

Optimisation algorithms in Statistics II, lecture 1

Frank Miller, Department of Statistics; Stockholm University

March 23, 2021

Course schedule

- Topic 1: **Stochastic gradient descent and quasi-Newton algorithms**

Lectures: March 23; Time 10-12, 13-15 (online, Zoom)

- Topic 2: Lectures: April 13; Time 10-12, 13-15
- Topic 3: Lectures: April 27; Time 10-12, 13-15

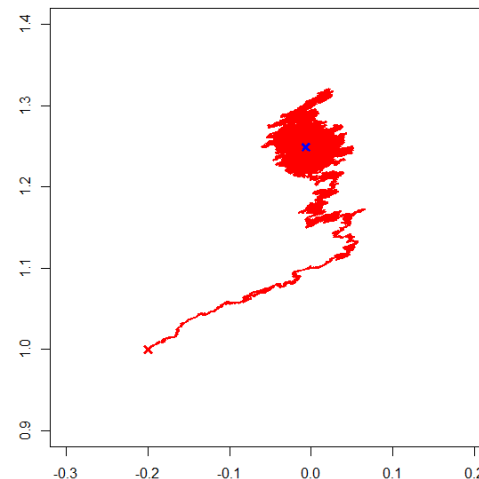
Course homepage:

<http://gauss.stat.su.se/phd/oasi/optimisation2.html>

Includes reading material, lecture notes, assignments

Today's schedule

- Stochastic gradient descent (SGD) algorithms
 - Quick recap
 - Convergence analysis for SGD
- Exercise session
- Quasi-Newton method
 - Idea and definition
 - The BFGS method

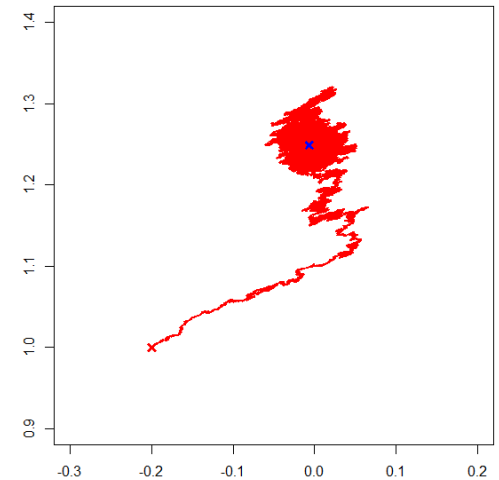


Multivariate optimisation – Gradient descent method

- Iteration: $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha_t \mathbf{g}'(\mathbf{x}^{(t)})$
- Optimisation problem (finite sum case):
 - \mathbf{x} p -dimensional vector, $g_i: \mathbb{R}^p \rightarrow \mathbb{R}$ functions
 - We search \mathbf{x}^* with $g(\mathbf{x}^*) = \min g(\mathbf{x})$ where $g = \frac{1}{n} \sum_{i=1}^n g_i$
- If n large: Takes time to evaluate gradient $\mathbf{g}' = \frac{1}{n} \sum_{i=1}^n \mathbf{g}'_i$

Stochastic gradient descent method

- Iteration:
 - Choose $i \in \{1, \dots, n\}$ randomly
 - $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha_t \mathbf{g}'_i(\mathbf{x}^{(t)})$
- α_t is a predefined sequence, either
 - constant step size $\alpha_t = \alpha$ or
 - decreasing step size e.g. $\alpha_t = \alpha/t$
- Convergence (to a local minimum) can be shown if step size fulfils $\sum_{t=1}^{\infty} \alpha_t = \infty$ and $\sum_{t=1}^{\infty} \alpha_t^2 < \infty$ (example: $\alpha_t = \alpha/t$)
 - > this claim will be made more specify shortly



Stochastic gradient descent method

- Function to be **minimised**: $g = \frac{1}{n} \sum_{i=1}^n g_i$
- Predefined sequence of step sizes: $\alpha_t, t = 1, 2, \dots$
- Starting value: $\mathbf{x}^{(0)}$
- Sequence of random numbers: $R^{(t)} \in \{1, \dots, n\}, t = 1, 2, \dots$
- Iteration: $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha_t \mathbf{g}'_{R^{(t)}}(\mathbf{x}^{(t)})$

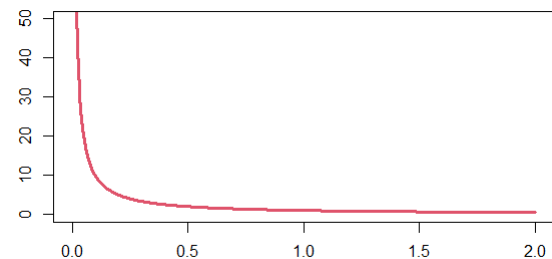
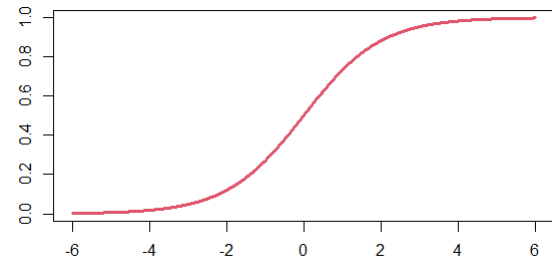
- Now: Convergence analysis
- We assume in the lecture:
 - $R^{(t)}$ uniformly distributed on $\{1, \dots, n\}$, all $R^{(t)}$ independent
- We note that $E \mathbf{g}'_{R^{(t)}}(\mathbf{x}^{(t)}) = \mathbf{g}'(\mathbf{x}^{(t)})$

Lipschitz continuous functions

- A function f is called *Lipschitz continuous* with Lipschitz constant $L > 0$, if for all $\mathbf{x}_1, \mathbf{x}_2$,

$$\|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2 \leq L \cdot \|\mathbf{x}_1 - \mathbf{x}_2\|_2.$$

- If $f: (a, b) \rightarrow \mathbb{R}$ is differentiable, the following is true:
f Lipschitz continuous with constant L if and only if $|f'(x)| \leq L$ for all x
- Example: $1/(1+\exp(-x))$ is Lipschitz continuous with $L=0.25$
- Example: $1/x$ is not Lipschitz continuous on $(0, \infty)$



Lipschitz continuous functions

- A function f is called *Lipschitz continuous* with Lipschitz constant $L > 0$, if for all $\mathbf{x}_1, \mathbf{x}_2$,

$$\|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2 \leq L \cdot \|\mathbf{x}_1 - \mathbf{x}_2\|_2.$$

- If f has a derivative (gradient) \mathbf{f}' which is Lipschitz continuous with $L > 0$, then f itself is called *L-smooth*. Further,

$$f(\mathbf{x}_1) - f(\mathbf{x}_2) \leq \mathbf{f}'(\mathbf{x}_2)^T (\mathbf{x}_1 - \mathbf{x}_2) + \frac{L}{2} \cdot \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2.$$

- If f has a Hessian matrix \mathbf{f}'' with a bounded spectral norm (by L), the gradient \mathbf{f}' is Lipschitz continuous with L :

$$\|\mathbf{f}''(\mathbf{x})\|_{\text{spectral}} \leq L \text{ for all } \mathbf{x} \Rightarrow \mathbf{f}' \text{ Lipschitz continuous with } L$$



Matrix norms

- Most often when writing $\|\cdot\|$, we have a norm for a vector inside the norm-signs (and $\|x\|_2$ can be interpreted as length of vector x)
- There are also matrix-norms, and the spectral norm is one example: $\|A\|_{\text{spectral}} = \sqrt{\lambda_{\max}(A^T A)}$ where $\lambda_{\max}(\cdot)$ is the largest eigenvalue of the matrix inside
- Spectral norm and Euclidian norm are *compatible* in the sense that for any $A \in \mathbb{R}^{n \times n}$ and $x \in \mathbb{R}^n$, we have

$$\|Ax\|_2 \leq \|A\|_{\text{spectral}} \|x\|_2$$

SGD's expected decrease per iteration

- Minimisation of $g = \frac{1}{n} \sum_{i=1}^n g_i$ with SGD

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha_t \mathbf{g}'_{R(t)}(\mathbf{x}^{(t)}) \quad (\text{SGD})$$

- Lemma 1 (Bottou et al): Let g be L -smooth with $L > 0$. Given $\mathbf{x}^{(t)}$, the expected decrease in an SGD iteration is bounded:

$$E[g(\mathbf{x}^{(t+1)})] - g(\mathbf{x}^{(t)}) \leq -\alpha_t \|\mathbf{g}'(\mathbf{x}^{(t)})\|_2^2 + \alpha_t^2 \frac{L}{2} E \|\mathbf{g}'_{R(t)}(\mathbf{x}^{(t)})\|_2^2$$

SGD's expected decrease per iteration

- Detailed proof of Lemma 1 (Bottou et al):

- We have:

- $$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha_t \mathbf{g}'_{R(t)}(\mathbf{x}^{(t)}) \quad \text{(SGD)}$$

- $$g(\mathbf{x}_1) - g(\mathbf{x}_2) \leq \mathbf{g}'(\mathbf{x}_2)^T (\mathbf{x}_1 - \mathbf{x}_2) + \frac{L}{2} \cdot \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 \text{ for all } \mathbf{x}_1, \mathbf{x}_2 \quad \text{(Lsmooth)}$$

- $$R(t) \text{ uniformly distributed on } \{1, \dots, n\} \quad \text{(R)}$$

- Using (Lsmooth) for $\mathbf{x}_1 = \mathbf{x}^{(t+1)}$ and $\mathbf{x}_2 = \mathbf{x}^{(t)}$ (conditional on $R(t)$ and $\mathbf{x}^{(t)}$),

$$g(\mathbf{x}^{(t+1)}) - g(\mathbf{x}^{(t)}) \leq \mathbf{g}'(\mathbf{x}^{(t)})^T (\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}) + \frac{L}{2} \cdot \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2$$

- Using (SGD),

$$= \mathbf{g}'(\mathbf{x}^{(t)})^T \left(-\alpha_t \mathbf{g}'_{R(t)}(\mathbf{x}^{(t)}) \right) + \frac{L}{2} \cdot \|\alpha_t \mathbf{g}'_{R(t)}(\mathbf{x}^{(t)})\|_2^2$$

- $$= -\alpha_t \mathbf{g}'(\mathbf{x}^{(t)})^T \mathbf{g}'_{R(t)}(\mathbf{x}^{(t)}) + \alpha_t^2 \frac{L}{2} \cdot \|\mathbf{g}'_{R(t)}(\mathbf{x}^{(t)})\|_2^2$$

- Take expectation over $R(t)$ given $\mathbf{x}^{(t)}$

$$E[g(\mathbf{x}^{(t+1)})] - g(\mathbf{x}^{(t)}) \leq -\alpha_t \mathbf{g}'(\mathbf{x}^{(t)})^T E \left[\mathbf{g}'_{R(t)}(\mathbf{x}^{(t)}) \right] + \alpha_t^2 \frac{L}{2} E \left[\|\mathbf{g}'_{R(t)}(\mathbf{x}^{(t)})\|_2^2 \right]$$

- $$= -\alpha_t \|\mathbf{g}'(\mathbf{x}^{(t)})\|_2^2 + \alpha_t^2 \frac{L}{2} E \|\mathbf{g}'_{R(t)}(\mathbf{x}^{(t)})\|_2^2 \text{ since}$$

- $$E \left[\mathbf{g}'_{R(t)}(\mathbf{x}^{(t)}) \right] = \frac{1}{n} \sum_{i=1}^n \mathbf{g}'_i(\mathbf{x}^{(t)}) = \mathbf{g}'(\mathbf{x}^{(t)}) \text{ due to (R).}$$

□



SGD's expected decrease per iteration

- Minimisation of $g = \frac{1}{n} \sum_{i=1}^n g_i$ with SGD

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha_t \mathbf{g}'_{R(t)}(\mathbf{x}^{(t)}) \quad (\text{SGD})$$

- Lemma 1 (Bottou et al): Let g be L -smooth with $L > 0$. Given $\mathbf{x}^{(t)}$, the expected decrease in an SGD iteration is bounded:

$$E[g(\mathbf{x}^{(t+1)})] - g(\mathbf{x}^{(t)}) \leq -\alpha_t \|\mathbf{g}'(\mathbf{x}^{(t)})\|_2^2 + \alpha_t^2 \frac{L}{2} E \|\mathbf{g}'_{R(t)}(\mathbf{x}^{(t)})\|_2^2$$

- Proof idea: We apply the **consequence of L-smoothness** for $\mathbf{x}_1 = \mathbf{x}^{(t+1)}$ and $\mathbf{x}_2 = \mathbf{x}^{(t)}$ (conditional on $R^{(t)}$ and $\mathbf{x}^{(t)}$),
 $g(\mathbf{x}^{(t+1)}) - g(\mathbf{x}^{(t)}) \leq \mathbf{g}'(\mathbf{x}^{(t)})^T (\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}) + \frac{L}{2} \cdot \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2$
We **use Equation (SGD)** above, **take expectation over $R^{(t)}$** (given $\mathbf{x}^{(t)}$ or the history $R^{(t-1)}, R^{(t-2)}, \dots$) and **replace $E \mathbf{g}'_{R(t)}(\mathbf{x}^{(t)})$ by $\mathbf{g}'(\mathbf{x}^{(t)})$** . This shows the claim. □



SGD's expected decrease per iteration

- Minimisation of $g = \frac{1}{n} \sum_{i=1}^n g_i$ with SGD

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha_t \mathbf{g}'_{R(t)}(\mathbf{x}^{(t)})$$

- Lemma 2 (Bottou et al): Let g be L -smooth with $L > 0$ and we have following second moment condition:

$$E \|\mathbf{g}'_{R(t)}(\mathbf{x}^{(t)})\|_2^2 \leq s + w \|\mathbf{g}'(\mathbf{x}^{(t)})\|_2^2 \text{ for all } t.$$

Given $\mathbf{x}^{(t)}$, the expected decrease in an SGD iteration is bounded:

$$E[g(\mathbf{x}^{(t+1)})] - g(\mathbf{x}^{(t)}) \leq -\alpha_t(1 - \alpha_t Lw/2) \|\mathbf{g}'(\mathbf{x}^{(t)})\|_2^2 + \alpha_t^2 \frac{Ls}{2}$$

- Proof: Follows directly from Lemma 1:

- $$\begin{aligned} E[g(\mathbf{x}^{(t+1)})] - g(\mathbf{x}^{(t)}) &\leq -\alpha_t \|\mathbf{g}'(\mathbf{x}^{(t)})\|_2^2 + \alpha_t^2 \frac{L}{2} E \|\mathbf{g}'_{R(t)}(\mathbf{x}^{(t)})\|_2^2 \\ &\leq -\alpha_t \|\mathbf{g}'(\mathbf{x}^{(t)})\|_2^2 + \alpha_t^2 \frac{L}{2} (s + w \|\mathbf{g}'(\mathbf{x}^{(t)})\|_2^2) \end{aligned}$$



$$\|y\|_2 = \sqrt{y^2} = |y|$$

$$\|y\|_2^2 = y^2$$

Set $y = x_1 - x_2$:

$$\|x_1 - x_2\|_2^2 = (x_1 - x_2)^2$$

Strongly convex functions

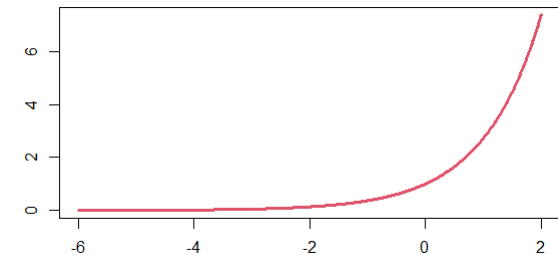
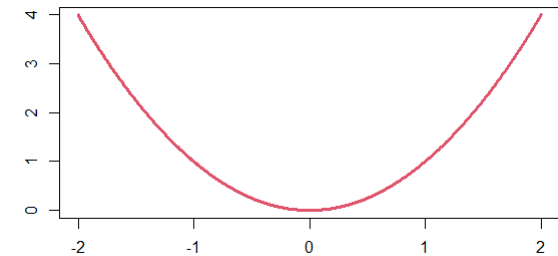
- A differentiable function f is called *m-strongly convex* with $m > 0$, if for all $\mathbf{x}_1, \mathbf{x}_2$,

$$(f'(x_1) - f'(x_2))^T (x_1 - x_2) \geq m \cdot \|x_1 - x_2\|_2^2.$$

- For one-dimensional functions:

$$(f'(x_1) - f'(x_2))/(x_1 - x_2) \geq m \text{ for all } x_1, x_2.$$

- The function $f(x) = x^2$ is m-strongly convex with $m = 2$
- The function $f(x) = \exp(x)$ is convex but not m-strongly convex since for $x \rightarrow -\infty$, smaller and smaller m would be necessary; no $m > 0$ can be found to fulfil condition above



Strongly convex functions

- A differentiable function f is called *m-strongly convex* with $m > 0$, if for all $\mathbf{x}_1, \mathbf{x}_2$,

$$(\mathbf{f}'(\mathbf{x}_1) - \mathbf{f}'(\mathbf{x}_2))^T (\mathbf{x}_1 - \mathbf{x}_2) \geq m \cdot \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2.$$

- An equivalent condition is

$$(f(\mathbf{x}_1) - f(\mathbf{x}_2)) \geq \mathbf{f}'(\mathbf{x}_2)^T (\mathbf{x}_1 - \mathbf{x}_2) + \frac{m}{2} \cdot \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2.$$

- An m -strongly convex function f has a unique minimum \mathbf{x}^* and the following holds true:

$$2m(f(\mathbf{x}) - f(\mathbf{x}^*)) \leq \|\mathbf{f}'(\mathbf{x})\|_2^2.$$



Assumptions

- Optimisation problem (finite sum case):
 - \mathbf{x} p -dimensional vector, $g_i: \mathbb{R}^p \rightarrow \mathbb{R}$ functions
 - We search \mathbf{x}^* with $g(\mathbf{x}^*) = \min g(\mathbf{x})$ where $g = \frac{1}{n} \sum_{i=1}^n g_i$
- Assumptions (A):
 - g is differentiable and L -smooth with $L > 0$,
 - g is m -strongly convex with $m > 0$
 - For all \mathbf{x} : $E \|\mathbf{g}'_{R(t)}(\mathbf{x})\|_2^2 \leq s + w \|\mathbf{g}'(\mathbf{x})\|_2^2$
- Assumptions (B):
 - g_i are differentiable and L -smooth with $L_i > 0$,
 - g is m -strongly convex with $m > 0$
 - $E \|\mathbf{g}'_{R(t)}(\mathbf{x}^*)\|_2^2 = s$



Convergence analysis for fixed step size

- Theorem 1 (Bottou et al): Consider the finite sum case of the optimization problem, assume Assumptions (A) and that the step size is constant, $\alpha_t = \alpha \leq 1/\{L \max(w, 1)\}$. Then, we have the following convergence result:

$$\mathbb{E}[g(\mathbf{x}^{(t)})] - g(\mathbf{x}^*) \leq \frac{\alpha L s}{2m} + (1 - \alpha m)^t \left\{ g(\mathbf{x}^{(0)}) - g(\mathbf{x}^*) - \frac{\alpha L s}{2m} \right\}$$

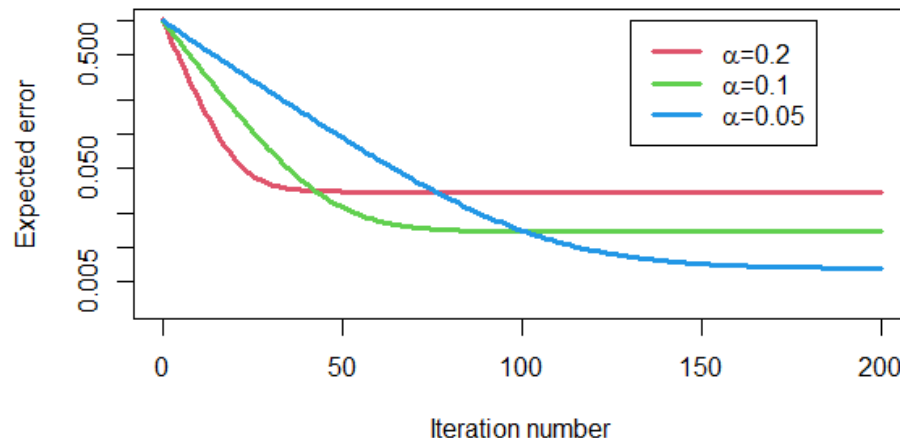
- Proof: Based on Lemma 2, see Bottou et al (2018), <https://arxiv.org/pdf/1606.04838.pdf>

Convergence analysis for fixed step size

- Theorem 2 (Needell et al): Consider the finite sum case of the optimization problem, assume Assumptions (B) and that the step size is constant, $\alpha_t = \alpha \leq \frac{1}{\max(L_i)}$. Then, we have the following convergence result:

$$\mathbb{E} \|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 \leq \frac{\alpha s}{m\{1 - \alpha \max(L_i)\}} + (1 - \alpha m\{1 - \alpha \max(L_i)\})^t \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2$$

- Proof: See Needell et al (2016), here an arXiv version: <https://arxiv.org/abs/1310.5715>



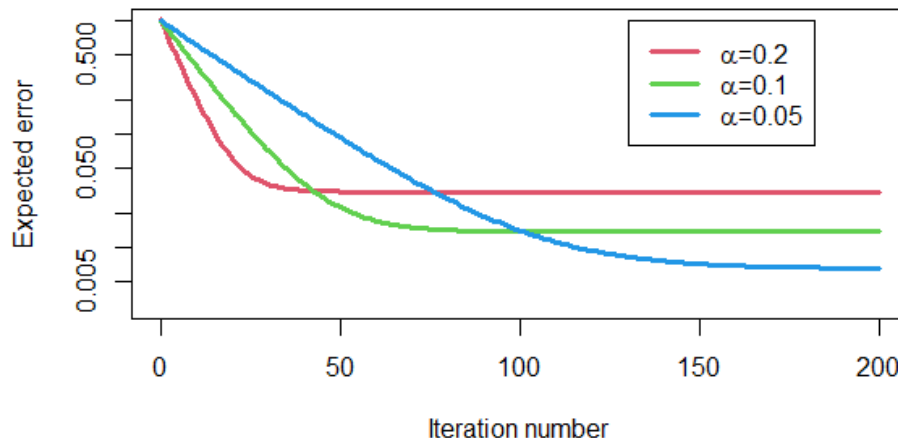
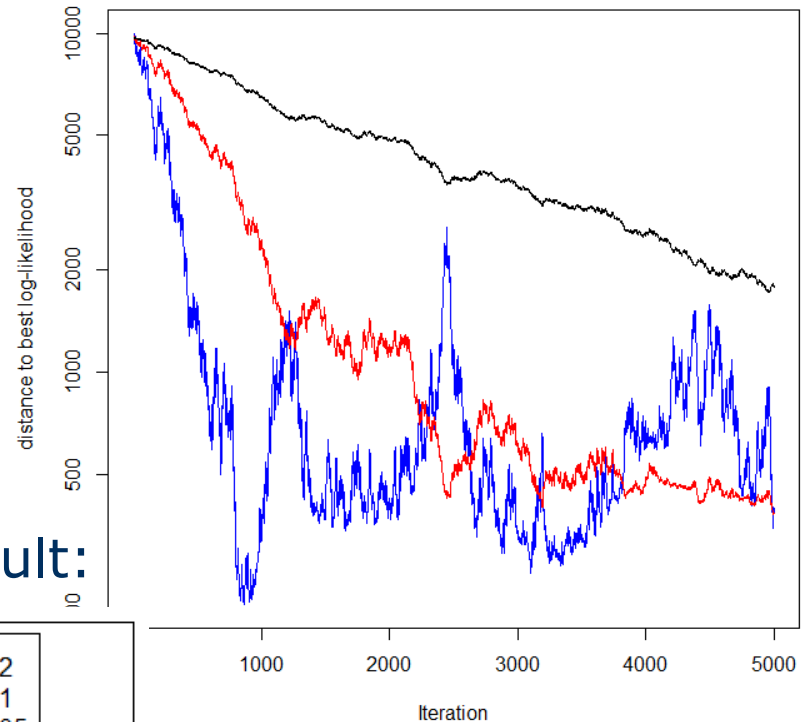
Theoretical behaviour of above bound for expected distance to optimum for $s=0.5$, $m=2$,

$$\max(L_i)=2, \varepsilon_0 = \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 = 1$$



Stochastic gradient descent method – empirical examples from OASI-L2

- Constant step size $\alpha^{(t)} = \alpha$
- Step size
 - $\alpha = 0.0006$ (black)
 - $\alpha = 0.002$ (red)
 - $\alpha = 0.006$ (blue)
- Compare with theoretical result:

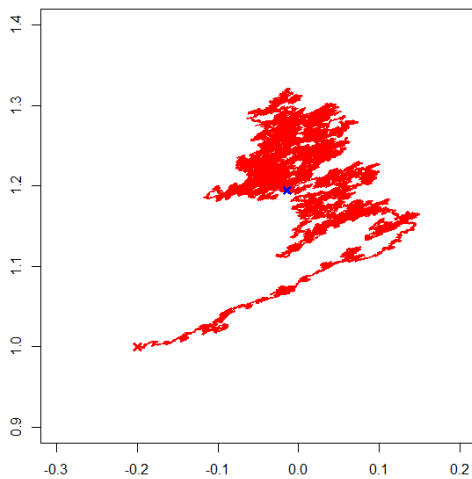


Theoretical behaviour of bound in Theorem 2 for expected distance to optimum for $s=0.5$, $m=2$, $\max(L_i)=2$, $\varepsilon_0=1$

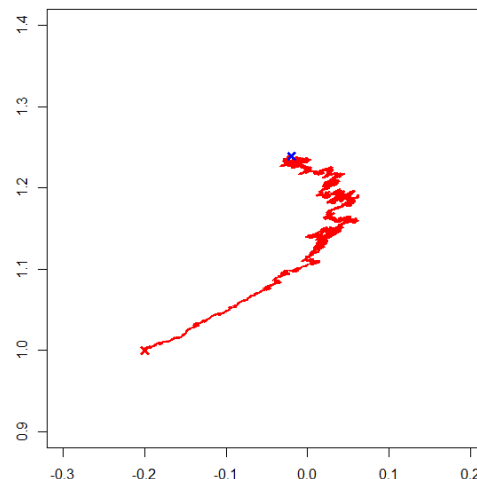


Stochastic gradient descent method – empirical examples from OASI-L2

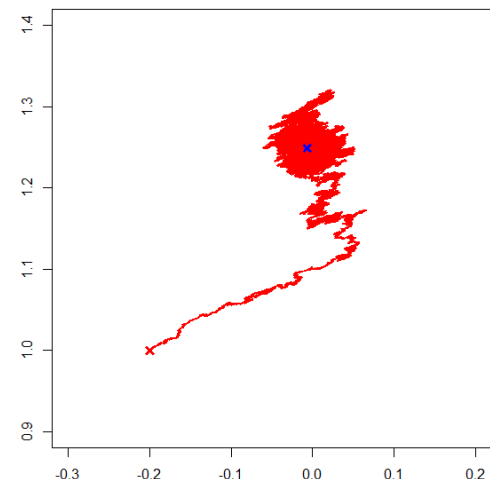
- Problem 1.3 with $n = 1\,000\,000$ (data duplicated 100 000 x); start at $(\beta_0, \beta_1) = (-0.2, 1)$; best $(\hat{\beta}_0, \hat{\beta}_1) = (-0.009, 1.263)$
- 100 000 iterations, $\alpha = 0.002$, $\alpha = 0.0006$
- 1 000 000 iterations, $\alpha = 0.0006$



• $(-0.014, 1.195)$



$(-0.020, 1.239)$



$(-0.007, 1.248)$



Convergence analysis for fixed step size

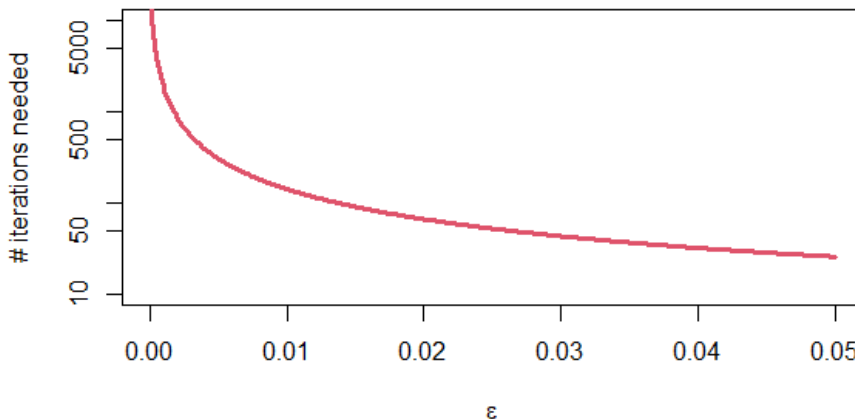
- Corollary:** Consider the finite sum case of the optimization problem, assume Assumptions (B) and set $\varepsilon_0 = \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2$. For any desired $\varepsilon > 0$, using constant step size, $\alpha_t = \alpha =$

$$\frac{m\varepsilon}{2m\varepsilon \max(L_i) + 2s'}$$

we have after

$$t = 2 \log(2\varepsilon_0/\varepsilon) \left(\frac{\max(L_i)}{m} + \frac{s}{m^2\varepsilon} \right)$$

iterations, $E \left[\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2^2 \right] \leq \varepsilon$.



Number of iterations needed to bring expected error below ε for $s=0.5$, $m=2$, $\max(L_i)=2$, $\varepsilon_0=1$

eps	# iter	alpha
0.00001	305177	0.00002
0.00002	143935	0.00004
0.00004	67646	0.00008
0.00008	31666	0.00016
0.00016	14759	0.00032
0.00032	6846	0.00064
0.00064	3160	0.00127
0.00128	1452	0.00253
0.00256	664	0.00502

Convergence analysis for decreasing step size

- Theorem 3 (Bottou et al): Consider the finite sum case of the optimization problem, assume Assumptions (A) and that the step size is decreasing as $\alpha_t = \frac{\beta}{t+\gamma}$ with $\beta > \frac{1}{m}$, $\gamma > 0$, $\alpha_0 \leq 1/\{L \max(w, 1)\}$. Then, we have the following convergence result:

$$\mathbb{E}[g(\mathbf{x}^{(t)})] - g(\mathbf{x}^*) \leq v/(\gamma + t),$$

where

$$v = \max \left\{ \frac{\beta^2 L s}{2(\beta m - 1)}, (\gamma + 1)(g(\mathbf{x}^{(0)}) - g(\mathbf{x}^*)) \right\}.$$

- Proof: See Bottou et al (2018),
<https://arxiv.org/pdf/1606.04838.pdf>

Convergence analysis for decreasing step size

- Note that

$$\mathbb{E}[g(\mathbf{x}^{(t)})] - g(\mathbf{x}^*) \approx v/(\gamma + t),$$

means sublinear convergence since

$$\frac{\{\mathbb{E}[g(\mathbf{x}^{(t+1)})] - g(\mathbf{x}^*)\}}{\{\mathbb{E}[g(\mathbf{x}^{(t)})] - g(\mathbf{x}^*)\}} \approx \frac{\gamma+t}{\gamma+t+1} \rightarrow 1 \quad (\text{for } t \rightarrow \infty)$$

- A bound like

$$\mathbb{E}[g(\mathbf{x}^{(t)})] - g(\mathbf{x}^*) \leq v^t \quad \text{with } 0 < v < 1$$

would lead to linear convergence

- So, SGD with $\alpha_t = \frac{\beta}{t+\gamma}$ gives only sublinear convergence
- But in reality, we never go to infinity with t...



SGD convergence analysis – exercise

Optimisation in a least squares situation,

- $g(\mathbf{b}) = \frac{1}{n} \|\mathbf{X}\mathbf{b} - \mathbf{y}\|_2^2 = \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{b})$ with $g_i(\mathbf{b}) = (\mathbf{x}_i^T \mathbf{b} - y_i)^2$
- $\mathbf{g}'(\mathbf{b}) = \frac{2}{n} \mathbf{X}^T (\mathbf{X}\mathbf{b} - \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n \mathbf{g}'_i(\mathbf{b})$ with $\mathbf{g}'_i(\mathbf{b}) = 2(\mathbf{x}_i^T \mathbf{b} - y_i) \mathbf{x}_i$
- $\mathbf{g}''(\mathbf{b}) = \frac{2}{n} \mathbf{X}^T \mathbf{X}$

R uniformly distributed on $\{1, \dots, n\}$

Compute for (i) general \mathbf{x}_i , (ii) $\mathbf{x}_i = \begin{pmatrix} 1 \\ w_i \end{pmatrix}$ (straight line regr.):

a) $\|\mathbf{g}'_i(\mathbf{b})\|_2^2 = 4 (\mathbf{x}_i^T \mathbf{b} - y_i)^2 \mathbf{x}_i^T \mathbf{x}_i = 4 (b_1 + b_2 w_i - y_i)^2 (1 + w_i^2)$

b) $E\|\mathbf{g}'_R(\mathbf{b})\|_2^2 = \frac{1}{n} \sum_{i=1}^n \mathbf{g}'_i(\mathbf{b}) = \frac{1}{n} \sum_{i=1}^n 4 (\mathbf{x}_i^T \mathbf{b} - y_i)^2 \mathbf{x}_i^T \mathbf{x}_i = \dots$

Compute for general \mathbf{x}_i, \mathbf{X} :

c) $E[\mathbf{g}'_R(\mathbf{b})] = \frac{1}{n} \sum_{i=1}^n \mathbf{g}'_i(\mathbf{b}) = \mathbf{g}'(\mathbf{b})$

d) $\|\mathbf{g}''(\mathbf{b})\|_{\text{spectral}} = \frac{2}{n} \sqrt{\lambda_{\max}((\mathbf{X}^T \mathbf{X})^T (\mathbf{X}^T \mathbf{X}))} = \frac{2}{n} \lambda_{\max}(\mathbf{X}^T \mathbf{X})$

e) $\|\mathbf{g}'(\mathbf{b})\|_2^2 = \left(\frac{2}{n} \mathbf{X}^T (\mathbf{X}\mathbf{b} - \mathbf{y}) \right)^T \frac{2}{n} \mathbf{X}^T (\mathbf{X}\mathbf{b} - \mathbf{y})$
 $= \frac{4}{n^2} (\mathbf{X}\mathbf{b} - \mathbf{y})^T \mathbf{X} \mathbf{X}^T (\mathbf{X}\mathbf{b} - \mathbf{y}) = \dots$



Quasi-Newton method

Quasi-Newton method

- Steepest descent and Newton method have iteration

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - (\mathbf{M}^{(t)})^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$

with $\mathbf{M}^{(t)} = \mathbf{g}''(\mathbf{x}^{(t)})$ for the Newton method and

with $(\mathbf{M}^{(t)})^{-1} = \alpha_t \mathbf{I}$ for the steepest descent method

- A disadvantage of Newton is the need to calculate the Hessian $\mathbf{g}''(\mathbf{x}^{(t)})$ in each iteration
- A disadvantage of steepest descent is that no information about the curvature is used
- We can monitor the computed gradients $\mathbf{g}'(\mathbf{x}^{(t)})$ and their change gives information about the curvature of g

Quasi-Newton method

- Steepest descent and Newton method have iteration

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - (\mathbf{M}^{(t)})^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$

- Newton ($\mathbf{M}^{(t)} = \mathbf{g}''(\mathbf{x}^{(t)})$) was motivated with the multidimensional Taylor expansion

$$\mathbf{g}'(\mathbf{x}^*) \approx \mathbf{g}'(\mathbf{x}^{(t)}) + \mathbf{g}''(\mathbf{x}^{(t)})(\mathbf{x}^* - \mathbf{x}^{(t)})$$

or

$$\mathbf{g}'(\mathbf{x}^*) - \mathbf{g}'(\mathbf{x}^{(t)}) \approx \mathbf{g}''(\mathbf{x}^{(t)})(\mathbf{x}^* - \mathbf{x}^{(t)})$$

- We want to use approximations $\mathbf{M}^{(t+1)}$ to $\mathbf{g}''(\mathbf{x}^{(t)})$ which fulfil this relation when \mathbf{x}^* is replaced by $\mathbf{x}^{(t+1)}$:

$$\mathbf{g}'(\mathbf{x}^{(t+1)}) - \mathbf{g}'(\mathbf{x}^{(t)}) = \mathbf{M}^{(t+1)}(\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)})$$

- This condition is called secant condition
- There are multiple solutions to the secant condition



Quasi-Newton method

- Steepest descent and Newton method have iteration

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - (\mathbf{M}^{(t)})^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$

- Secant condition:

$$\mathbf{g}'(\mathbf{x}^{(t+1)}) - \mathbf{g}'(\mathbf{x}^{(t)}) = \mathbf{M}^{(t+1)}(\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)})$$

- Or, with $\mathbf{y}^{(t)} = \mathbf{g}'(\mathbf{x}^{(t+1)}) - \mathbf{g}'(\mathbf{x}^{(t)})$ and $\mathbf{z}^{(t)} = \mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}$:

$$\mathbf{y}^{(t)} = \mathbf{M}^{(t+1)} \mathbf{z}^{(t)}$$

- Suggestion from Broyden, Fletcher, Goldfarb, and Shanno (BFGS; 4 publications in 1970) fulfilling secant condition:

$$\mathbf{M}^{(t+1)} = \mathbf{M}^{(t)} - \frac{\mathbf{M}^{(t)} \mathbf{z}^{(t)} (\mathbf{M}^{(t)} \mathbf{z}^{(t)})^T}{\mathbf{z}^{(t)T} \mathbf{M}^{(t)} \mathbf{z}^{(t)}} + \frac{\mathbf{y}^{(t)} \mathbf{y}^{(t)T}}{\mathbf{y}^{(t)T} \mathbf{z}^{(t)}}$$



Quasi-Newton method

- The BFGS (quasi-Newton) method has iteration

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - (\mathbf{M}^{(t)})^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$

and

$$\mathbf{M}^{(t+1)} = \mathbf{M}^{(t)} - \frac{\mathbf{M}^{(t)} \mathbf{z}^{(t)} (\mathbf{M}^{(t)} \mathbf{z}^{(t)})^T}{\mathbf{z}^{(t)T} \mathbf{M}^{(t)} \mathbf{z}^{(t)}} + \frac{\mathbf{y}^{(t)} \mathbf{y}^{(t)T}}{\mathbf{y}^{(t)T} \mathbf{z}^{(t)}}$$

- Descent is not ensured but stepsize-halving (“backtracking”) can be used as for steepest descent to ensure it:

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha_t (\mathbf{M}^{(t)})^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$

- The `R` function `optim` includes the quasi-Newton BFGS
- Convergence of quasi-Newton methods is faster than linear but slower than quadratic (some assumptions necessary; see e.g. Nocedal and Wright, 2006, Theorem 3.7)

$$\lim_{t \rightarrow \infty} \frac{\|\mathbf{x}^{(t+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(t)} - \mathbf{x}^*\|} = 0, \quad \frac{\|\mathbf{x}^{(t+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(t)} - \mathbf{x}^*\|^2} \text{ diverging}$$

Quasi-Newton method

- The BFGS (quasi-Newton) method has iteration

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - (\mathbf{M}^{(t)})^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$

and

$$\mathbf{M}^{(t+1)} = \mathbf{M}^{(t)} - \frac{\mathbf{M}^{(t)} \mathbf{z}^{(t)} (\mathbf{M}^{(t)} \mathbf{z}^{(t)})^T}{\mathbf{z}^{(t)T} \mathbf{M}^{(t)} \mathbf{z}^{(t)}} + \frac{\mathbf{y}^{(t)} \mathbf{y}^{(t)T}}{\mathbf{y}^{(t)T} \mathbf{z}^{(t)}}$$

- For higher dimensional problems, computation of inverse could be computationally intensive, but can be avoided by applying the Sherman-Morrison-Woodbury formula
- Defining $\mathbf{L}^{(t)} = (\mathbf{M}^{(t)})^{-1}$, the matrix update can be done by
$$\mathbf{L}^{(t+1)} = \left(\mathbf{I} - r_k \mathbf{z}^{(t)} (\mathbf{y}^{(t)})^T \right) \mathbf{L}^{(t)} \left(\mathbf{I} - r_k \mathbf{y}^{(t)} (\mathbf{z}^{(t)})^T \right) + r_k \mathbf{z}^{(t)} \mathbf{z}^{(t)T}$$
with $r_k = 1/(\mathbf{y}^{(t)T} \mathbf{z}^{(t)})$
- Starting value could be $\mathbf{L}^{(0)} = \mathbf{M}^{(0)} = \mathbf{I}$

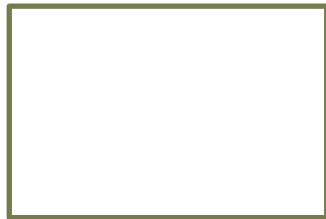


Quasi-Newton method

- Convergence can be extremely fast in praxis
- Example from Nocedal and Wright (2006), chapter 6:
Rosenbrock function $g(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$, starting point $(-1.2, 1)$, optimum at $(1,1)$.
#iterations until error $< 10^{-5}$:
 - Steepest descent 5264
 - BFGS 34
 - Newton 21

Quasi-Newton method

- In high dimensional problems, storage and update of $M^{(t)}$ might be problematic
- L-BFGS modification deals with it (saves not $M^{(t)}$ but sequence of $y^{(t)}$ and $z^{(t)}$)
- Simple box-constraints are handled with BFGS-B



- Together: L-BFGS-B which is option in `optim`

Convergence order for deterministic algorithms

- Recall: Convergence order and convergence rate

$$\frac{\{E[g(\mathbf{x}^{(t+1)})] - g(\mathbf{x}^*)\}}{\{E[g(\mathbf{x}^{(t)})] - g(\mathbf{x}^*)\}^p} \rightarrow c \quad (\text{for } t \rightarrow \infty)$$

- p is convergence order ($p=1$, $0 < c < 1$ linear; $p=2$, $0 < c < 1$ quadratic)
- c is convergence rate
- Under certain assumption, we have following orders:

unidimensional	Bisection order = roughly 1^*	Secant order = $(1 + \sqrt{5})/2$	Newton order = 2
multidimensional	Steepest descent order = 1	Quasi-Newton order $> 1^{**}$	Newton order = 2

*strictly, the above criterion cannot be proven for bisection

**criterion above fulfilled for $p=1$ and $c=0$; “superlinear”



Assignments

- Topic 1: March 23 until April 12
- Topic 2: April 13 until April 26
- Topic 3: April 27 until Mai 11
- Second chance for Topic 1-3: until **August 31 (no extension!)**
- I might collect all submissions which have missed the first deadline and look at them in September