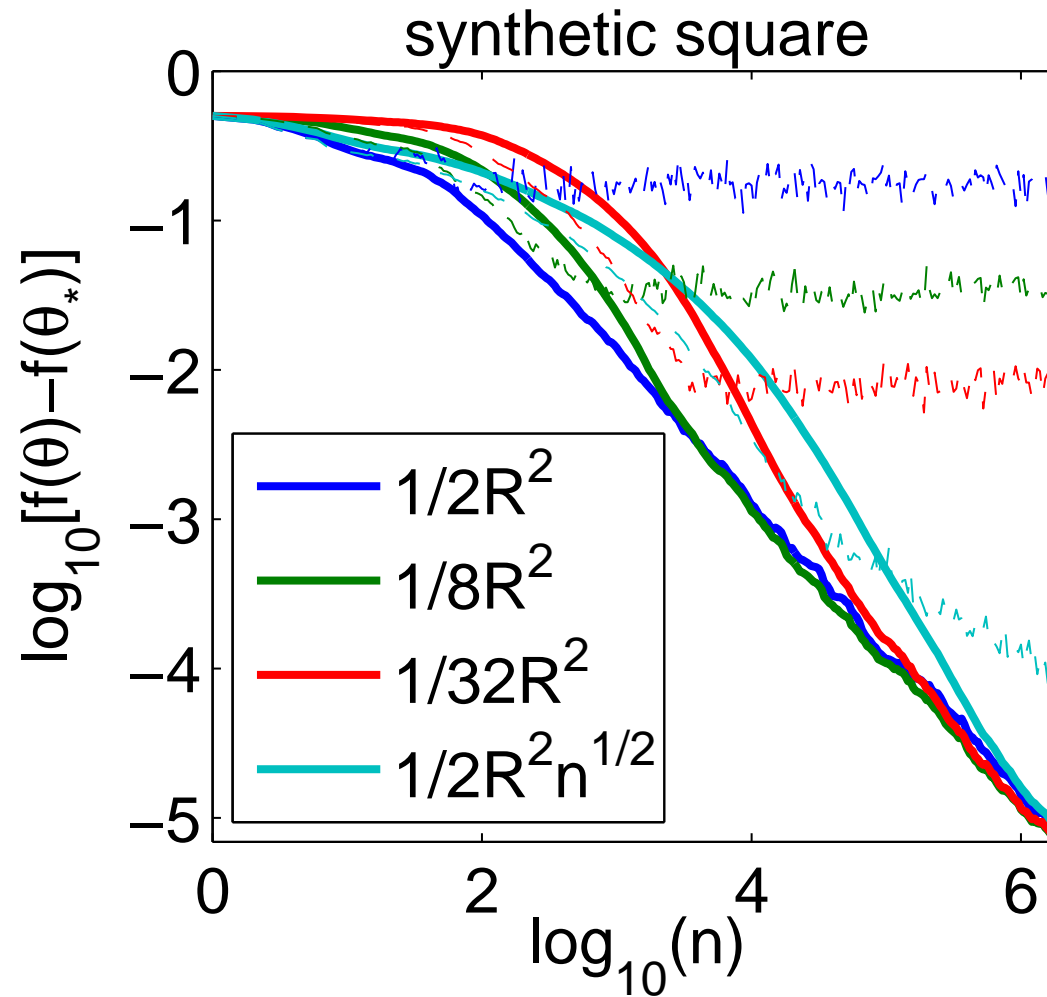# Markov chain interpretation of constant step sizes

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

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

- The sequence $(\theta_n)_n$ is a <span style="color:red">homogeneous Markov chain</span>

  - convergence to a stationary distribution $\pi_\gamma$
  - with expectation $\bar{\theta}_\gamma \overset{\text{def}}{=} \int \theta \pi_\gamma(\mathrm{d}\theta)$

# Markov chain interpretation of constant step sizes

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

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

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

  – $\theta_n$ does not converge to $\theta_*$ but oscillates around it
  – oscillations of order $\sqrt{\gamma}$
  – cf. Kaczmarz method (Strohmer and Vershynin, 2009)

- **Ergodic theorem:**

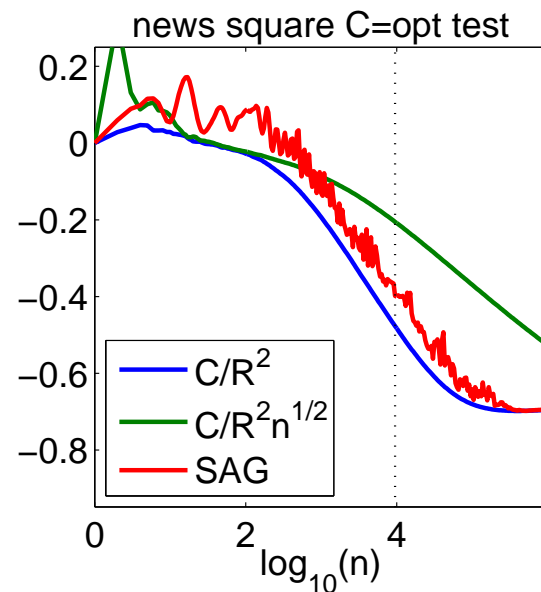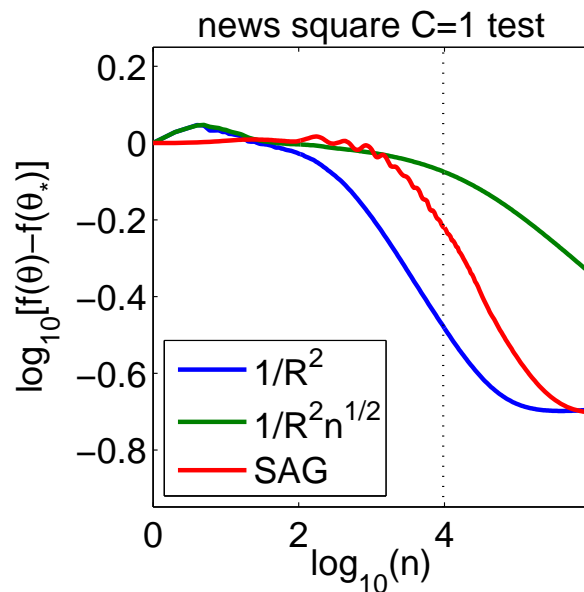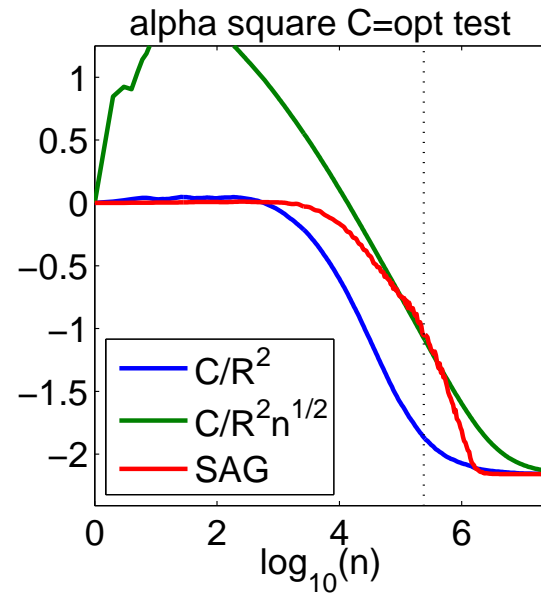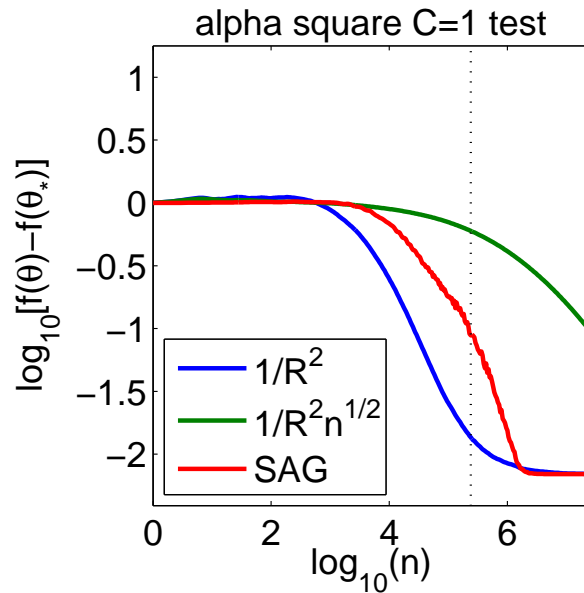  – Averaged iterates converge to $\bar{\theta}_\gamma = \theta_*$ at rate $O(1/n)$

# Simulations - synthetic examples

- Gaussian distributions - $p = 20$



synthetic square

- $1/2R^2$
- $1/8R^2$
- $1/32R^2$
- $1/2R^2n^{1/2}$

x-axis: $\log_{10}(n)$

y-axis: $\log_{10}[f(\theta) - f(\theta_*)]$

# Simulations - benchmarks

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

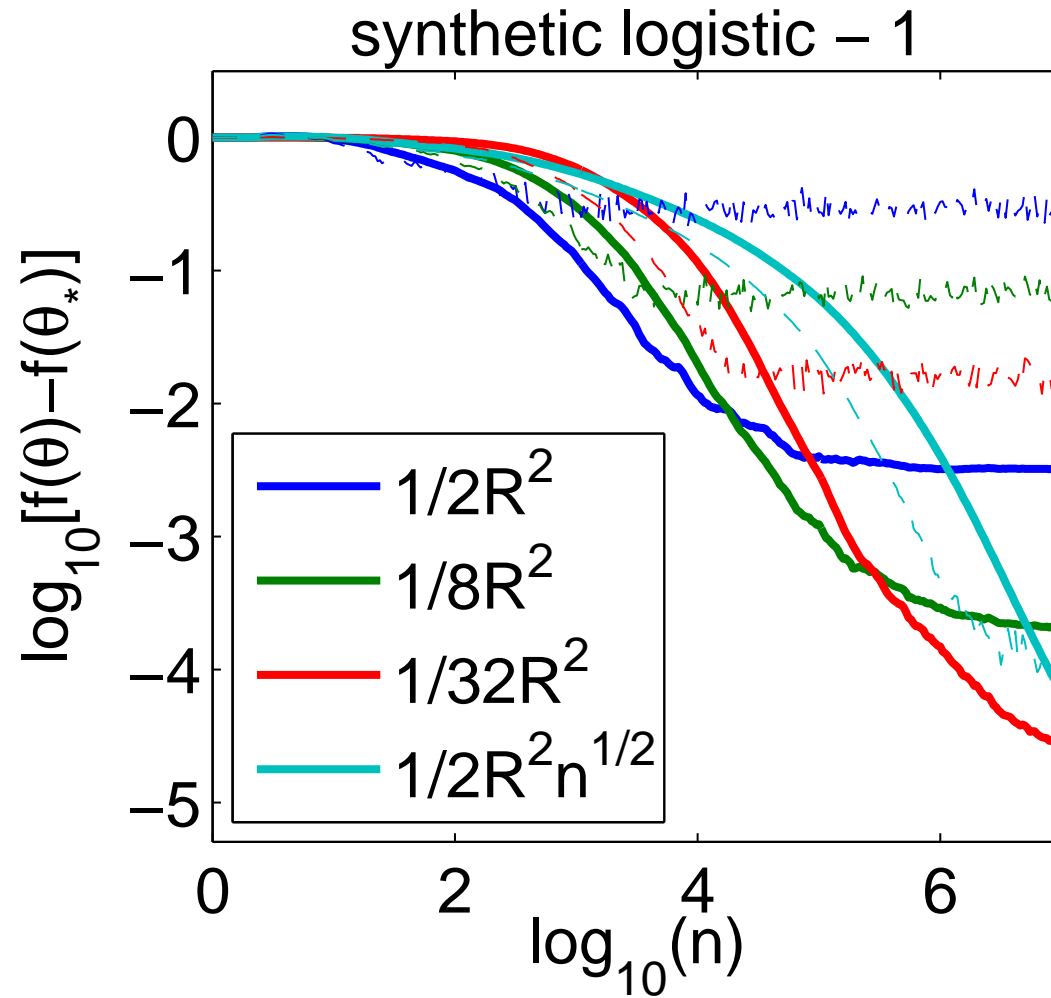# Beyond least-squares - Markov chain interpretation

- Recursion $\theta_n = \theta_{n-1} - \gamma f'_n(\theta_{n-1})$ also defines a Markov chain

  – Stationary distribution $\pi_\gamma$ 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 $\int f'(\theta)\pi_\gamma(\mathrm{d}\theta) = 0$
  - When $f'$ is not linear, $f'(\int \theta\pi_\gamma(\mathrm{d}\theta)) \neq \int f'(\theta)\pi_\gamma(\mathrm{d}\theta) = 0$

- **$\theta_n$ 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)

- NB: coherent with earlier results by Nedic and Bertsekas (2000)

# Simulations - synthetic examples

- Gaussian distributions - $p = 20$



synthetic logistic – 1

Legend:
- $1/2R^2$
- $1/8R^2$
- $1/32R^2$
- $1/2R^2n^{1/2}$

y-axis: $\log_{10}[f(\theta) - f(\theta_*)]$

x-axis: $\log_{10}(n)$

# 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 $\gamma_n$ 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

# 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 $\gamma_n$ 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

- **Online Newton step**

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

# Restoring convergence through online Newton steps

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

$$g(\theta) = f(\tilde{\theta}) + \langle f'(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \tfrac{1}{2}\langle \theta - \tilde{\theta}, f''(\tilde{\theta})(\theta - \tilde{\theta}) \rangle$$

$$= f(\tilde{\theta}) + \langle \mathbb{E} f_n'(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \tfrac{1}{2}\langle \theta - \tilde{\theta}, \mathbb{E} f_n''(\tilde{\theta})(\theta - \tilde{\theta}) \rangle$$

$$= \mathbb{E}\left[ f(\tilde{\theta}) + \langle f_n'(\tilde{\theta}), \theta - \tilde{\theta} \rangle + \tfrac{1}{2}\langle \theta - \tilde{\theta}, f_n''(\tilde{\theta})(\theta - \tilde{\theta}) \rangle \right]$$

# Restoring convergence through online Newton steps

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

$$g(\theta) = f(\tilde{\theta}) + \langle f'(\tilde{\theta}), \theta - \tilde{\theta}\rangle + \tfrac{1}{2}\langle\theta - \tilde{\theta}, f''(\tilde{\theta})(\theta - \tilde{\theta})\rangle$$

$$= f(\tilde{\theta}) + \langle\mathbb{E}f_n'(\tilde{\theta}), \theta - \tilde{\theta}\rangle + \tfrac{1}{2}\langle\theta - \tilde{\theta}, \mathbb{E}f_n''(\tilde{\theta})(\theta - \tilde{\theta})\rangle$$

$$= \mathbb{E}\Big[f(\tilde{\theta}) + \langle f_n'(\tilde{\theta}), \theta - \tilde{\theta}\rangle + \tfrac{1}{2}\langle\theta - \tilde{\theta}, f_n''(\tilde{\theta})(\theta - \tilde{\theta})\rangle\Big]$$

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

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

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

# 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: $|\varphi'''(u)| \leqslant 2\varphi''(u)^{3/2}$

- 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 $\tilde{\theta}$
(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

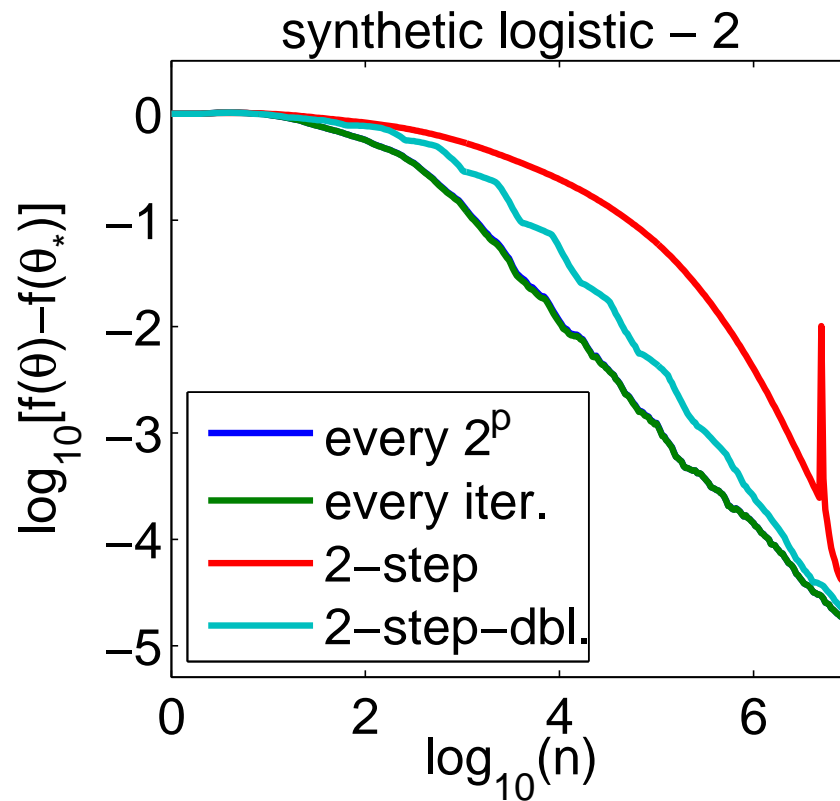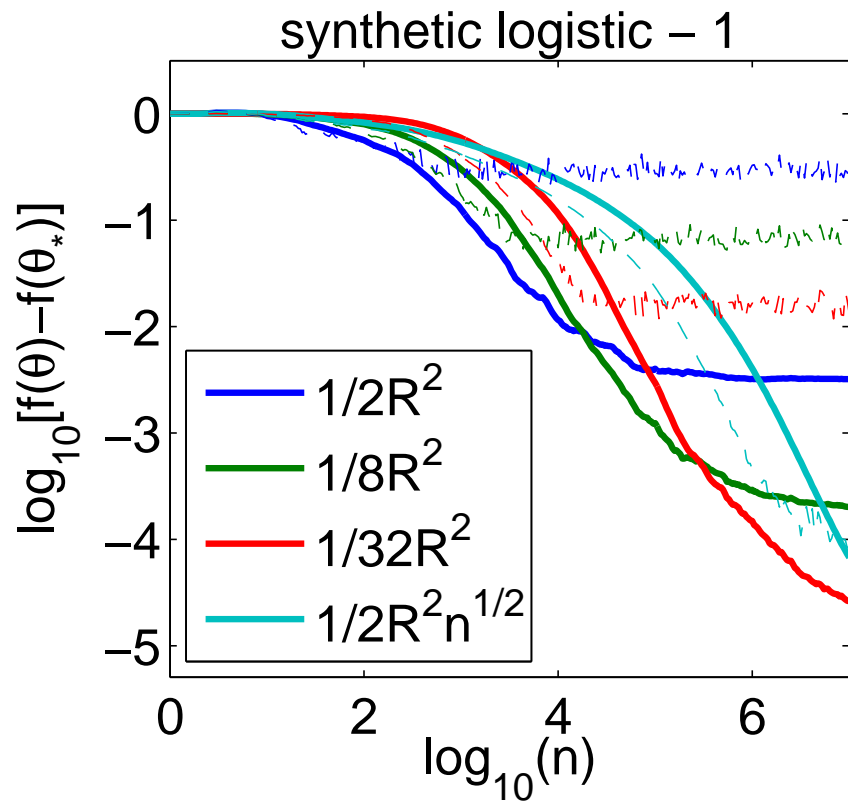- **Update at each iteration using the current averaged iterate**

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

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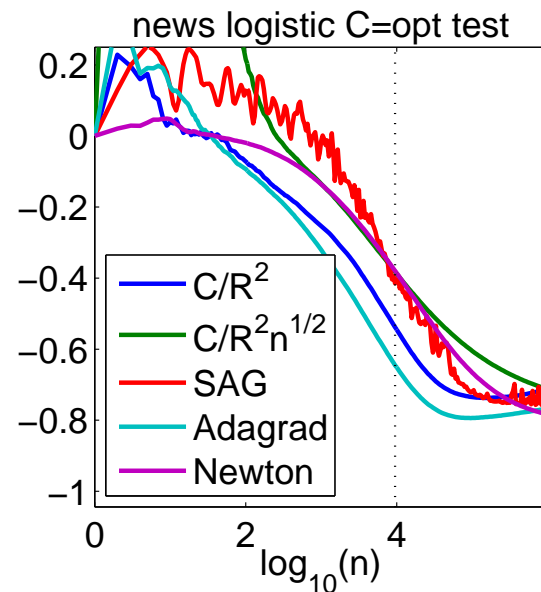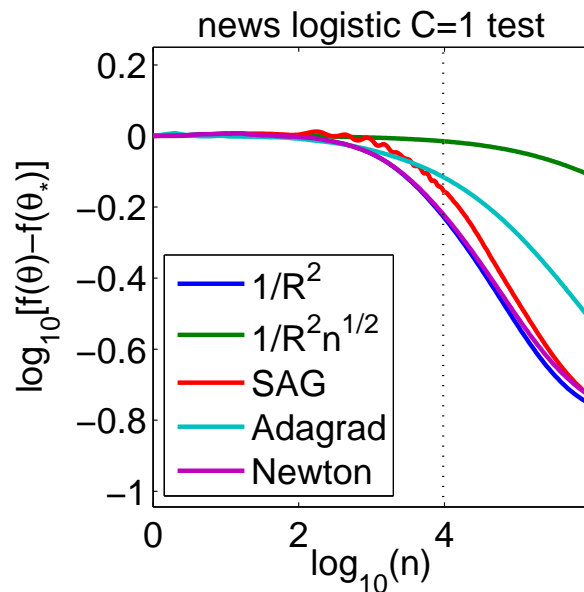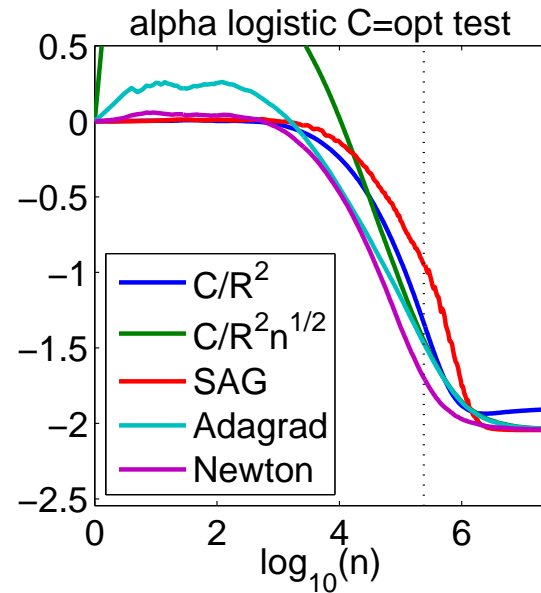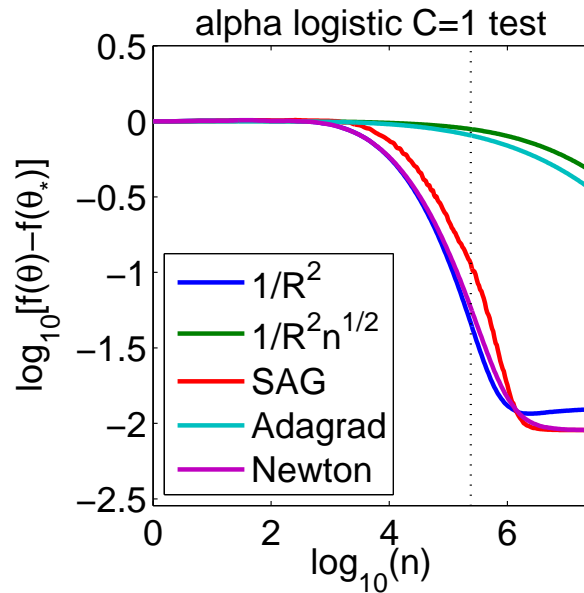
# Simulations - synthetic examples

- 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.

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-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.

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.

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.

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.

H. J. Kushner and G. G. Yin. *Stochastic approximation and recursive algorithms and applications*. Springer-Verlag, second edition, 2003.

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 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.

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 N. Srebro. SVM optimization: inverse dependence on training set size. In *Proc. ICML*, 2008.

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.

Thomas Strohmer and Roman Vershynin. A randomized kaczmarz algorithm with exponential convergence. *Journal of Fourier Analysis and Applications*, 15(2):262–278, 2009.

A. B. Tsybakov. Optimal rates of aggregation. In *Proc. COLT*, 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.