



Advanced computational statistics, lecture 2

Frank Miller, Department of Computer and Information Science,
Linköping University

Department of Statistics; Stockholm University

March 17, 2023

Course schedule

- Topic 1: **Gradient based optimisation**
- **Topic 2: Stochastic gradient based optimisation**
- Topic 3: **Gradient free optimisation**
- Topic 4: **Optimisation with constraints**
- Topic 5: **EM algorithm and bootstrap**
- Topic 6: **Simulation of random variables**
- Topic 7: **Importance sampling**

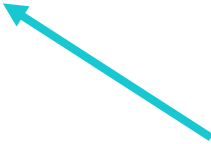
Course homepage:

<http://www.adoptdesign.de/frankmillereu/adcompstat2023.html>

Includes schedule, reading material, lecture notes, assignments

Today's schedule

- Stochastic steepest descent (SSD; Stochastic gradient descent; SGD)
 - Idea and issues
 - Choice of step size
 - Mini-batches
 - Convergence analysis
- Exercise session



Note: Changed to descent and minimisation problem here to correspond to most literature, but this is no essential change.

Steepest descent

- Optimisation problem:
 - \mathbf{x} p -dimensional vector, $g: \mathbb{R}^p \rightarrow \mathbb{R}$ function
 - We search \mathbf{x}^* with $g(\mathbf{x}^*) = \min g(\mathbf{x})$
- Steepest descent:
 - Iteration: $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha^{(t)} \mathbf{g}'(\mathbf{x}^{(t)})$

Steepest descent

- 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 $g = \sum_{i=1}^n g_i$
- If n large: Takes time to evaluate gradient $\mathbf{g}' = \sum_{i=1}^n \mathbf{g}'_i$

Stochastic steepest descent

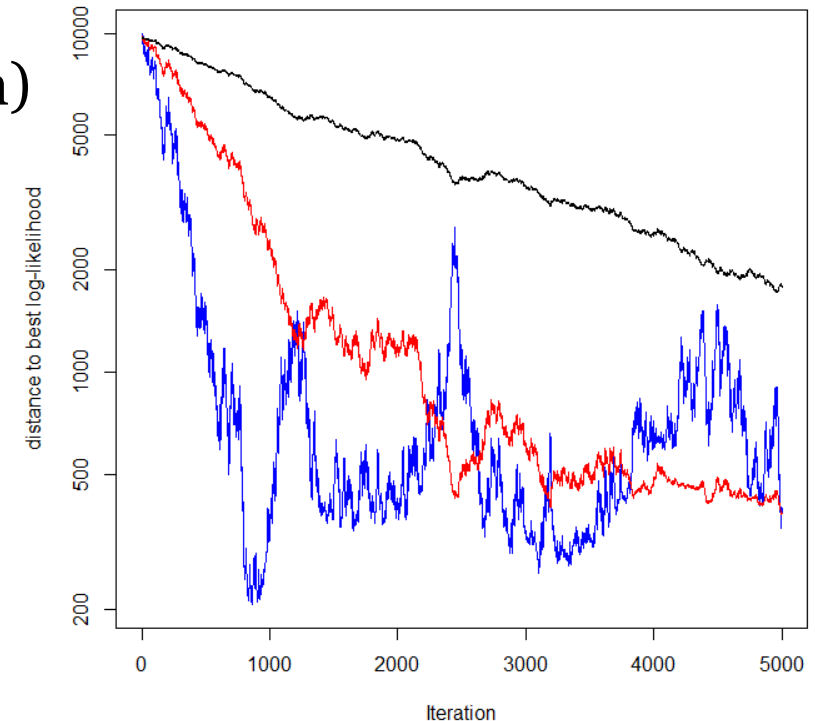
- Iteration:
 - Choose $i \in \{1, \dots, n\}$ randomly
 - $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha^{(t)} \mathbf{g}'_i(\mathbf{x}^{(t)})$
- $\alpha^{(t)}$ is a predefined sequence, either
 - constant step size $\alpha^{(t)} = \alpha$ or
 - decreasing step size e.g. $\alpha^{(t)} = \alpha/t$
-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

Stochastic steepest descent

- Constant step size $\alpha^{(t)} = \alpha$ can still make sense if
 - Another algorithm is run afterwards, or
 - If good but not necessarily best solution desired
- Choice of step size is critical
- Example: Two-parameter MLE computation (large n)
Computation of MLE for a model with two parameters and $n = 1\,000\,000$.
Starting value is not too good (has some distance to correct MLE). We monitor:
 - distance of current log likelihood to maximal log likelihood,
 - search path in 2d parameter space.

Stochastic steepest descent: choice of step size

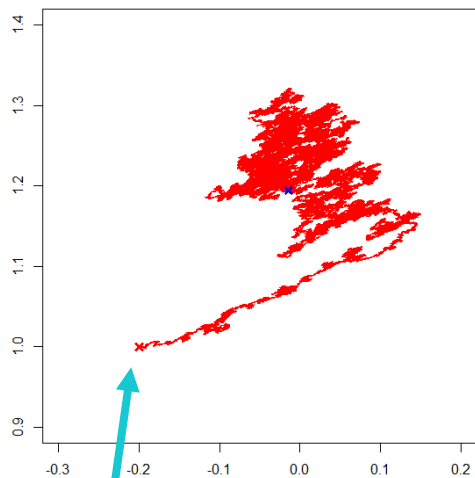
- Example: Two-parameter MLE computation (large n)
- Constant step size $\alpha^{(t)} = \alpha$
- Choice of step size is critical
- Step size here:
 - $\alpha = 0.0006$ (black)
 - $\alpha = 0.002$ (red)
 - $\alpha = 0.006$ (blue)
- If you have time for 5000 iterations: which step size is best?
- If you have only time for 500 iterations?
- If you have time for 50000 iterations?



Stochastic steepest descent: choice of step size

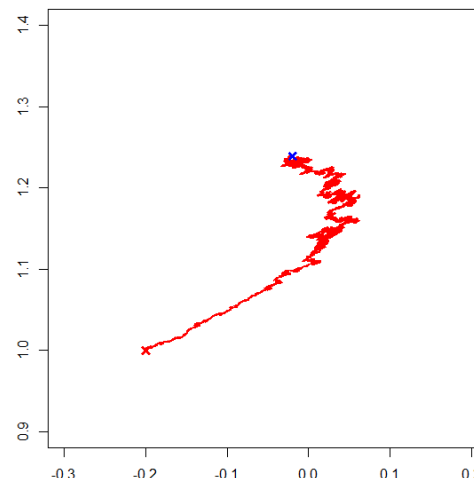
- Example: Two-parameter MLE computation (large n)
Search path in the 2d parameter space

- 100 000 iterations,
 $\alpha = 0.002$,

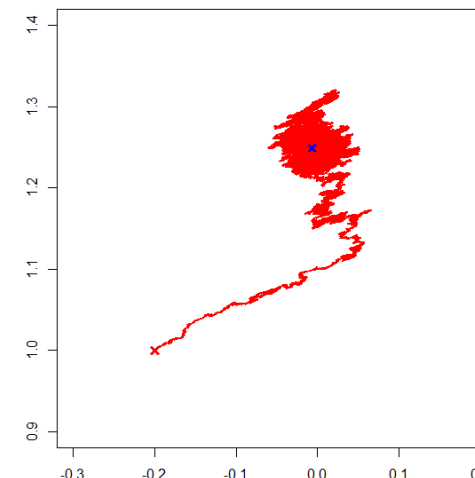


Starting value

- 1 000 000 iterations,
 $\alpha = 0.0006$



- 1 000 000 iterations,
 $\alpha = 0.0006$



Stochastic steepest descent

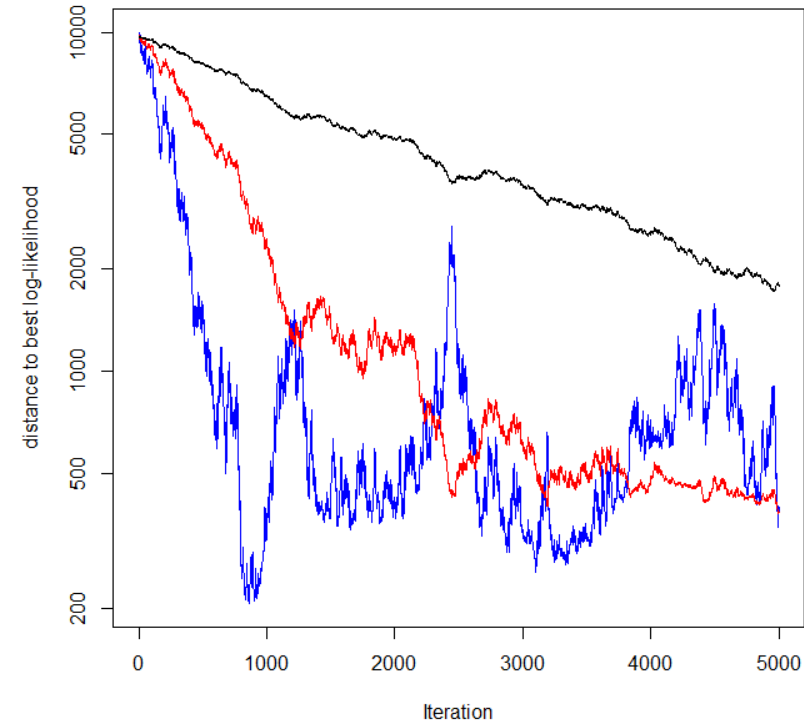
- Influence of step size can be investigated at http://fa.bianp.net/teaching/2018/COMP-652/stochastic_gradient.html (Fabia Pedregosa, Nov 2018)

Stochastic steepest descent: choice of step size

- [Goodfellow et al., 2016](#), Chapter 8.3.1 (notation adjusted):
- “In practice, it is common to decay the learning rate [=step size $\alpha^{(t)}$] linearly until iteration τ : $\alpha^{(t)} = (1 - \gamma)\alpha_0 + \gamma\alpha_\tau$ with $\gamma = \frac{t}{\tau}$. After iteration τ , it is common to leave α constant.”
- Choice of step size “is more of an art than a science, and most guidance on this subject should be regarded with some skepticism.”
- Choose τ “to make a few hundred passes through the training set.”
- $\alpha_\tau \approx \alpha_0/100$
- Choose α_0 avoiding violent oscillations and too low learning rate

(Stochastic) steepest descent: running time

- Example: Two-parameter MLE computation (large n)
- Stochastic steepest ascent: 50000 iterations took 7 s
- Steepest ascent with alpha-halving: 112 iterations took 52 s
- Stochastic steepest ascent could run 3320 iterations when steepest ascent could run 1 iteration



Stochastic steepest descent: mini-batches

- Instead of sampling a single i , a batch of size m can be sampled in each iteration
- Iteration:
 - Choose $\{i_1, \dots, i_m\} \subseteq \{1, \dots, n\}$ randomly
 - $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha^{(t)} \sum_{j=1}^m \mathbf{g}'_{i_j}(\mathbf{x}^{(t)})$
- Decreases risk of large random oscillations
- Especially interesting when algorithm performed on a parallel computer

Accelerated stochastic steepest descent (adding momentum)

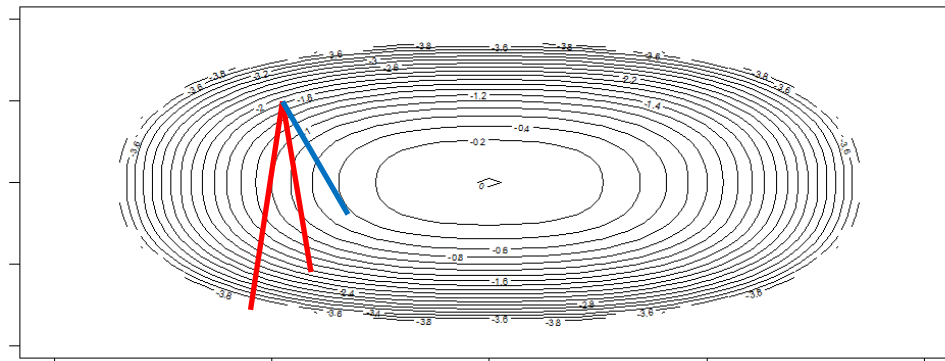
- Stochastic steepest descent $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha \mathbf{g}'_i(\mathbf{x}^{(t)})$ can be combined with momentum method (see [Goodfellow, Bengio, Courville, 2016](#), Chapter 8.3.2)
- Iteration:
 - Choose $i \in \{1, \dots, n\}$ randomly
 - $\mathbf{v}^{(t+1)} = \beta \mathbf{v}^{(t)} - \mathbf{g}'_i(\mathbf{x}^{(t)})$
 - $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha \mathbf{v}^{(t+1)}$
- Advantages:
 - Momentum advantages (handling ill-conditioning, accelerating)
 - Information from previous gradients contribute (variance of stochastic gradient reduced)

Accelerated stochastic steepest ascent (adding momentum)

- Both hyperparameters α, β may depend on iteration number
- Iteration:
 - Choose $i \in \{1, \dots, n\}$ randomly
 - $\mathbf{v}^{(t+1)} = \beta^{(t)} \mathbf{v}^{(t)} - \mathbf{g}'_i(\mathbf{x}^{(t)})$
 - $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha^{(t)} \mathbf{v}^{(t+1)}$
- Changing hyperparameters:
 - $\beta^{(t)}$ usually increased with t , common values 0.5 to 0.99
 - $\alpha^{(t)}$ is decreased with t
 - Decreasing $\alpha^{(t)}$ more important than changing $\beta^{(t)}$

Stochastic steepest descent: adaptive step sizes

- Stochastic steepest descent:
 - Choose $i \in \{1, \dots, n\}$ randomly
 - $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha^{(t)} \mathbf{g}'_i(\mathbf{x}^{(t)})$
- $\alpha^{(t)}$ is now adapted automatically based on previous iterations and separately for each dimension
- If previous gradients in a dimension were large, we want to reduce step size more



Stochastic steepest descent: adaptive step sizes

- **AdaGrad:**

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \text{diag}(\boldsymbol{\alpha}^{(t)}) \mathbf{g}'_i(\mathbf{x}^{(t)}) \quad \text{with vector } \boldsymbol{\alpha}^{(t)}$$

- $\alpha_j^{(t)} = \alpha / \sqrt{\epsilon + \sum_{k=1}^t (g'_{j,k})^2}$

- $g'_{j,k}$ is j^{th} partial derivative of gradient in iteration k

- ϵ is small constant (e.g. $1e-8$)

- A default value $\alpha = 0.01$ is a popular choice

Stochastic steepest descent: adaptive step sizes

- **AdaGrad:**

- Choose $i \in \{1, \dots, n\}$ randomly

- $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \text{diag}(\boldsymbol{\alpha}^{(t)}) \mathbf{g}'_i(\mathbf{x}^{(t)})$

- $\alpha_j^{(t)} = \alpha / \sqrt{\epsilon + \sum_{k=1}^t (g'_{j,k})^2}$

- Disadvantage: $\alpha_j^{(t)}$ can only decrease

- **AdaDelta:**

- $\alpha_j^{(t)} = \alpha / \sqrt{\epsilon + h_j^{(t)}}$

- $h_j^{(t)} = \gamma h_j^{(t-1)} + (1 - \gamma)(g'_{j,t})^2$

(exponential smoothing of earlier partial derivatives; popular choice of γ is around 0.9)

Stochastic steepest descent: adaptive step sizes

- AdaDelta:
 - Random $i \in \{1, \dots, n\}$; $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \text{diag}(\boldsymbol{\alpha}^{(t)})\mathbf{g}'_i(\mathbf{x}^{(t)})$
 - $\alpha_j^{(t)} = \alpha / \sqrt{\epsilon + h_j^{(t)}}$; $h_j^{(t)} = \gamma h_j^{(t-1)} + (1 - \gamma)(g'_{j,t})^2$
- **Adam** ("Adaptive moment estimation")
 - Random $i \in \{1, \dots, n\}$; $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \text{diag}(\boldsymbol{\alpha}^{(t)})\hat{\mathbf{m}}_t$;
 - $\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1)\mathbf{g}'_{i(t)}(\mathbf{x}^{(t)})$ $\hat{\mathbf{m}}_t = \mathbf{m}_t / (1 - \beta_1^t)$
 - $\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2)((g'_{j,t})^2)_{j=1,\dots,p}$ $\hat{\mathbf{v}}_t = \mathbf{v}_t / (1 - \beta_2^t)$
 - $\alpha_j^{(t)} = \alpha / \sqrt{\epsilon + \hat{v}_{j,t}}$
- Default values $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$

Stochastic steepest descent: adaptive step sizes

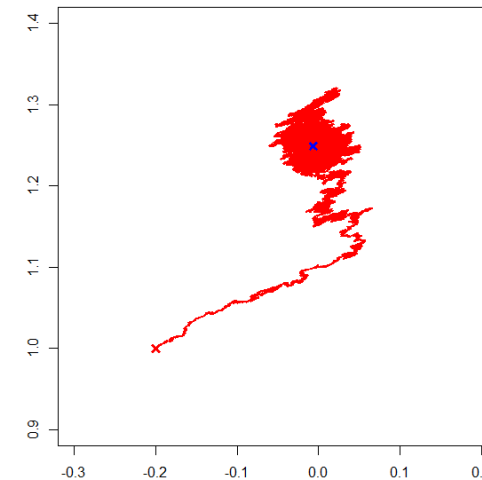
- Momentum method can be added to AdaGrad and AdaDelta
- AdaGrad works well for convex functions
- AdaDelta handles non-convex functions better
- In Adam, momentum method already included

Steepest descent - comparisons of methods

- Animated comparisons:
 - <https://imgur.com/a/Hqolp>

Stochastic steepest descent

- Going back to the stochastic steepest descent (with non-adaptive step sizes)
- 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 maximum) can be shown if step size fulfils $\sum_{t=1}^{\infty} \alpha_t = \infty$ and $\sum_{t=1}^{\infty} \alpha_t^2 < \infty$ (example: $\alpha_t = \alpha/t$)
- Now: Looking closer into the convergence properties, "Convergence analysis"



Stochastic steepest descent (SSD)

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

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

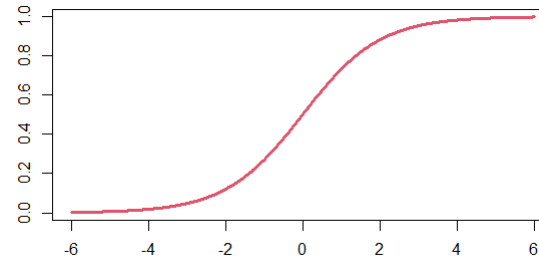
Lipschitz continuous functions

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

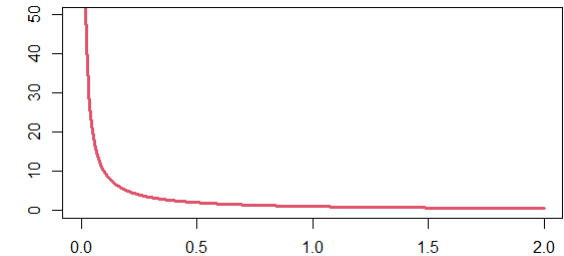
$$\|f(\mathbf{x}_1) - f(\mathbf{x}_2)\|_2 \leq L \cdot \|\mathbf{x}_1 - \mathbf{x}_2\|_2.$$

- 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 $(0, \infty)$



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

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

SSD's expected decrease per iteration

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

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

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

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

SSD's expected decrease per iteration

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

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

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

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

- 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 (SSD),

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

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

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

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

SSD's expected decrease per iteration

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

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

- 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 (SSD)** 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. \square

SSD's expected decrease per iteration

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

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

- 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 SSD 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:

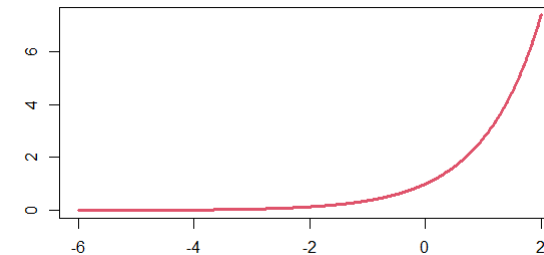
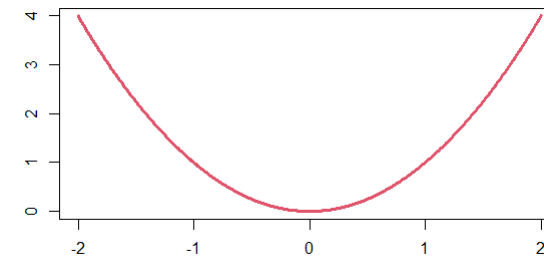
$$\begin{aligned} E[g(\mathbf{x}^{(t+1)})] - g(\mathbf{x}^{(t)}) &\leq -\alpha_t \|\mathbf{g}'(\mathbf{x}^{(t)})\|_2^2 + \alpha_t^2 \frac{L}{2} E \|\mathbf{g}'_{R(t)}(\mathbf{x}^{(t)})\|_2^2 \\ &\leq -\alpha_t \|\mathbf{g}'(\mathbf{x}^{(t)})\|_2^2 + \alpha_t^2 \frac{L}{2} (s + w \|\mathbf{g}'(\mathbf{x}^{(t)})\|_2^2)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

Strongly convex functions

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

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

- For one-dimensional functions:
 $(f'(x_1) - f'(x_2))/(x_1 - x_2) \geq m$ 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

- 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>

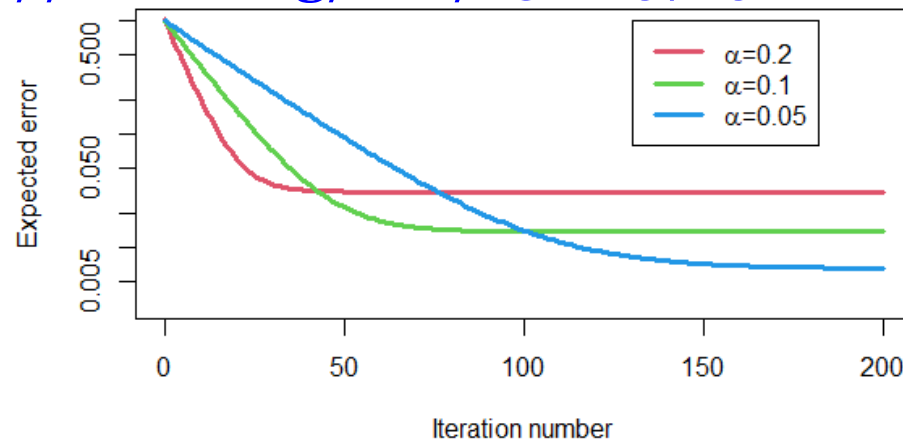
Convergence analysis for fixed step size

- 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:

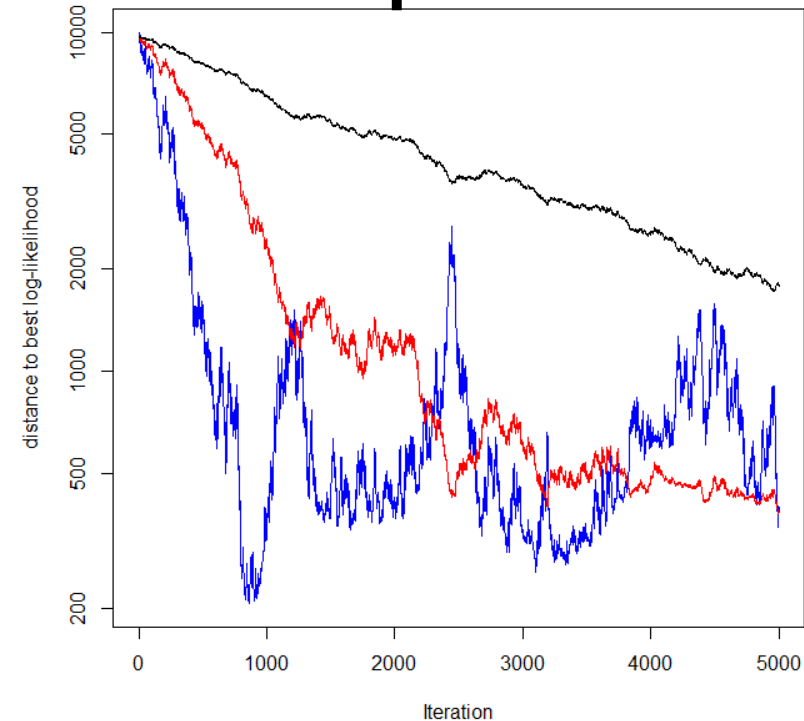
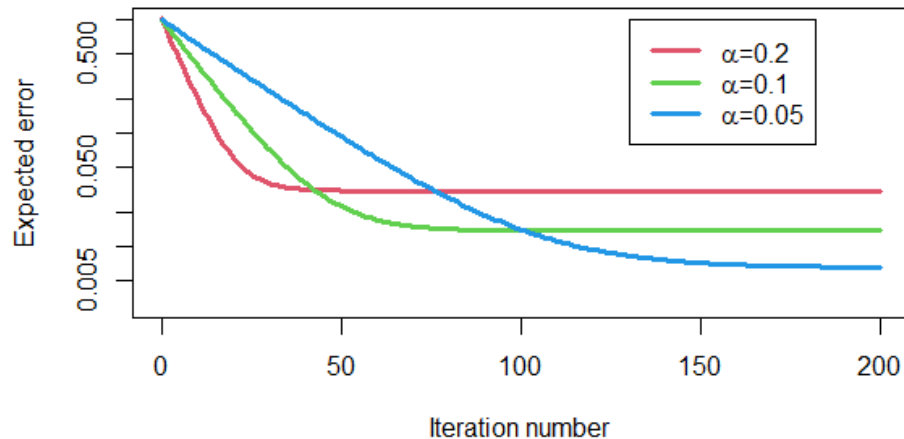
<https://arxiv.org/abs/1310.5715>



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 steepest descent: empirical examples

- 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 above bound for expected distance to optimum for $s=0.5$, $m=2$, $\max(L_i)=2$, $\epsilon_0 = \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 = 1$

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

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

means sublinear convergence since

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

- A bound like

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

would lead to linear convergence

- So, SSD with $\alpha_t = \frac{\beta}{t+\gamma}$ gives only sublinear convergence

Lipschitz continuous functions and matrix norms

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

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

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

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

SSD 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 regression):

a) $\|\mathbf{g}'_i(\mathbf{b})\|_2^2 =$

b) $E\|\mathbf{g}'_R(\mathbf{b})\|_2^2 =$

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

c) $E[\mathbf{g}'_R(\mathbf{b})] =$

d) $\|\mathbf{g}''(\mathbf{b})\|_{\text{Spectral}} =$

SSD 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 regression):

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

$$b) \quad E\|\mathbf{g}'_R(\mathbf{b})\|_2^2 = \frac{1}{n} \sum_i \|\mathbf{g}'_i(\mathbf{b})\|_2^2 = \dots$$

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

$$c) \quad E[\mathbf{g}'_R(\mathbf{b})] = \frac{1}{n} \sum_i \mathbf{g}'_i(\mathbf{b}) = \mathbf{g}'(\mathbf{b}) = \dots$$

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

Maximum likelihood estimator (MLE)

- The MLE is solution of $g(\hat{\boldsymbol{\beta}}) = \max g(\mathbf{b})$ with
 $g(\hat{\boldsymbol{\beta}}) = \log\text{-likelihood}(\hat{\boldsymbol{\beta}}, \mathbf{X}, \mathbf{y}) = \sum_{i=1}^n \log\text{-likelihood}(\hat{\boldsymbol{\beta}}, \mathbf{x}_i, y_i)$
(the latter equation requires independence of observations)
- In the simple case of normally distributed observations, MLE=LSE and we have an algebraic solution
- Otherwise, we need usually iterative methods to compute the MLE
- If the data is from an exponential family, the function g is concave ($-g$ is convex)