# Lecture 5: Steepest descent methods

Coralia Cartis, Mathematical Institute, University of Oxford

C6.2/B2: Continuous Optimization

# A generic linesearch method (Lecture 2)

(UP):   minimize $f(x)$ subject to $x \in \mathbb{R}^n$, where $f \in \mathcal{C}^1$ or $\mathcal{C}^2$.

**A Generic Linesearch Method (GLM)**

Choose $\epsilon > 0$ and $x^0 \in \mathbb{R}^n$. For $k \geq 0$, do:
While $\|\nabla f(x^k)\| > \epsilon$, REPEAT:

- compute a <u>descent</u> search direction $s^k \in \mathbb{R}^n$,

$$\nabla f(x^k)^T s^k < 0;$$

- compute a stepsize $\alpha^k > 0$ along $s^k$ such that

$$f(x^k + \alpha^k s^k) < f(x^k);$$

- set $x^{k+1} := x^k + \alpha^k s^k$ and $k := k + 1$.   □

Recall property of descent directions (Lemma 1, Lecture 1).

# Global convergence of GLM (Lecture 4)

**Theorem 4.** Let $f \in \mathcal{C}^1(\mathbb{R}^n)$ be bounded below on $\mathbb{R}^n$ by $f_{\mathrm{low}}$. Let $\nabla f$ Lipschitz continuous. Apply GLM with bArmijo linesearch to minimizing $f$ with $\epsilon := 0$. Then either

there exists $l \geq 0$ such that $\nabla f(x^l) = 0$

or

$$\lim_{k \to \infty} \min \left\{ \frac{|\nabla f(x^k)^T s^k|}{\|s^k\|}, |\nabla f(x^k)^T s^k| \right\} = 0. \qquad \text{(conv)}$$

Note that the limit (conv) is equivalent to

$$\lim_{k \to \infty} \|\nabla f(x^k)\| \cdot \cos \theta_k \cdot \min\{1, \|s^k\|\} = 0,$$

where $\cos \theta_k = \frac{(-\nabla f(x^k))^T s^k}{\|\nabla f(x^k)\| \cdot \|s^k\|}$.

# Steepest descent method

Steepest descent (SD) direction:  set $s^k := -\nabla f(x^k),\ \ k \geq 0,$ in Generic Linesearch Method (GLM).

# Steepest descent method

Steepest descent (SD) direction: set $s^k := -\nabla f(x^k),\ k \geq 0,$ in Generic Linesearch Method (GLM).

- $s^k$ <u>descent</u> direction whenever $\nabla f(x^k) \neq 0$:

$$\nabla f(x^k)^T s^k < 0 \iff \nabla f(x^k)^T(-\nabla f(x^k)) < 0 \iff -\|\nabla f(x^k)\|^2 < 0.$$

- $s^k$ <u>steepest</u> descent: unique global solution of

$$\text{minimize}_{s \in \mathbb{R}^n}\ f(x^k) + s^T \nabla f(x^k) \quad \text{subject to} \quad \|s\| = \|\nabla f(x^k)\|.$$

Cauchy-Schwarz: $|s^T \nabla f(x^k)| \leq \|s\| \cdot \|\nabla f(x^k)\|, \forall s,$ with equality iff $s$ is proportional to $\nabla f(x^k)$.

# Steepest descent methods

Method of steepest descent (SD): GLM with $s^k == SD$ direction; any linesearch.

**Steepest Descent (SD) Method**

Choose $\epsilon > 0$ and $x^0 \in \mathbb{R}^n$. While $\|\nabla f(x^k)\| > \epsilon$, REPEAT:

- ■ compute $s^k = -\nabla f(x^k)$.
- ■ compute a stepsize $\alpha^k > 0$ along $s^k$ such that

$$f(x^k + \alpha^k s^k) < f(x^k);$$

- ■ set $x^{k+1} := x^k + \alpha^k s^k$ and $k := k + 1$. □

- ■ SD-e :== SD method with exact linesearches;
- ■ SD-bA :== SD method with bArmijo linesearches.

# Global convergence of steepest descent methods

- $f \in \mathcal{C}^1(\mathbb{R}^n)$; $\nabla f$ is Lipschitz continuous (on $\mathbb{R}^n$) iff $\exists L > 0$,
$$\|\nabla f(y) - \nabla f(x)\| \leq L\|y - x\|, \quad \forall x, y \in \mathbb{R}^n.$$

Theorem 5    Let $f \in \mathcal{C}^1(\mathbb{R}^n)$ be bounded below on $\mathbb{R}^n$. Let $\nabla f$ be Lipschitz continuous. Apply the SD-e or the SD-bA method to minimizing $f$ with $\epsilon := 0$.
Then both variants of the SD method have the property:

either
$\qquad$ there exists $l \geq 0$ such that $\nabla f(x^l) = 0$

or
$$\|\nabla f(x^k)\| \to 0 \text{ as } k \to \infty.$$

**Proof for SD-bA.**    Let $s^k = -\nabla f(x^k)$ for all $k$ in Th 4.    □

SD methods have excellent global convergence properties (under weak assumptions).

# Some disadvatanges of steepest descent methods

- SD methods are scale-dependent.

  poorly scaled problem/variables $\Longrightarrow$ SD direction gives little progress.

- Usually, SD methods converge very slowly to solution, asymptotically.

# The scale-dependence of steepest descent

Example of a poorly scaled quadratic.

$$f(x) = \frac{1}{2}(ax_1^2 + x_2^2) = \frac{1}{2}x^T \begin{pmatrix} a & 0 \\ 0 & 1 \end{pmatrix} x, \quad x = (x_1 \ x_2)^T, \quad (\Diamond)$$

where $a > 0$. Note $x^* = (0 \ 0)^T$ unique global minimizer.
- $a \gg 1 \longrightarrow \quad f$ poorly scaled (or poorly conditioned).

# The scale-dependence of steepest descent

Example of a poorly scaled quadratic.

$$f(x) = \frac{1}{2}(ax_1^2 + x_2^2) = \frac{1}{2}x^T \begin{pmatrix} a & 0 \\ 0 & 1 \end{pmatrix} x, \quad x = (x_1 \ x_2)^T, \quad (\diamond)$$

where $a > 0$. Note $x^* = (0 \ 0)^T$ unique global minimizer.

■ $a \gg 1 \ \longrightarrow \ f$ poorly scaled (or poorly conditioned).

■ apply SD-e to ($\diamond$) starting at $x^0 := (1 \ a)^T$. Then[see Pb Sheet 2]

$$x^k = \left(\frac{a-1}{a+1}\right)^k \begin{pmatrix} (-1)^k \\ a \end{pmatrix}, \quad k \geq 0.$$

# The scale-dependence of steepest descent

Example of a poorly scaled quadratic.

$$f(x) = \frac{1}{2}(ax_1^2 + x_2^2) = \frac{1}{2}x^T \begin{pmatrix} a & 0 \\ 0 & 1 \end{pmatrix} x, \quad x = (x_1 \ x_2)^T, \quad (\diamond)$$

where $a > 0$. Note $x^* = (0 \ 0)^T$ unique global minimizer.

- $a \gg 1 \longrightarrow f$ poorly scaled (or poorly conditioned).
- apply SD-e to $(\diamond)$ starting at $x^0 := (1 \ a)^T$. Then[see Pb Sheet 2]

$$x^k = \left(\frac{a-1}{a+1}\right)^k \begin{pmatrix} (-1)^k \\ a \end{pmatrix}, \quad k \geq 0.$$

$\implies x^k \to 0$ as $k \to \infty$, linearly with $\rho := |(a-1)/(a+1)|$ convergence factor.

- $a \gg 1 \implies \rho$ closer to $1 \implies$ SD-e converges very slowly.
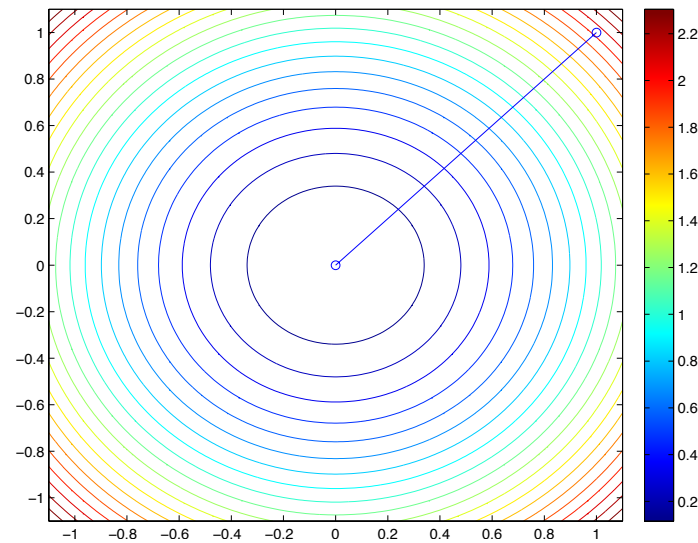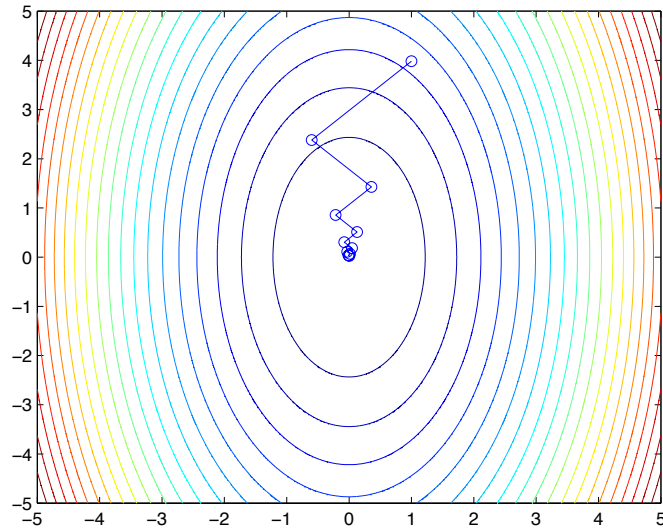
# The scale-dependence of steepest descent

Example of a well-scaled quadratic.

Linear transformation of variables:

$$y = \begin{pmatrix} a^{1/2} & 0 \\ 0 & 1 \end{pmatrix} x.$$

- let $\overline{f}(y) := f(x(y))$, namely $f$ in the new coordinates $y$.

$\implies \overline{f}(y) = \frac{1}{2} y^T y = \frac{1}{2}(y_1^2 + y_2^2).$

$\longrightarrow \quad \overline{f}$ well-scaled.

- $y^* = (0\ 0)^T$ unique global minimizer.

- apply SD-e to $\overline{f}$ from any $y^0 \in \mathbb{R}^2$: $y^1 = (0\ 0)^T = y^*$.

# The scale-dependence of steepest descent



The effect of problem scaling on SD-e performance.
Left figure: $a = 10^{0.6}$ (mildly poor scaling).
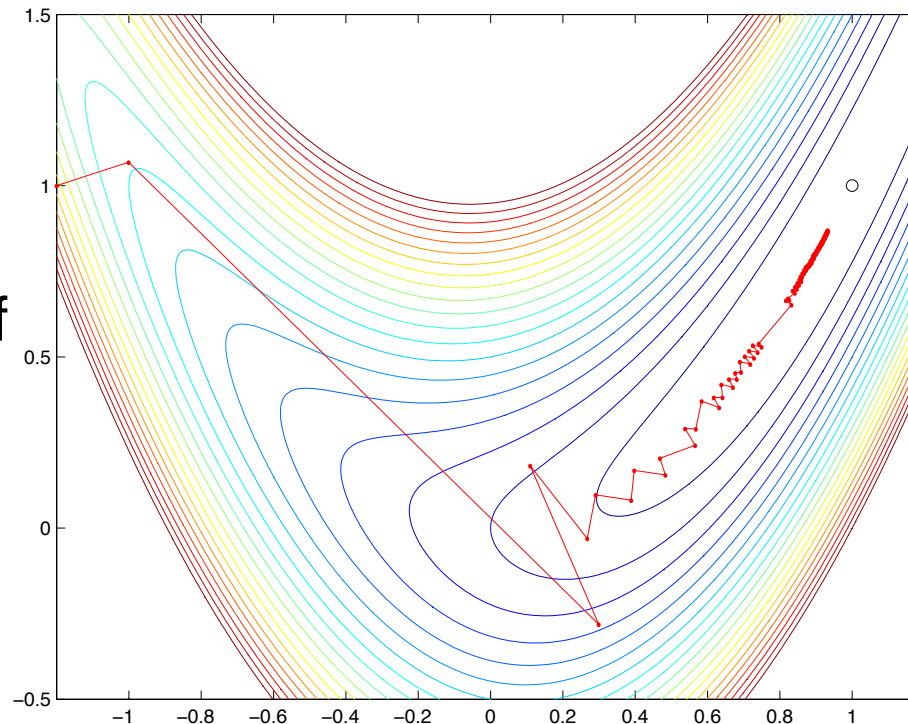Right figure: $a = 1$ ("perfect" scaling).

# Local rate of convergence for steepest descent

- ■ Usually, SD methods converge very slowly to solution, asymptotically.

theory: very slow conv.

numerics: break-down (cumulation of round-off and ill-conditioning).

$$f(x_1, x_2) = 10(x_2 - x_1^2)^2 + (x_1 - 1)^2.$$



SD-bA applied to the Rosenbrock function $f$.

# Local rate of convergence for steepest descent

Asymptotically, SD converges <u>linearly</u> to a solution. Namely, if $x^k \to x^*$, as $k \to \infty$, then

$$\|x^{k+1} - x^*\| \le \rho \|x^k - x^*\|, \ \forall k \text{ suff. large}$$

BUT

# Local rate of convergence for steepest descent

Asymptotically, SD converges <u>linearly</u> to a solution. Namely, if $x^k \to x^*$, as $k \to \infty$, then
$$\|x^{k+1} - x^*\| \le \rho \|x^k - x^*\|, \; \forall k \text{ suff. large}$$

BUT  convergence factor $\rho$ v. close to $1$ usually!

# Local rate of convergence for steepest descent

Asymptotically, SD converges <u>linearly</u> to a solution. Namely, if $x^k \to x^*$, as $k \to \infty$, then
$$\|x^{k+1} - x^*\| \le \rho \|x^k - x^*\|, \ \forall k \text{ suff. large}$$

BUT  convergence factor $\rho$ v. close to $1$ usually!

Theorem 6  $f \in \mathcal{C}^2$; $x^*$ local minimizer of $f$ with $\nabla^2 f(x^*)$ positive definite $\longrightarrow \lambda^*_{\max}, \lambda^*_{\min}$ eigenvalues.
Apply SD-e to $\min f$. If $x^k \to x^*$ as $k \to \infty$, then $x^k$ converges linearly to $x^*$
$$\rho \le \frac{\kappa(x^*) - 1}{\kappa(x^*) + 1} := \rho_{SD},$$
where $\kappa(x^*) = \lambda^*_{\max}/\lambda^*_{\min}$ condition number of $\nabla^2 f(x^*)$.

# Local rate of convergence for steepest descent

Asymptotically, SD converges <u>linearly</u> to a solution. Namely, if $x^k \to x^*$, as $k \to \infty$, then

$$\|x^{k+1} - x^*\| \le \rho\|x^k - x^*\|, \ \forall k \text{ suff. large}$$

BUT  convergence factor $\rho$ v. close to $1$ usually!

Theorem 6   $f \in \mathcal{C}^2$; $x^*$ local minimizer of $f$ with $\nabla^2 f(x^*)$ positive definite $\longrightarrow \lambda^*_{\max}, \lambda^*_{\min}$ eigenvalues.
Apply SD-e to $\min f$. If $x^k \to x^*$ as $k \to \infty$, then $x^k$ converges linearly to $x^*$

$$\rho \le \frac{\kappa(x^*)-1}{\kappa(x^*)+1} := \rho_{SD},$$

where $\kappa(x^*) = \lambda^*_{\max}/\lambda^*_{\min}$ condition number of $\nabla^2 f(x^*)$.

• practice: $\rho = \rho_{SD}$;
for Rosenbrock $f$: $\kappa(x^*) = 258.10$, $\rho_{SD} \approx 0.992$.

# Summary: steepest descent methods

- first-order method $\longrightarrow$ inexpensive.

- global convergence under weak assumptions, but no second-order optimality guarantees for the generated solution.

- scale-dependent; too expensive, or impossible, to make a function well-scaled.

- when the objective is poorly scaled, very very slow convergence to a solution; hence, not used in general.

- useful sometimes: for example, for some convex problems with special structure that are very well conditioned (compressed sensing, etc).