# ECE 8101: Nonconvex Optimization for Machine Learning

### Lecture Note 2-3: Gradient Descent

Jia (Kevin) Liu

Associate Professor
Department of Electrical and Computer Engineering
The Ohio State University, Columbus, OH, USA

Autumn 2024

# Outline

In this lecture:

- Convergence rate concept

- Gradient descent method

- Convergence performance of gradient descent

- Step size selection strategies

# Iterative Algorithms for Optimization

We consider the following iterative algorithms:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + s_k \mathbf{d}_k,$$

where $s_k$ is step-size, and $\mathbf{d}_k$ is search direction depending on $(\mathbf{x}_k, \mathbf{x}_{k-1}, \ldots)$.

For now: assume $f$ smooth, $f(\mathbf{x}_k)$ and $\nabla f(\mathbf{x}_k)$ is easy to evaluate

Complications from ML:

- Nonconvex $f$
- Nonsmooth $f$
- $f$ not available (or too expensive to evaluate exactly)
- Only an estimate of $\nabla f(\mathbf{x}_k)$ is available
- A constraint $\mathbf{x} \in \Omega$ (usually a relatively simple $\Omega$, e.g., ball, box, simplex...)
- Nonsmooth regularization, i.e., instead of $f(\mathbf{x})$, we want $\min f(\mathbf{x}) + \tau \psi(\mathbf{x})$

# How to Evaluate the Speed of an Iterative Algorithm?
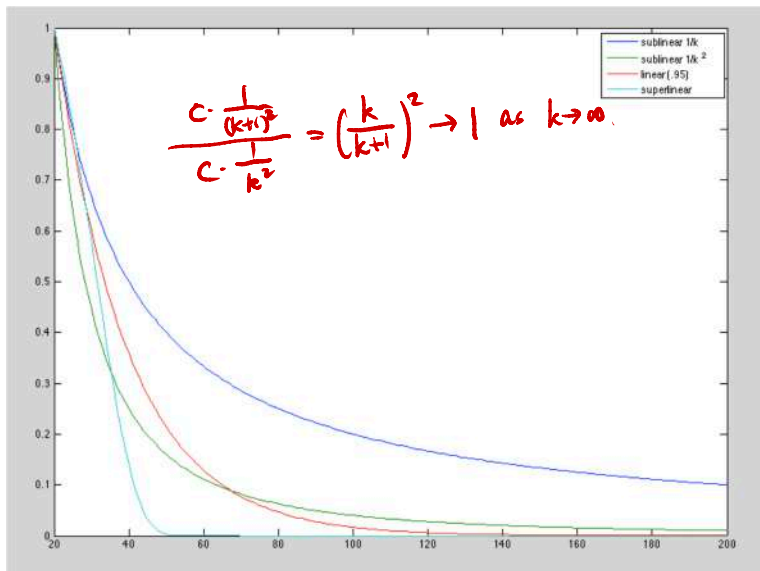
## Definition 1 (Convergence rate)

A sequence $\{r_k\} \to r^*$ and $r_k \neq r^*$ for all $k$. The rate (or order) of convergence $p$ is a nonnegative number satisfying

$$\limsup_{k \to \infty} \frac{\|r_{k+1} - r^*\|}{\|r_k - r^*\|^p} = \beta < \infty.$$

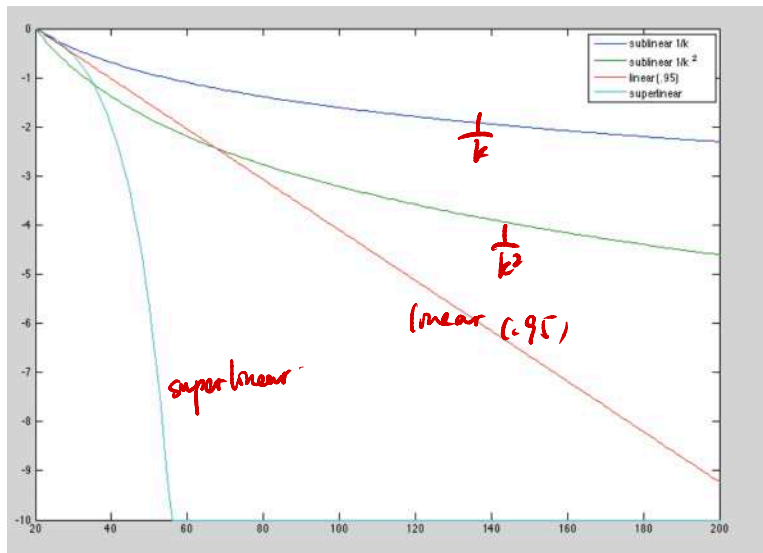$\frac{\|r_{k+1} - r^*\|}{\|r_k - r^*\|} \to 1$ as $k \to \infty$.

- Sublinear: $p = 1$ and $\beta = 1$ (e.g., $O(1/k)$ rate, kind of slow but still OK)

  $\frac{C/(k+1)}{C/k} = \frac{k}{k+1} \to 1$. Desired $\varepsilon > 0$, $\frac{C}{k} \leq \varepsilon \Rightarrow k \geq \frac{C}{\varepsilon} = O(\frac{1}{\varepsilon})$   contraction.

- Linear or geometric: $p = 1$ and $0 < \beta < 1$ (i.e., $\|r_{k+1} - r^*\| \leq \beta\|r_k - r^*\|$ $\leq \beta^2 \|r_{k-1} - r^*\|$
  for some $\beta \in (0,1)$, or $\|r_k - r^*\| = O(\beta^k)$, which is quite fast) $\leq \beta^k \|r_1 - r^*\| = O(\beta^k)$

  Desired $\varepsilon$: $C\beta^k \leq \varepsilon \Rightarrow k \geq C\log(\varepsilon^{-1})$. Need $O(\log(\varepsilon^{-1}))$ iter.

- Superlinear: $p > 1$ and $\beta < \infty$, or $p = 1$ and $\beta = 0$ (i.e., $\frac{\|r_{k+1} - r^*\|}{\|r_k - r^*\|} \to 0$, $= O(\beta^k)$
  that's very fast!) Not only a contraction, but also the rate of convergence
  is decelerating

- Quadratic: $p = 2$ and $\beta < \infty$ ($\|r_{k+1} - r^*\| \leq \beta\|r_k - r^*\|^2$, # of correct
  significant digits doubles per iteration. Rarely need anything faster than this!)

For $\varepsilon$-accuracy: Need $O(\log\log(\varepsilon^{-1}))$ iter. $\leftarrow$ almost const.

# Convergence Rates Comparisons



$$\frac{C \cdot \frac{1}{(k+1)^2}}{C \cdot \frac{1}{k^2}} = \left(\frac{k}{k+1}\right)^2 \to 1 \quad \text{as} \quad k \to \infty.$$

# Convergence Rates Comparisons: Log-Scale

# Gradient Descent

Back to the unconstrained optimization problem, with $f$ smooth and convex:

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

Denote the optimal value as $f^* = \min_{\mathbf{x}} f(\mathbf{x}^*)$ and an optimal solution as $\mathbf{x}^*$

---

### Gradient Descent

Choose initial point $\mathbf{x}_0 \in \mathbb{R}^n$. Repeat:

$$\mathbf{x}_k = \mathbf{x}_{k-1} - s_k \nabla f(\mathbf{x}_{k-1}), \quad k = 1, 2, 3, \ldots$$

Stop if some stopping criterion is satisfied. $\quad$ e.g., $\|\nabla f(z_k)\| \leq \varepsilon$

$$\|z_{k+1} - z_k\| \leq \varepsilon.$$

stop after fixed # of iter.
(finite-time conv. analysis).

---

# Gradient Descent: Geometric Interpretation

Gradient descent is a first-order method: Consider the following quadratic Taylor approximation:

**SO - approx**

$$f(\mathbf{y}) \approx \underbrace{f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y}-\mathbf{x})}_{\text{FO-approx}} + \frac{1}{2}(\mathbf{y}-\mathbf{x})^\top \nabla^2 f(\mathbf{x})(\mathbf{y}-\mathbf{x}) + o(\|\mathbf{y}-\mathbf{x}\|^2)$$

$$\frac{1}{s}\mathbf{I}$$

No, we replace Hessian $\nabla^2 f(\mathbf{x})$ by $\frac{1}{s}\mathbf{I}$ to obtain:

$$f(\mathbf{y}) \approx \underbrace{f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y}-\mathbf{x})}_{\text{FO-approx}} + \boxed{\frac{1}{2s}\|\mathbf{y}-\mathbf{x}\|^2}$$
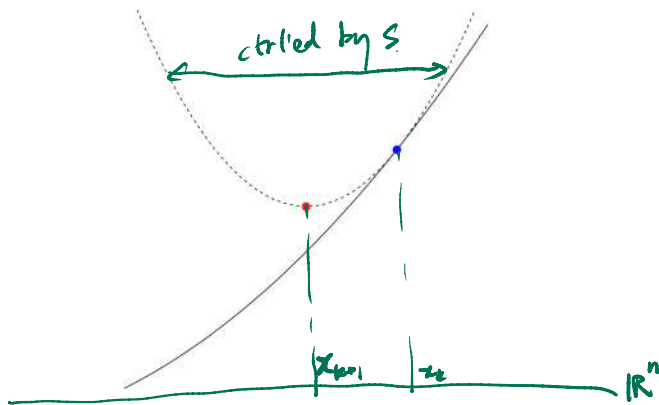
**"proximity penalty"**

Can be viewed as a linear approximation to $f$, with proximity term to $\mathbf{x}$ weighted by $\frac{1}{2s}$. Choose next point $\mathbf{y} = \mathbf{x}^+$ to minimize this approximation:

$$\mathbf{x}^+ = \mathbf{x} - s\nabla f(\mathbf{x})$$

**Quadratic fn of y. set grad → 0, solve for y**

$$\nabla f(y) = 0 \rightarrow \nabla f(\mathbf{x}) + \frac{1}{s}(y-\mathbf{x}) \Rightarrow \boxed{y = \mathbf{x} - s\nabla f(\mathbf{x})}$$

# Gradient Descent: Geometric Interpretation



$$\mathbf{x}^+ = \arg\min_y f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2s}\|\mathbf{y} - x\|_2^2$$
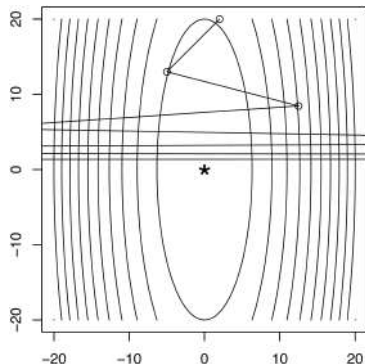
Questions:

- How to choose step sizes $\{s_k\}$?
- What is the according convergence rate? Or does it depend on $\{s_k\}$?
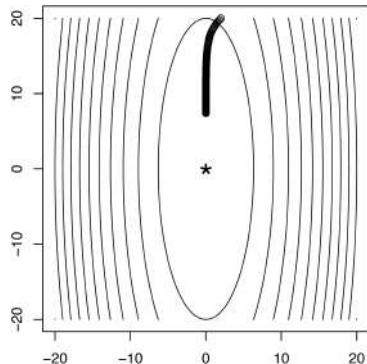
# Strategy 1: Fixed Step Size

Simply set $s_k = s$ for all $k = 1, 2, 3, \ldots$.

Limitations: May diverge if $s$ is too large, Can be slow if $s$ is too small.

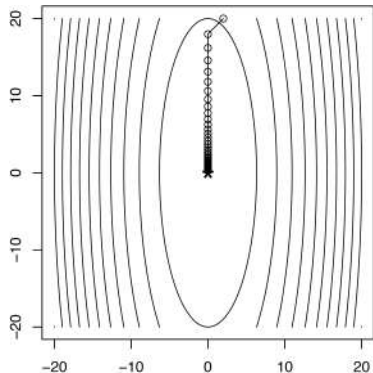Example: Consider $f(\mathbf{x}) = (10x_1^2 + x_2^2)/2$: $\Rightarrow (x_1^*, x_2^*) = (0, 0)$



8 iterations

100 iterations

# Strategy 1: Fixed Step Size

Converges nicely when $s$ is "just right." Same example, GD after 40 iterations:



Will be clear what we mean by "just right" in convergence rate analysis later

# Convergence Rate Analysis (Convex): Fixed Step Size

Assume that $f$ is convex & differentiable, with $\mathrm{dom}(f) = \mathbb{R}^n$ and additionally

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 \le L\|\mathbf{y} - \mathbf{x}\|_2, \quad \forall \mathbf{x}, \mathbf{y}$$

*L-smooth.*

That is, $\nabla f$ is Lipschitz continuous with constant $L > 0$ ($L$-Lipschitz continuous)

$h : D \subseteq \mathbb{R}^n$. $h$ is Lip. cont. $\exists L > 0$ s.t. $\|h(y) - h(x)\| \le L \|y - x\|$.

## Theorem 1 (Optimality Gap)

*If $f$ is convex, differentiable, and $L$-smooth, gradient descent with fixed step size $s \le 1/L$ satisfies*

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \le \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2sk},$$

*i.e., gradient descent method has sublinear convergence rate $O(1/k)$.*

Remark:
- To get $f(\mathbf{x}_k) - f(\mathbf{x}^*) \le \epsilon$, it takes $O(1/\epsilon)$ iterations.

> ## Theorem 1 (Optimality Gap)
>
> If $f$ is convex, differentiable, and $L$-smooth, gradient descent with fixed step size $s \leq 1/L$ satisfies
>
> $$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2sk}, \quad = O\left(\frac{1}{k}\right)$$
>
> i.e., gradient descent method has *sublinear* convergence rate $O(1/k)$.

Proof. Step ① Claim : If $\nabla f$ is Lipshitz.. then — descent lemma

$$f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2}\|y-x\|^2, \quad \forall x, y \in \mathbb{R}^n. \quad (1)$$

To show (1), we consider — dir. der.

$$f(y) = f(x) + \int_0^1 \underline{f'(x + \tau(y-x))}\, d\tau$$

$$= f(x) + \int_0^1 \nabla f(x + \tau(y-x))^T (y-x)\, d\tau \quad \text{(chain rule)}$$

add & sub
$$\overset{\text{add}}{=} f(x) + \int_0^1 \left[ \left[\nabla f(x+\tau(y-x))\right]^T (y-x) + \nabla f(x)^T(y-x) - \nabla f(x)^T(y-x) \right] d\tau$$

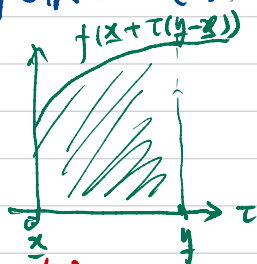$$= f(x) + \nabla f(x)^T(y-x) + \int_0^1 \left[ \nabla f(x+\tau(y-x)) - \nabla f(x)\right]^T (y-x)\, d\tau$$

By rearranging & taking abs. val. on both sides

$$\left| f(y) - f(x) - \nabla f(x)^T (y-x) \right| = \left| \int_0^1 \left[ \nabla f(x+\tau(y-x)) - \nabla f(x)\right]^T (y-x)\, d\tau \right|$$

$$\leq \int_0^1 \left| \left[ \nabla f(x+\tau(y-x)) - \nabla f(x)\right]^T (y-x) \right| d\tau \quad \left( \begin{array}{c} \text{Triangle Ineq} \\ \|a+b\| \leq \|a\| + \|b\| \end{array} \right)$$

$$\leq \int_0^1 \underline{\|\nabla f(x+\tau(y-x)) - \nabla f(x)\|} \cdot \|y-x\|\, d\tau \quad \left( \begin{array}{c} \text{Cauchy-Schwtz} \\ |a^T b| \leq \|a\| \cdot \|b\| \end{array} \right)$$

$$\underset{L - \text{Lipshitz}}{} \leq L\tau \|y-x\|$$

$$\leq \int_0^1 L\tau \|y-x\|^2 d\tau = L\|y-x\|^2 \int_0^1 \tau d\tau = \frac{L}{2}\|y-x\|^2.$$

$$\underbrace{}_{=\frac{1}{2}}$$

(1) is proved.

Step ② : WTS "Descent property of GD":

$$x_{k+1} = x_k - s_k \nabla f(x_k). \quad \text{Plug this into} \quad (1).$$

$$f(x_{k+1}) \leq f(x_k) + \nabla f(x_k)^T \underbrace{(x_{k+1} - x_k)}_{-s_k \nabla f(x_k)} + \frac{L}{2} \underbrace{\|x_{k+1} - x_k\|^2}_{-s_k \nabla f(x_k)}$$

$$= f(x_k) - s\|\nabla f(x_k)\|^2 + \frac{Ls^2}{2}\|\nabla f(x_k)\|^2$$

General:

$$= f(x_k) - s(1 - \frac{Ls}{2})\|\nabla f(x_k)\|^2 \quad \longleftarrow \text{ starting pt.}$$

(2)

$$\geq 0 \, . = 0 \quad \textcolor{red}{\text{Lyapunov.}}$$

Step ③ : Consider $\{\|x_k - x^*\|^2\}_{k=1}^{\infty}$

$$V_k - V_{k-1} \leq -\delta_{k-1}$$

check: $\|x_{k+1} - x^*\|^2 - \|x_k - x^*\|^2$

$$\vdots$$

$$= \|x_k - s\nabla f(x_k) - x^*\|^2 - \|x_k - x^*\|^2$$

$$V_1 - V_0 \leq -\delta_0$$

$$V_k - V_0 \leq -\sum_{i=0}^{k-1} \delta_i$$

$$= \|x_k - x^*\|^2 - 2s\nabla f(x_k)^T (x_k - x^*) + s^2\|\nabla f(x_k)\|^2 - \|x_k - x^*\|^2$$

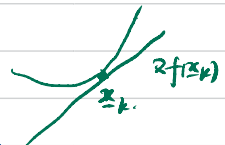$$= -2s\nabla f(x_k)^T (x_k - x^*) + s^2\|\nabla f(x_k)\|^2$$

Due to convexity: $f(x^*) \geq f(x_k) + \nabla f(x_k)^T (x^* - x_k)$

$$\Rightarrow f(x_k) \leq f(x^*) + f(x_k)^T (x_k - x^*).$$

(3)

Plugging (3) into (2):

$$f(z_{k+1}) \leq \underbrace{f(z^k) + \nabla f(z_k)^T (z_k - z^*)}_{f(z_k) \leq} - s\left(1 - \frac{Ls}{2}\right) \|\nabla f(z_k)\|^2. \quad (4)$$

Note that $s \in (0, \frac{1}{L}]$. Then $0 < s \leq \frac{1}{L} \Rightarrow 0 < Ls \leq 1$.

$$\Rightarrow -\frac{1}{2} \leq -\frac{Ls}{2} < 0 \Rightarrow \frac{1}{2} \leq 1 - \frac{Ls}{2} \leq 1 \Rightarrow -s \leq -s\left(1 - \frac{Ls}{2}\right) \leq -\frac{s}{2}.$$

Using this in (4):

$$f(z_{k+1}) - f(z^*) \leq \nabla f(z_k)^T (z_k - z^*) - \frac{s}{2} \|\nabla f(z_k)\|^2$$

$$\Rightarrow -2s \nabla f(z_k)^T (z_k - z^*) \leq -2s (f(z_{k+1}) - f(z^*)) - s^2 \|\nabla f(z_k)\|^2.$$

$$\Rightarrow 2s (f(z_{k+1}) - f(z^*)) \leq 2s \nabla f(z_k)^T (z_k - z^*) - s^2 \|\nabla f(z_k)\|^2$$

$$\Rightarrow f(z_{k+1}) - f(z^*) \leq \frac{1}{2s}\left[\|z_k - z^*\|^2 - \|z_{k+1} - z^*\|^2\right] \quad (5)$$

step ④: Summing (5) from 1 to $k$.

$$\sum_{i=1}^{k} (f(z_i) - f(z^*)) \leq \frac{1}{2s}\left(\|z_0 - z^*\|^2 - \|z_k - z^*\|^2\right).$$

$$\leq \frac{1}{2s} \|z_0 - z^*\|^2$$

Since $\{f(z_k)\}$ is mono. non-incr. (GD descent prop), we have

$$f(z_k) - f(z^*) \leq \frac{1}{k} \sum_{i=1}^{k}[f(z_i) - f(z^*)] \leq \frac{\|z_0 - z^*\|^2}{2sk} = O(\frac{1}{k}). \quad ▣$$

" classical result: $O(\frac{1}{k})$ of GD ".

(HW)
prove: $\|z_k - z^*\| = O(t)$

## Convergence Rate Analysis (Convex): Fixed Step Size

*Proof Sketch.*

- (**Descent Lemma**): $\nabla f$ is $L$-Lipschitz $\Rightarrow$

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2}\|\mathbf{y} - \mathbf{x}\|_2^2, \quad \forall \mathbf{x}, \mathbf{y}$$

- Plugging in $\mathbf{x}_{k+1} = \mathbf{x}_k - s\nabla f(\mathbf{x}_k)$ to obtain:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \left(1 - \frac{Ls}{2}\right) s\|\nabla f(\mathbf{x}_k)\|_2^2$$

- Using the convexity of $f$ and taking $0 < s \leq 1/L$, and , we have

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}^*) + \nabla f(\mathbf{x}_k)^\top (\mathbf{x}_k - \mathbf{x}^*) - \frac{s}{2}\|\nabla f(\mathbf{x}_k)\|_2^2$$

$$= f(\mathbf{x}^*) + \frac{1}{2s}\left(\|\mathbf{x}_k - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|_2^2\right)$$

# Convergence Rate Analysis (Convex): Fixed Step Size

- Summing over iterations & after telescoping:

$$\sum_{i=1}^{k} \big(f(\mathbf{x}_i) - f(\mathbf{x}^*)\big) \leq \frac{1}{2s}\big(\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_k - \mathbf{x}^*\|_2^2\big)$$

$$\leq \frac{1}{2s}\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$$

- Since $f(\mathbf{x}_k)$ is non-increasing, we have

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{1}{k}\sum_{i=1}^{k}\big(f(\mathbf{x}_i) - f(\mathbf{x}^*)\big) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2sk}. \qquad \square$$

# Convergence Rate Analysis (Nonconvex): Fixed Step Size

Assume that $f$ is nonconvex & differentiable, and $L$-smooth

### Theorem 2 (Stationarity Gap)

*If $f$ is nonconvex, differentiable, and $L$-smooth, then gradient descent with fixed step size $s \leq 1/L$ satisfies*

$$\min_{t=0,\ldots,k-1} \|\nabla f(\mathbf{x}_t)\|_2^2 \leq \frac{2(f(\mathbf{x}_0) - f^*)}{sk}$$

*i.e., gradient descent method has sublinear convergence rate $O(1/k)$.*

Remark:

- To get $\|\nabla f(\mathbf{x}_k)\|_2 \leq \epsilon$ for some $k$, it takes $O(\epsilon^{-2})$ iterations.

### Theorem 2 (Stationarity Gap)

If $f$ is nonconvex, differentiable, and $L$-smooth, then gradient descent with fixed step size $s \leq 1/L$ satisfies

$$\min_{t=0,\ldots,k-1} \|\nabla f(\mathbf{x}_t)\|_2^2 \leq \frac{2(f(\mathbf{x}_0) - f^*)}{sk}$$

i.e., gradient descent method has sublinear convergence rate $O(1/k)$.

Proof. We know that:

$$f(x_{k+1}) \leq f(x_k) - s\left(1 - \frac{Ls}{2}\right)\|\nabla f(x_k)\|^2 \quad \Bigg\}$$

$$\text{Also, } 0 < s \leq \frac{1}{L} \Rightarrow -s\left(1 - \frac{Ls}{2}\right) \leq -\frac{s}{2} \quad \Bigg\} \Rightarrow$$

$$f(x_{k+1}) - f(x_k) \leq \frac{-s}{2}\|\nabla f(x_k)\|^2 \qquad (1)$$

Summing (1) from 0 to k-1 :

$$\underbrace{f(x_k) - f(x_0)}_{\geq f(x^*) - f(x_0)} \leq -\frac{s}{2}\sum_{t=0}^{k-1}\|\nabla f(x_t)\|^2 \leq -\frac{sk}{2}\min_{t=0,\ldots,k-1}\|\nabla f(x_k)\|^2$$

Let $f^* = \inf_{x \in \mathbb{R}^n} f(x) > -\infty$. Then :

$$\min_{t=0,\ldots,k-1}\|\nabla f(x_t)\|^2 \leq \frac{2(f(x_0) - f^*)}{sk} = O\left(\frac{1}{k}\right) \quad \square$$

# Strategy 2: Exact Line Search

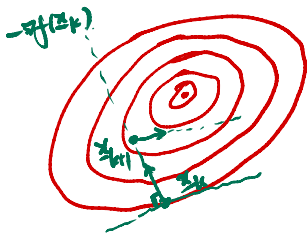Choose the step size $s$ to do the "best" we can along the direction of $-\nabla f(\mathbf{x})$:

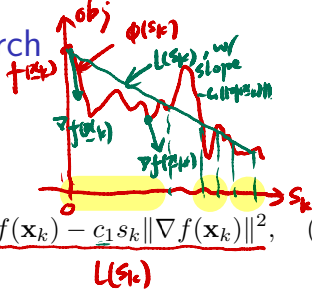$$s = \arg\min_{t \geq 0} f(\mathbf{x}_k - t\nabla f(\mathbf{x}_k))$$

*dir fn*

Limitations:

$$-\nabla f(\mathbf{z}_{k+1})^\top \nabla f(\mathbf{z}_k) = -\nabla f(\mathbf{z}_k - t\nabla f(\mathbf{z}_k))^\top \nabla f(\mathbf{z}_k) = 0$$

- Usually it's too expensive to do this in each iteration.

# Strategy 3: Inexact Line Search

Seek $s_k$ that satisfies Wolfe conditions:

- "Sufficient decrease" in $f$:
$$f(\mathbf{x}_{k+1}) = f(\mathbf{x}_k - s_k \nabla f(\mathbf{x}_k)) \leq f(\mathbf{x}_k) - c_1 s_k \|\nabla f(\mathbf{x}_k)\|^2, \quad (0 < c_1 \ll 1)$$

- "Not zigzagging too badly":
$$\nabla f(\mathbf{x}_{k+1})^\top \nabla f(\mathbf{x}_k) \leq c_2 \|\nabla f(\mathbf{x}_k)\|^2, \quad (c_1 < c_2 < 1)$$

Main features:

- Can show that accumulation points $\bar{\mathbf{x}}$ of $\{\mathbf{x}_k\}$ are stationary: $\nabla f(\bar{\mathbf{x}}) = \mathbf{0}$ (thus minimizer if $f$ is convex)

- Can do 1-dim line search for $s_k$, taking minima of quadratic or cubic interpolations of $f$ and $\nabla f$ at the last two values tried. Use brackets for reliability. Often finds suitable $s_k$ within 3 attempts (see [Nocedal & Wright, 2006, Ch. 3])

# Strategy 3: Inexact Line Search – Backtracking

One way to adaptively choose step size is to use backtracking line search

1. First fix parameters $0 < \beta < 1$ and $0 < \alpha \leq \frac{1}{2}$
2. At each iteration, start with $s = 1$, and while

$$f(\mathbf{x} - s\nabla f(\mathbf{x})) > f(\mathbf{x}) - \alpha s\|\nabla f(\mathbf{x})\|_2^2$$

shink $\underbrace{s = \beta s}_{shrinking}$. Else, perform gradient descent update:
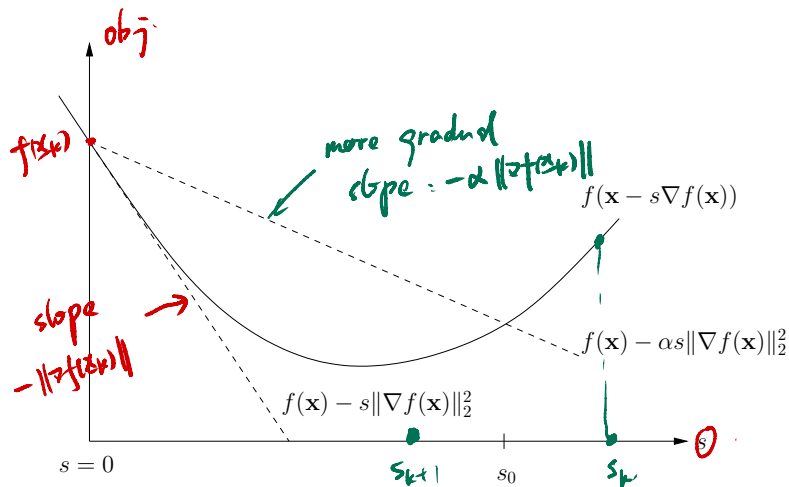
$$\mathbf{x}^+ = \mathbf{x} - s\nabla f(\mathbf{x})$$
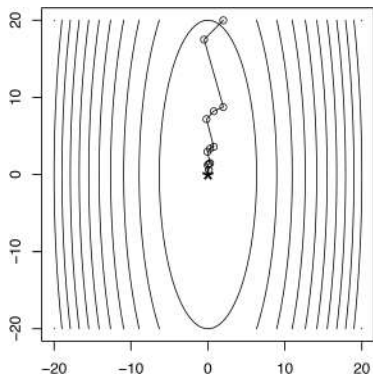
Remarks:

- Simple and tends to work well in practice (further simplification: just take $\alpha = \beta = 1/2$). But doesn't work for $f$ nonsmooth
- Also referred to as Armijo's rule. Step size shrinking very aggressively
- Not checking the second Wolfe condition: the $s_k$ thus identified is "within striking distance" of an $s$ that's not too large

# Backtracking Interpretation

# Backtracking Example

Backtracking picks up roughly the right step size (12 outer iterations, 40 iterations in total):



$$O\left(\frac{1}{k}\right)$$

# Next Class

Stochastic Gradient Descent