# Math 273a: Optimization
# Gradient descent

Instructor: Wotao Yin
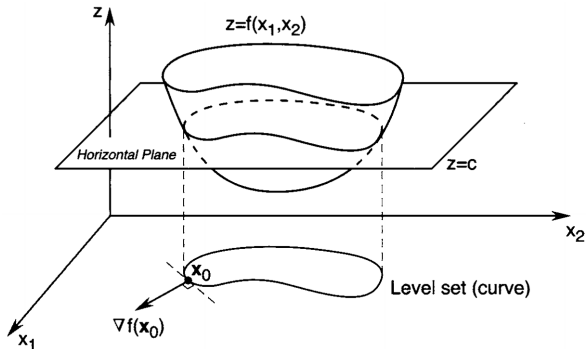
Department of Mathematics, UCLA

Fall 2015

slides based on Chong-Zak, 4th Ed.

online discussions on piazza.com

# Main features of gradient methods

- The most popular methods (in continuous optimization)

- simple and intuitive

- work under very few assumptions

  (although they cannot directly handle nondifferentiable objectives and constraints, without applying smoothing techniques)

- work together with many other methods: *duality*, *splitting*, *coordinate descent*, *alternating direction*, *stochastic*, *online*, etc.

- suitable for large-scale problems, e.g., easy to parallelize for problems with many terms in the objective

# Gradients

- We let $\nabla f(x_0)$ denote the gradient of $f$ at point $x_0$.

- $\nabla f(x_0) \perp$ tangent of the levelset curve of $f$ passing $x_0$, pointing outward (recall: level set $\mathcal{L}_f(c) := \{x : f(x) = c\}$)

- $\nabla f(x_0)$ is <u>max-rate ascending direction</u> of $f$ at $x_0$ (for a small displacement), and $\|\nabla f(x_0)\|$ is the rate.

  **Reason:** pick any direction $d$ with $\|d\| = 1$. The rate of change at $x$ is

  $$\langle \nabla f(x), d \rangle \leq \|\nabla f(x)\| \cdot \|d\| = \|\nabla f(x)\|.$$

  If we set $d = \nabla f(x)/\|\nabla f(x)\|$, then

  $$\langle \nabla f(x), d \rangle = \|\nabla f(x)\|.$$

- Therefore, $-\nabla f(x)$ is the <u>max-rate descending direction</u> of $f$ and a <u>good search direction</u>.

## A negative gradient step can decrease the objective

- Let $x^{(0)}$ be any initial point.

- First-order Taylor expansion for candidate point $x = x^{(0)} - \alpha \nabla f(x^{(0)})$:

$$f(x) - f(x^{(0)}) = -\alpha \|\nabla f(x^{(0)})\|^2 + o(\alpha)$$

- Hence, if $\nabla f(x^{(0)}) \neq 0$ (the first-order necessary condition is not met) and $\alpha$ is <u>sufficiently small</u>, we have

$$f(x) < f(x^{(0)}).$$

- Therefore, for sufficiently small $\alpha$, $x$ is an improvement over $x^{(0)}$.

## Gradient descent algorithm

- Given initial $\boldsymbol{x}^{(0)}$, the gradient descent algorithm uses the following update to generate $\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \ldots$, until a stopping condition is met:
  from the current point $\boldsymbol{x}^{(k)}$, generate the next point $\boldsymbol{x}^{(k+1)}$ by

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha_k \nabla f(\boldsymbol{x}^{(k)}),$$

- $\alpha_k$ is called the <u>step size</u>

**Alternative interpretation**:

- notice that

$$\boldsymbol{x}^{(k+1)} = \arg\min_{\boldsymbol{x}} \frac{1}{2\alpha_k} \left\| \boldsymbol{x} - \left( \boldsymbol{x}^{(k)} - \alpha_k \nabla f(\boldsymbol{x}^{(k)}) \right) \right\|^2$$

$$= \arg\min_{\boldsymbol{x}} f(\boldsymbol{x}^{(k)}) + \langle \nabla f(\boldsymbol{x}^{(k)}), \boldsymbol{x} - \boldsymbol{x}^{(k)} \rangle + \frac{1}{2\alpha_k} \|\boldsymbol{x} - \boldsymbol{x}^{(k)}\|^2$$

  (2nd "=" follows from that adding constants and multiplying a positive constant do *not* change the set of minimizers or "$\arg\min$")

- Hence, $\boldsymbol{x}^{(k+1)}$ is obtained by minimizing the <u>linearization</u> of $f$ at $\boldsymbol{x}^{(k)}$ and <u>a proximal term</u> that keeps $x^{k+1}$ close to $\boldsymbol{x}^{(k)}$.

- The reformulation is useful to develop the extensions of gradient descent:
    - projected gradient method
    - proximal-gradient method
    - accelerated gradient method
    - ......

## When to stop the iteration

The first-order necessary condition $\|\nabla f(\boldsymbol{x}^{(k+1)})\| = 0$ is not practical.

Practical conditions:

- gradient condition $\|\nabla f(\boldsymbol{x}^{(k+1)})\| < \epsilon$
- successive objective condition $|f(\boldsymbol{x}^{(k+1)}) - f(\boldsymbol{x}^{(k)})| < \epsilon$ or the relative one

$$\frac{|f(\boldsymbol{x}^{(k+1)}) - f(\boldsymbol{x}^{(k)})|}{|f(\boldsymbol{x}^{(k)})|} < \epsilon$$

- successive point difference $\|\boldsymbol{x}^{(k+1)} - \boldsymbol{x}^{(k)}\| < \epsilon$ or the relative one

$$\frac{\|\boldsymbol{x}^{(k+1)} - \boldsymbol{x}^{(k)}\|}{\|\boldsymbol{x}^{(k)}\|} < \epsilon$$

- to avoid division by tiny numbers (unstable division), we can replace the denominators by $\max\{1, |f(\boldsymbol{x}^{(k)})|\}$ and $\max\{1, \|\boldsymbol{x}^{(k)}\|\}$, respectively

# Small versus large step sizes $\alpha_k$

**Small step size:**

- Pros: iterations are more likely converge, closely traces max-rate descends
- Cons: need more iterations and thus evaluations of $\nabla f$

**Large step size:**

- Pros: better use of each $\nabla f(x^{(k)})$, may reduce the total iterations
- Cons: can cause overshooting and zig-zags, too large $\Rightarrow$ diverged iterations

# Small versus large step sizes $\alpha_k$

**Small step size:**

- Pros: iterations are more likely converge, closely traces max-rate descends
- Cons: need more iterations and thus evaluations of $\nabla f$

**Large step size:**

- Pros: better use of each $\nabla f(x^{(k)})$, may reduce the total iterations
- Cons: can cause overshooting and zig-zags, too large $\Rightarrow$ diverged iterations
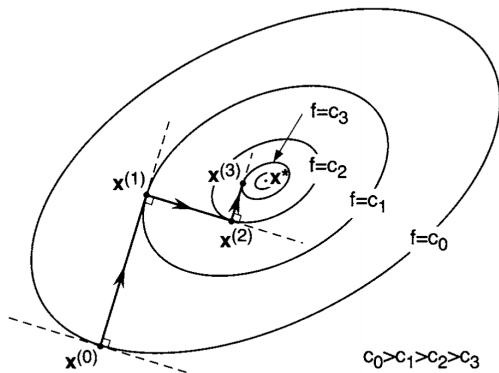
In practice, step sizes are often chosen

- as a fixed value if $\nabla f$ is Lipschitz (rate of change is bounded) with the constant known or an upper bound of it known
- by line search
- by a method called Barzilai-Borwein with nonmonotone line search

## Steepest descent method
## (gradient descent with exact line search)

Step size $\alpha_k$ is determined by exact minimization

$$\alpha_k = \arg\min_{\alpha \geq 0} \ f(\boldsymbol{x}^{(k)} - \alpha \nabla f(\boldsymbol{x}^{(k)})).$$

It is used mostly for quadratic programs (with $\alpha_k$ in a closed form) and some problems with inexpensive evaluation values but expensive gradient evaluation; otherwise it is not worth the effort to solve this subproblem exactly.

Labels in figure: $f = c_3$, $f = c_2$, $f = c_1$, $f = c_0$, $\mathbf{x}^{(0)}$, $\mathbf{x}^{(1)}$, $\mathbf{x}^{(2)}$, $\mathbf{x}^{(3)}$, $\mathbf{x}^*$, $c_0 > c_1 > c_2 > c_3$

**Proposition 8.1** *If $\{\boldsymbol{x}^{(k)}\}_{k=0}^{\infty}$ is a steepest descent sequence for a given function $f : \mathbb{R}^n \to \mathbb{R}$, then for each $k$ the vector $\boldsymbol{x}^{(k+1)} - \boldsymbol{x}^{(k)}$ is orthogonal to the vector $\boldsymbol{x}^{(k+2)} - \boldsymbol{x}^{(k+1)}$.* $\square$

## Steepest descent for quadratic programming

**Assume** that $Q$ is symmetric and positive definite ($x^T Q x > 0$ for any $x \neq 0$).

Consider the quadratic program

$$f(x) = \frac{1}{2} x^T Q x - b^T x$$

with

$$\nabla f(x) = Q x - b.$$

**Steepest descent iteration:** start from any $x^{(0)}$, set

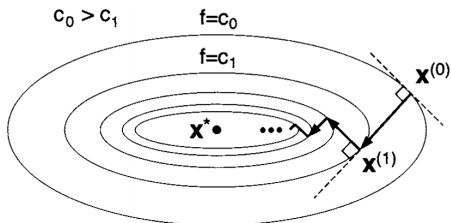$$x^{(k+1)} = x^{(k)} - \alpha_k g^{(k)}, \quad k = 0, 1, 2, \dots$$

where $g^{(k)} := \nabla f(x^{(k)})$ and

$$
\begin{aligned}
\alpha_k &= \arg\min_{\alpha \geq 0} \; f(x^{(k)} - \alpha g^{(k)}) \\
&= \frac{g^{(k)T} g^{(k)}}{g^{(k)T} Q g^{(k)}}.
\end{aligned}
$$

# Examples

**Example 1:** $f(\boldsymbol{x}) = x_1^2 + x_2^2$. Steepest descent arrives at $\boldsymbol{x}^* = \boldsymbol{0}$ in 1 iteration.

**Example 1:** $f(\boldsymbol{x}) = \frac{1}{5}x_1^2 + x_2^2$. Steepest descent makes progress in a narrow valley

# Performance of steepest descent

- **Per-iteration cost:** dominated by *two* matrix-vector multiplications:
    - $g^{(k)} = Qx^{(k)} - b$
    - computing $\alpha_k$ involves $Qg^{(k)}$

  but they can be easily reduced to *one* matrix-vector multiplication.

- **Convergence speed:** determined by the initial point and the spectral condition of $Q$. To analyze them, we
    - define solution error: $e^{(k)} = x^{(k)} - x^*$ (not known, an analysis tool)
    - have property: $g^{(k)} = Qx^{(k)} - b = Qe^{(k)}$.

**Good cases**:

- $e^{(k)}$ is an eigenvector of $Q$ with eigenvalue $\lambda$

$$e^{(k+1)} = e_k - \alpha_k g^{(k)} = e^{(k)} - \frac{g^{(k)T}g^{(k)}}{g^{(k)T}Qg^{(k)}}(Qe^{(k)})$$

$$= e^{(k)} + \frac{g^{(k)T}g^{(k)}}{\lambda\, g^{(k)T}g^{(k)}}(-\lambda e^{(k)}) = 0.$$

- $Q$ has only one distinct eigenvalue (the level sets of $Q$ are circles)

**The general case:** define $\|e\|_A := \sqrt{e^T A e}$ and $\kappa := \lambda_{\max}(Q)/\lambda_{\min}(Q)$, then we have

$$\|e^{(k)}\|_A \leq \left(\frac{\kappa - 1}{\kappa + 1}\right)^k \|e^{(0)}\|_A.$$

A example from *An Introduction to CG method* by Shewchuk

## Gradient descent with <u>fixed</u> step size

- Iteration:
$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha \boldsymbol{g}^{(k)}$$

- We assume that $\boldsymbol{x}^*$ exists

- Check distance to solution:
$$\begin{aligned}
\|\boldsymbol{x}^{(k+1)} - \boldsymbol{x}^*\|^2 &= \|\boldsymbol{x}^{(k)} - \boldsymbol{x}^* - \alpha \boldsymbol{g}^{(k)}\|^2 \\
&= \|\boldsymbol{x}^{(k)} - \boldsymbol{x}^*\|^2 - 2\alpha \langle \boldsymbol{g}^{(k)}, \boldsymbol{x}^{(k)} - \boldsymbol{x}^* \rangle + \alpha^2 \|\boldsymbol{g}^{(k)}\|^2.
\end{aligned}$$

- Therefore, in order to have $\|\boldsymbol{x}^{(k+1)} - \boldsymbol{x}^*\| \leq \|\boldsymbol{x}^{(k)} - \boldsymbol{x}^*\|$, we must have
$$\frac{\alpha}{2} \|\boldsymbol{g}^{(k)}\|^2 \leq \langle \boldsymbol{g}^{(k)}, \boldsymbol{x}^{(k)} - \boldsymbol{x}^* \rangle.$$

Since $\boldsymbol{g}^* := \nabla f(\boldsymbol{x}^*) = \boldsymbol{0}$, the condition is equivalent to
$$\frac{\alpha}{2} \|\boldsymbol{g}^{(k)} - \boldsymbol{g}^*\|^2 \leq \langle \boldsymbol{g}^{(k)} - \boldsymbol{g}^*, \boldsymbol{x}^{(k)} - \boldsymbol{x}^* \rangle.$$

## Special case: convex and Lipschitz differentiable $f$

- **Definition:** A function $f$ is $L$-<u>Lipschitz differentiable</u>, $L \geq 0$, if $f \in \mathcal{C}^1$ and

$$\|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y})\| \leq L\|\boldsymbol{x} - \boldsymbol{y}\|, \quad \forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^n$$

  (the maximum rate of change of $\nabla f$ is $L$)

- **Baillon-Haddad theorem**: if $f \in \mathcal{C}^1$ is a convex function, then it is $L$-Lipschitz differentiable if and only if

$$\|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y})\|^2 \leq L\langle \nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y}), \boldsymbol{x} - \boldsymbol{y}\rangle.$$

  (such $\nabla f$ is called $1/L$-cocoercive)

- **Theorem:** Let $f \in \mathcal{C}^1$ be a convex function and $L$-Lipschitz differentiable. If $0 < \alpha \le 2/L$, then

$$\frac{\alpha}{2}\|\boldsymbol{g}^{(k)} - \boldsymbol{g}^*\|^2 \le \langle \boldsymbol{g}^{(k)} - \boldsymbol{g}^*, \boldsymbol{x}^{(k)} - \boldsymbol{x}^* \rangle$$

and thus $\|\boldsymbol{x}^{(k+1)} - \boldsymbol{x}^*\| \le \|\boldsymbol{x}^{(k)} - \boldsymbol{x}^*\|$ for $k = 0, 1 \ldots$. The iteration stays bounded.

- **Theorem:** Let $f \in \mathcal{C}^1$ be a convex function and $L$-Lipschitz differentiable. If $0 < \alpha < L/2$, then
    - both $f(\boldsymbol{x}^{(k)})$ and $\|\nabla f(\boldsymbol{x}^{(k)})\|$ are monotonically decreasing,
    - $f(\boldsymbol{x}^{(k)}) - f(\boldsymbol{x}^{(*)}) = O(\frac{1}{k})$,
    - $\|\nabla f(\boldsymbol{x}^{(k)})\| = o(\frac{1}{k})$.
      (one often writes $\|\nabla f(\boldsymbol{x}^{(k)})\|^2 = o(\frac{1}{k^2})$ since $\|\nabla f(\boldsymbol{x}^{(k)})\|^2$ naturally appears in most analysis.)

## Gradient descent with fixed step size
## for quadratic programming

**Assume** that $Q$ is symmetric and positive definite ($x^T Q x > 0$ for any $x \neq 0$).

Consider the quadratic program

$$f(x) = \frac{1}{2} x^T Q x - b^T x$$

**Theorem 8.3** *For the fixed-step-size gradient algorithm, $x^{(k)} \to x^*$ for any $x^{(0)}$ if and only if*

$$0 < \alpha < \frac{2}{\lambda_{\max}(Q)}.$$

# Summary

- Negative gradient $-\nabla f(\boldsymbol{x}^{(k)})$ is the max-rate descending direction

- For some small $\alpha_k$, $\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - \alpha_k \nabla f(\boldsymbol{x}^{(k)})$ improves over $\boldsymbol{x}^{(k)}$

- There are practical rules to determine when to stop the iteration

- Exact line search works for quadratic program with $\boldsymbol{Q} > 0$. Zig-zag occurs if $\boldsymbol{x}^{(0)} - \boldsymbol{x}^*$ is away from an eigenvector and spectrum of $\boldsymbol{Q}$ is spread

- Fixed step gradient descent works for convex and Lipschitz-differentiable $f$

- To keep the discussion short and informative, we have omitted much other convergence analysis.