

# ECE 8101: Nonconvex Optimization for Machine Learning

Lecture Note 2-3: Gradient Descent

Jia (Kevin) Liu

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

Spring 2022

# 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  $(\mathbf{x}_k, \mathbf{x}_{k-1}, \dots)$ .

**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?

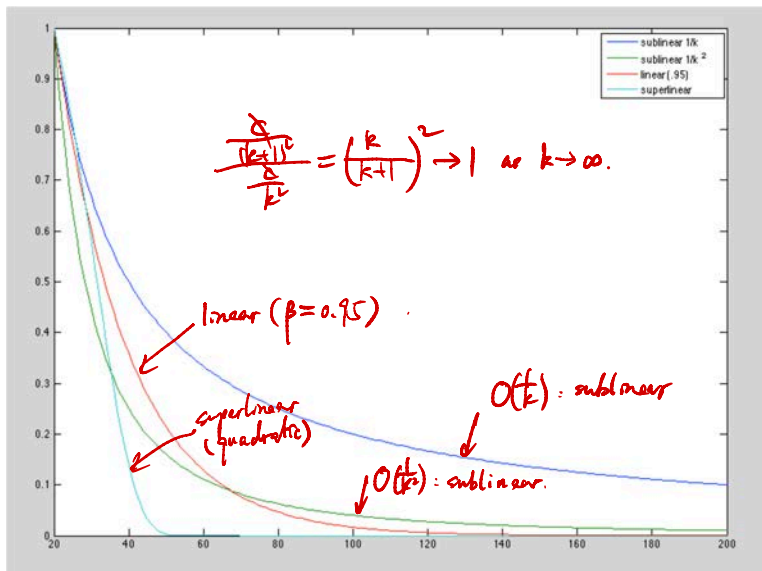
## Definition 1 (Convergence rate)

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

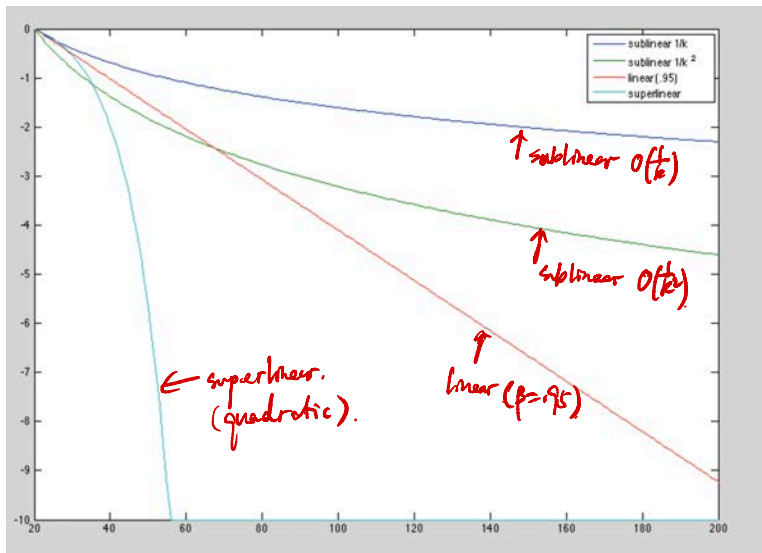
$$\limsup_{k \rightarrow \infty} \frac{\|r_{k+1} - r^*\|}{\|r_k - r^*\|^p} = \beta < \infty.$$

- **Sublinear:**  $p = 1$  and  $\beta = 1$  (e.g.,  $O(1/k)$  rate, kind of slow but still OK)  
*Handwritten:*  $\frac{\|r_k - r^*\|}{\|r_k - r^*\|} \rightarrow 1$ .  
 $\frac{\epsilon}{k} \rightarrow \frac{\epsilon}{k} \rightarrow 0$  as  $k \rightarrow \infty$ .  $r_k - r^* \leq \frac{\epsilon}{k}$ . Desired  $\epsilon > 0$ ,  $\frac{\epsilon}{k} \leq \epsilon \Rightarrow k \geq \frac{\epsilon}{\epsilon} \Rightarrow O(\frac{1}{\epsilon})$ .
- **Linear or geometric:**  $p = 1$  and  $0 < \beta < 1$  (i.e.,  $\|r_{k+1} - r^*\| \leq \beta \|r_k - r^*\|$  for some  $\beta \in (0, 1)$ , or  $\|r_k - r^*\| = O(\beta^k)$ , which is quite fast)  
*Handwritten:*  $\leq \beta^k \|r_k - r^*\| \dots$
- **Superlinear:**  $p > 1$  and  $\beta < \infty$ , or  $p = 1$  and  $\beta = 0$  (i.e.,  $\frac{\|r_{k+1} - r^*\|}{\|r_k - r^*\|} \rightarrow 0$ , that's very fast!)  
*Handwritten:*  $\leq \beta^k \|r_k - r^*\| = O(\beta^k)$ .  
*Handwritten:* also a contraction mapping. Desired  $\epsilon > 0$ ,  $c \cdot \beta^k \leq \epsilon \Rightarrow k \geq \log(\frac{\epsilon}{c}) \Rightarrow O(\log(\frac{1}{\epsilon}))$  iter.
- **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!)  
*Handwritten:* Not only a contraction mapping, but also the rate of contraction is accelerating.  
Desired  $\epsilon > 0$ : Need  $O(\log \log(\frac{1}{\epsilon}))$  iter.  $\leftarrow$  almost const.

# Convergence Rates Comparisons



# 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:

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

Stop if some stopping criterion is satisfied.

$$\left( \begin{array}{l} \text{e.g., } \|\nabla f(\mathbf{z}_k)\| \leq \epsilon, \\ \|\mathbf{z}_{k+1} - \mathbf{z}_k\| \leq \epsilon, \\ \text{some \# of steps, } \dots \end{array} \right).$$

# Gradient Descent: Geometric Interpretation

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

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

SO - approx

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

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

$$f(\mathbf{y}) \approx \underbrace{f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})}_{\text{1st-order approx.}} + \underbrace{\left(\frac{1}{2s}\right) \|\mathbf{y} - \mathbf{x}\|^2}_{\text{proximity term}}$$

"proximity term":  
penalize moving too far away from  $\mathbf{x}$ .

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:

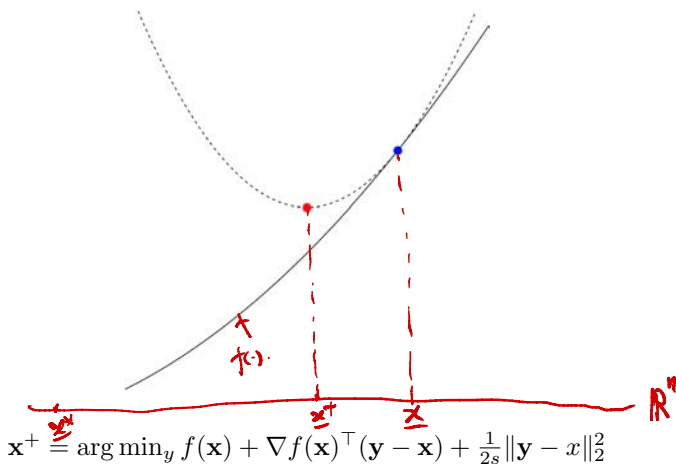
$$\mathbf{x}^+ = \mathbf{x} - s \nabla f(\mathbf{x})$$

Quad fn of  $\mathbf{y}$ : unconstr, set grad to 0, then solve for  $\mathbf{y}$ .

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



# Gradient Descent: Geometric Interpretation



## Questions:

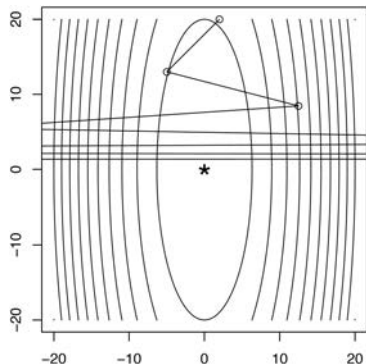
- How to choose step sizes  $\{s_k\}$ ?
- What is the according convergence rate? Or does it depend on  $\{s_k\}$ ?

## Strategy 1: Fixed Step Size

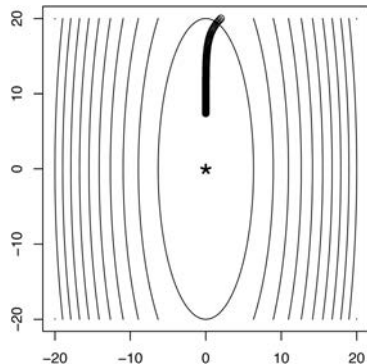
Simply set  $s_k = s$  for all  $k = 1, 2, 3, \dots$

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

**Example:** Consider  $f(\mathbf{x}) = (10x_1^2 + x_2^2)/2 \Rightarrow (\mathbf{x}_1^*, \mathbf{x}_2^*) = (0, 0)$ .



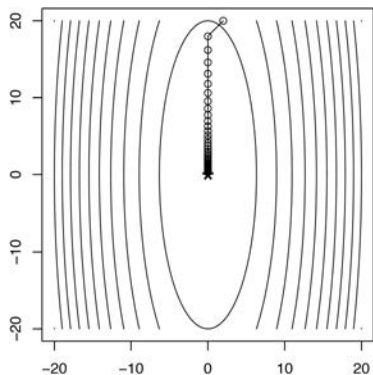
8 iterations



100 iterations

# Strategy 1: Fixed Step Size

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



Dynamic syst. View Pt:

$$\dot{\underline{z}} = -\nabla f(\underline{z}(t))$$

Lyapunov fn:

$$V(\underline{z}) = f(\underline{z}) - f(\underline{z}^*)$$

$$\frac{dV(\underline{z}(t))}{dt} = \nabla f(\underline{z})^T \frac{d\underline{z}(t)}{dt}$$

$$= -\nabla f(\underline{z}(t))^T \nabla f(\underline{z}(t))$$

$$= -\|\nabla f(\underline{z}(t))\|_2^2 \leq 0$$

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

Need info of the "Lipschitz const." of  $\nabla f(\underline{z})$  (smoothness).

# Convergence Rate Analysis (Convex): Fixed Step Size

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

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 \leq \underline{L\|\mathbf{y} - \mathbf{x}\|_2}, \quad \forall \mathbf{x}, \mathbf{y} \quad (L\text{-smoothness}).$$

That is,  $\nabla f$  is Lipschitz continuous with constant  $L > 0$  ( $L$ -Lipschitz continuous)  
 $h: \mathcal{D} \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ .  $h$  is call lip. cont.:  $\exists L > 0$ , s.t.  $\|h(\mathbf{x}) - h(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$ ,  $\forall \mathbf{x}, \mathbf{y} \in \mathcal{D}$ .

## Theorem 1 (Optimality Gap)

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

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2sk}, \quad = O\left(\frac{1}{k}\right).$$

*init. pt. dist.* ↙

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

Remark:

- To get  $f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \epsilon$ , 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

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2sk}, = O\left(\frac{1}{k}\right).$$

init pt. dist.

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

Proof: step 0: claim: If  $\nabla f$  is Lipschitz, then:

$$f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|_2^2, \quad \forall x, y \in \mathbb{R}^n.$$

To show (1), we start from: dir. der.

$$f(y) = f(x) + \int_0^1 \nabla f(x + \tau(y-x)) d\tau.$$

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

add & subtract

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

By some rearranging and take abs. value on both sides:

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

$$\leq \int_0^1 |[\nabla f(x + \tau(y-x)) - \nabla f(x)]^T (y-x)| d\tau$$

(Triangle ineq:  $\|a+b\| \leq \|a\| + \|b\|$ )

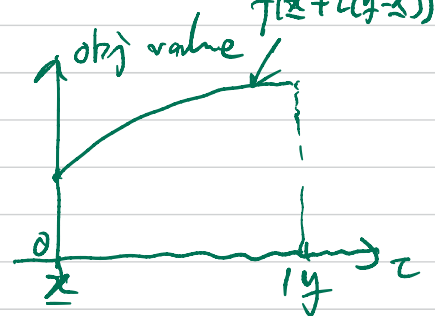
$$\leq \int_0^1 \underbrace{\|\nabla f(x + \tau(y-x)) - \nabla f(x)\|}_{L\text{-Lipschitz} \leq L\tau\|y-x\|} \cdot \|y-x\| \cdot d\tau.$$

(Cauchy-Schwarz Ineq:  $|a^T b| \leq \|a\| \cdot \|b\|$ )

$$\leq \int_0^1 L \cdot \tau \|y-x\|^2 \cdot d\tau = L \|y-x\|^2 \underbrace{\int_0^1 \tau d\tau}_{\frac{1}{2}} = \frac{L}{2} \|y-x\|^2. \quad (1) \text{ is proved.}$$

(descent lemma)

(1)



Step ②: WTS: "Descent property of GD":

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

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x}_{k+1} - \mathbf{x}_k) + \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2.$$

$$= f(\mathbf{x}_k) - s \|\nabla f(\mathbf{x}_k)\|_2^2 + \frac{Ls^2}{2} \|\nabla f(\mathbf{x}_k)\|_2^2$$

$$= f(\mathbf{x}_k) - s \left(1 - \frac{Ls}{2}\right) \|\nabla f(\mathbf{x}_k)\|_2^2. \quad (2)$$

Step ③: From convexity of  $f(\mathbf{x})$ , we have:

$$f(\mathbf{x}^*) \geq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x}^* - \mathbf{x}_k).$$

$$\Rightarrow f(\mathbf{x}_k) \leq f(\mathbf{x}^*) + \nabla f(\mathbf{x}_k)^T (\mathbf{x}_k - \mathbf{x}^*). \quad (3)$$

Plugging (3) into (2) yields:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}^*) + \nabla f(\mathbf{x}_k)^T (\mathbf{x}_k - \mathbf{x}^*) - s \left(1 - \frac{Ls}{2}\right) \|\nabla f(\mathbf{x}_k)\|_2^2 \quad (4)$$

Recall step-size  $s \in (0, \frac{1}{L}]$ . Then,

$$0 < s \leq \frac{1}{L} \Rightarrow 0 < Ls \leq 1 \Rightarrow -\frac{1}{2} \leq -\frac{Ls}{2} < 0 \Rightarrow \frac{1}{2} \leq 1 - \frac{Ls}{2} < 1.$$

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

Use the above in (4)  $\Rightarrow$

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}_k)^T (\mathbf{x}_k - \mathbf{x}^*) - \frac{s}{2} \|\nabla f(\mathbf{x}_k)\|_2^2 \quad (5)$$

step ④: Consider RHS of (5):

$$\begin{aligned} & \nabla f(x_k)^T (x_k - x^*) - \frac{s}{2} \|\nabla f(x_k)\|_2^2 \\ &= -\frac{1}{2s} \left[ s^2 \|\nabla f(x_k)\|_2^2 - 2s \nabla f(x_k)^T (x_k - x^*) + \|x_k - x^*\|_2^2 - \|x_k - x^*\|_2^2 \right] \\ &= -\frac{1}{2s} \left[ \|(x_k - x^*) - s \nabla f(x_k)\|_2^2 - \|x_k - x^*\|_2^2 \right] \\ &= -\frac{1}{2s} \left[ \|x_{k+1} - x^*\|_2^2 - \|x_k - x^*\|_2^2 \right] \end{aligned}$$

Therefore, (5)  $\Rightarrow$

$$f(x_{k+1}) - f(x^*) \leq \frac{1}{2s} \left( \|x_k - x^*\|_2^2 - \|x_{k+1} - x^*\|_2^2 \right). \quad (6)$$

step ⑤: Summing (6) from 0 to  $k-1$  (telescoping).

$$\begin{aligned} \sum_{i=0}^{k-1} (f(x_i) - f(x^*)) &\leq \frac{1}{2s} \left( \|x_0 - x^*\|_2^2 - \|x_k - x^*\|_2^2 \right) \\ &\leq \frac{1}{2s} \|x_0 - x^*\|_2^2 \end{aligned}$$

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

$$f(x_k) - f(x^*) \leq \frac{1}{k} \sum_{i=0}^{k-1} [f(x_i) - f(x^*)] \leq \frac{\|x_0 - x^*\|_2^2}{2sk} = O\left(\frac{1}{k}\right). \quad \square$$

Classic  $O\left(\frac{1}{k}\right)$  of GD.

# 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

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



# Convergence Rate Analysis (Convex): Fixed Step Size

- Summing over iterations & after telescoping:

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

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

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

□

# 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

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

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

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

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

Proof.. We know that,

$$\left. \begin{aligned} f(\mathbf{x}_{k+1}) &\leq f(\mathbf{x}_k) - s\left(1 - \frac{Ls}{2}\right) \|\nabla f(\mathbf{x}_k)\|_2^2 \\ \text{Also, } 0 < s \leq \frac{1}{L} &\Rightarrow -s\left(1 - \frac{Ls}{2}\right) \leq -\frac{s}{2} \end{aligned} \right\} \Rightarrow$$

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{s}{2} \|\nabla f(\mathbf{x}_k)\|_2^2. \text{ Summing from } 0 \text{ to } k-1.:$$

$$\begin{aligned} \underline{f(\mathbf{x}_k) - f(\mathbf{x}_0)} &\leq -\frac{s}{2} \sum_{t=0}^{k-1} \|\nabla f(\mathbf{x}_t)\|_2^2 \leq -\frac{sk}{2} \min_{t=0, \dots, k-1} \|\nabla f(\mathbf{x}_t)\|_2^2 \\ &\geq f(\mathbf{x}^*) - f(\mathbf{x}_0) \end{aligned}$$

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

$$\min_{t=0, \dots, k-1} \|\nabla f(\mathbf{x}_t)\|_2^2 \leq \frac{2(f(\mathbf{x}_0) - f^*)}{sk} = O\left(\frac{1}{k}\right). \quad \square$$

## Strategy 2: Exact Line Search

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

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

Limitations:

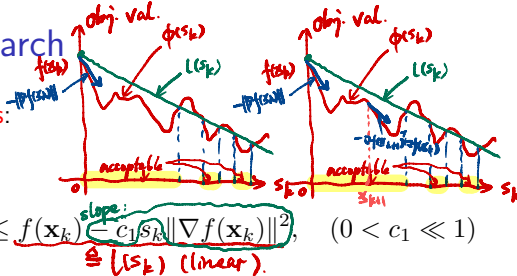
Take dir. der. w.r.t.  $t$ , set it to 0,  
solve for  $t$ .  $-\nabla f(\mathbf{x}_{k+1})^T \nabla f(\mathbf{x}_k) = 0$

- 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  $f$ :

$$\cong \phi(s_k)$$

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

$\cong L(s_k) \text{ (linear)}$

- “Not zigzagging too badly”:

$$-\nabla f(\mathbf{x}_{k+1})^\top \nabla f(\mathbf{x}_k) \geq -c_2 \|\nabla f(\mathbf{x}_k)\|^2, \quad (c_1 < c_2 < 1)$$

Main features:



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

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

shrink  $s = \beta s$ . Else, perform gradient descent update:

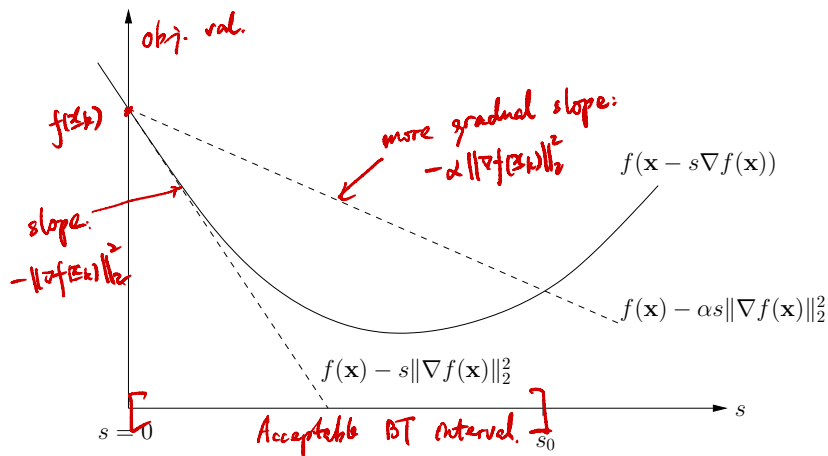
*shrink by  
a  $\beta$ -factor.*

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

Remarks:

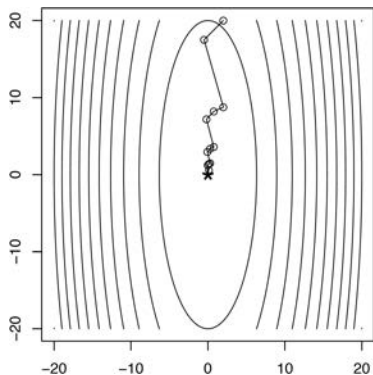
- Simple and tends to work well in practice (further simplification: just take  $\alpha = \beta = 1/2$ ). But doesn't work for  $f$  nonsmooth
- Also referred to as **Armijo's rule**. Step size shrinking very aggressively
- Not checking the second Wolfe condition: the  $s_k$  thus identified is “within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





Next Class

## Stochastic Gradient Descent