

# COM S 578X: Optimization for Machine Learning

## Lecture Note 6: Gradient Descent

Jia (Kevin) Liu

Assistant Professor  
Department of Computer Science  
Iowa State University, Ames, Iowa, USA

Fall 2019

# Outline

In this lecture:

- Convergence rate concept
- Gradient descent method
- Step size selection strategies
- Convergence performance of gradient descent

# First-Order Algorithms: Smooth Convex Functions

Consider an unconstrained optimization problem, with  $f$  smooth and convex:

$$\underline{A} \preceq \underline{B} \Leftrightarrow (\underline{B} - \underline{A}) \text{ is PSD.}$$

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

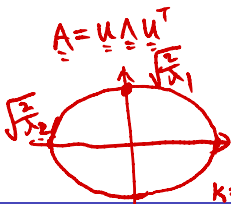
twice cont. diff.  $\mathcal{F}_L^{2,1}$   $\leftarrow$  1st der.  $\swarrow$   $\searrow$  UBed by  $L$ .

- Usually assume  $\mu \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L \mathbf{I}, \forall \mathbf{x}$ , with  $0 \leq \mu \leq L$  (using Nesterov's notation:  $\mathcal{F}_L^{2,1}, \mathcal{S}_{\mu,L}^{2,1}$ )
- If  $\mu > 0$ , then  $f$  is  $\mu$ -strongly convex, i.e.,

twice cont. diff.  $\mathcal{S}_{\mu,L}^{2,1}$   $\leftarrow$  1st der  $\swarrow$   $\searrow$  UBed by  $L$  strongly convex  $\uparrow$  LBed by  $\mu$ .

$$f(\mathbf{y}) \geq \underbrace{f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})}_{\text{FO approx.}} + \underbrace{\frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2}_{\geq 0}$$

- Condition number:  $\kappa = L/\mu$  (the larger  $\kappa$  is, the more ill-conditioned)
- In ML, people are often interested in convex quadratics, e.g.,



$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \underline{A} \mathbf{x}, \quad \underline{A} = \frac{1}{\|\tilde{\mathbf{x}}\|^2} (\tilde{\mathbf{x}} \tilde{\mathbf{x}})^\top \underline{A} (\tilde{\mathbf{x}} \tilde{\mathbf{x}}) \quad \tilde{\mathbf{x}} \in \mathbb{R}^2 \quad \begin{bmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{bmatrix} \begin{bmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{bmatrix}$$

$$\mu \mathbf{I} \preceq \underline{A} \preceq L \mathbf{I}$$

$$f(\mathbf{x}) = \frac{1}{2} \|\underline{A} \mathbf{x} - \mathbf{b}\|_2^2, \quad \mu \mathbf{I} \preceq \underline{A}^\top \underline{A} \preceq L \mathbf{I} = \frac{\tilde{x}_1^2}{\tilde{x}_1^2} + \frac{\tilde{x}_2^2}{\tilde{x}_2^2} = 1$$

an ellipse

# Iterative Algorithms

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:

- 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?

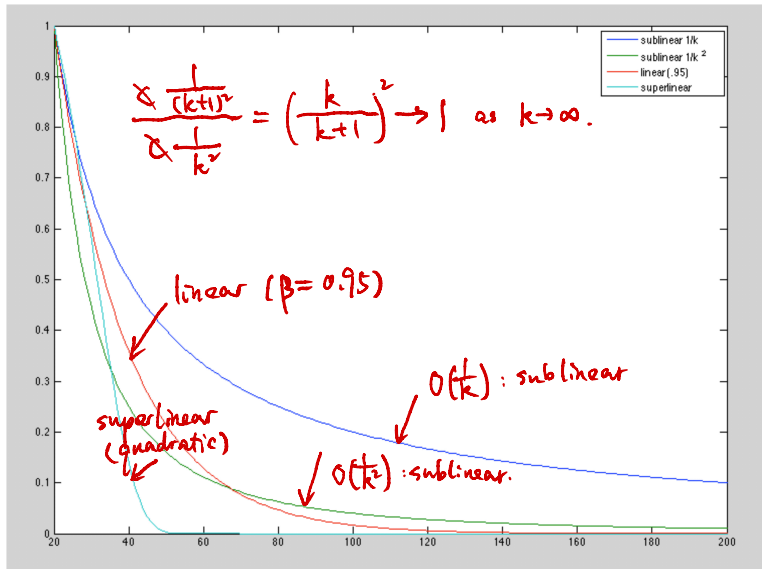
## Definition 1 (Convergence rate)

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

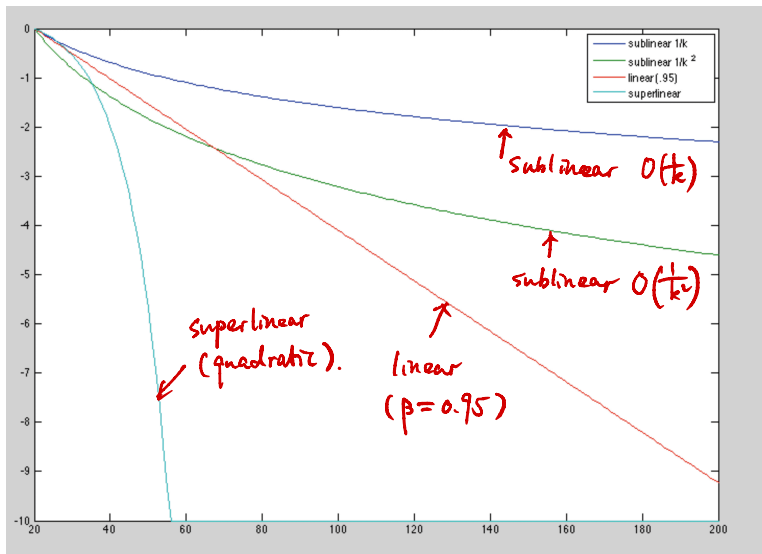
$$\limsup_{k \rightarrow \infty} \frac{\|r_{k+1} - r^*\|}{\|r_k - r^*\|^p} = \beta < \infty.$$

- Sublinear:**  $p = 1$  and  $\beta = 1$  (e.g.,  $O(1/k)$  rate) *i.e.,  $\frac{\|r_{k+1} - r^*\|}{\|r_k - r^*\|} \rightarrow 1$ , as  $k \rightarrow \infty$*   
*iteration complexity:  $\frac{c}{k}$*   
*Desired  $\epsilon$ :  $\frac{c}{k} \leq \epsilon \Rightarrow k \geq \frac{c}{\epsilon} = \Omega(\frac{1}{\epsilon})$*
- Linear or geometric:**  $p = 1$  and  $0 < \beta < 1$  (i.e.,  $\|r_{k+1} - r^*\| \leq \beta \|r_k - r^*\|$ ) *contraction mapping*  
 for some  $\beta \in (0, 1)$ , or  $\|r_k - r^*\| = O(\beta^k)$ , which is quite fast)  *$\leq \beta^k \|r_0 - r^*\| \dots \leq \beta^k \|r_0 - r^*\| = O(\beta^k)$*   
*Desired  $\epsilon$ :  $c\beta^k \leq \epsilon \Rightarrow k \geq c \log(\frac{1}{\epsilon})$ : Need  $O(\log \frac{1}{\epsilon})$  iter.*
- 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!) *Not only a contraction mapping but also, the rate of contraction is accelerating!*
- Quadratic:**  $p = 2$  and  $\beta < \infty$  (i.e.,  $\|r_{k+1} - r^*\| \leq \beta \|r_k - r^*\|^2$ , # of correct significant digits doubles each iteration. We rarely need anything faster than this!) *For  $\epsilon$ -accuracy: Need  $O(\log \log(\frac{1}{\epsilon}))$  iterations  $\leftarrow$  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.

(e.g.,  $\|\nabla f(\mathbf{x}_k)\| \leq \epsilon$ ,  
some # of iter, ---)



# Gradient Descent: Geometric Interpretation

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

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

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

$$f(\mathbf{y}) \approx \underbrace{f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})}_{F(\mathbf{y})} + \underbrace{\frac{1}{2s} \|\mathbf{y} - \mathbf{x}\|^2}_{\text{"proximity term": penalize moving too far from } \mathbf{x}}$$

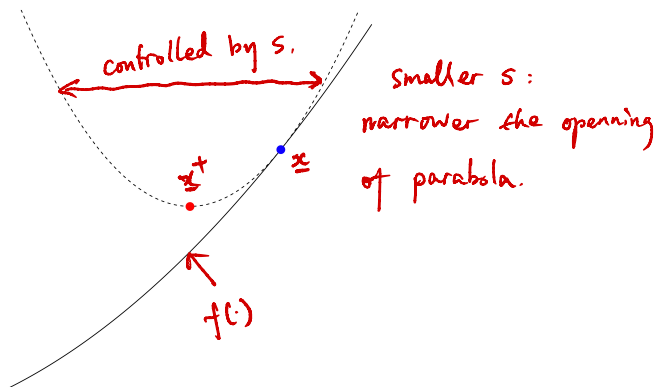
Can be viewed as a linear approximation to  $f$ , with proximity term to  $\mathbf{x}$  weighted by  $\frac{1}{2s}$ . Choose next point  $\mathbf{y} = \mathbf{x}^+$  to minimize this approximation:

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

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

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

# Gradient Descent: Geometric Interpretation



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

## Questions:

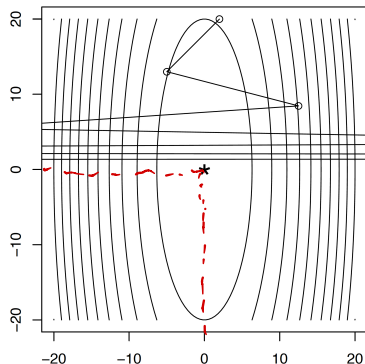
- How to choose step sizes  $\{s_k\}$ ?
- What is the according convergence rate? Or does it depend on  $\{s_k\}$ ? *Yes!*

## Strategy 1: Fixed Step Size

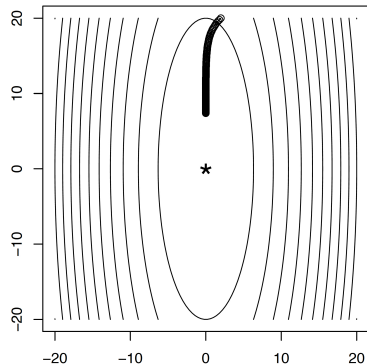
Simply set  $s_k = s$  for all  $k = 1, 2, 3, \dots$

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

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



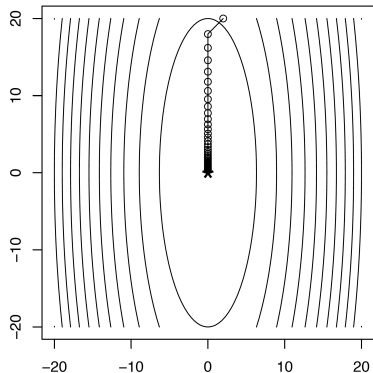
8 iterations



100 iterations

## Strategy 1: Fixed Step Size

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



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

*Need info of the “Lipschitz const.” of  $\nabla f(x)$ .*

## Strategy 2: Exact Line Search

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

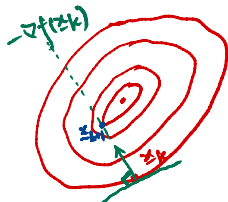
$$s = \arg \min_{t \geq 0} f(\mathbf{x} - t \nabla f(\mathbf{x}))$$

Limitations:

Take directional der., set it to 0, and solve for  $t$ :

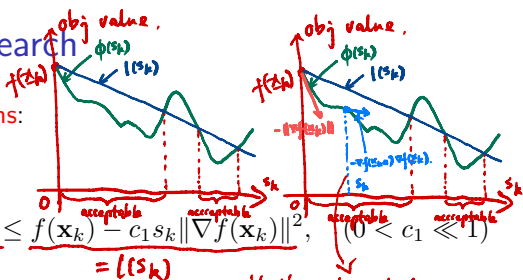
$$-\nabla f(\mathbf{z}_{k+1})^T \nabla f(\mathbf{z}_k) = -\nabla f(\mathbf{z}_k - t \nabla f(\mathbf{z}_k))^T \nabla f(\mathbf{z}_k) = 0$$

- Usually it's too expensive to do this in each iteration.
- **Spoiler:** Our convergence rate analysis later will also show that it's not worth the effort



# Strategy 3: Inexact Line Search

Seek  $s_k$  that satisfies **Wolfe conditions**:



- “Sufficient decrease” in  $f$ :

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

- “Not zigzagging too badly”:

$$\underbrace{-\nabla f(\mathbf{x}_{k+1})^\top \nabla f(\mathbf{x}_k)}_{\text{directional der. of } \phi(s_k)} \geq -c_2 \|\nabla f(\mathbf{x}_k)\|^2, \quad (c_1 < c_2 < 1) \text{ should STOP!}$$

Main features: w.r.t.  $s_k$  (i.e.,  $\phi'(s_k)$ ):  $\phi'(s_k) = -\nabla f(\mathbf{x}_{k+1})^\top \nabla f(\mathbf{x}_k)$

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

*or start w/  $s_{int}$*

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

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

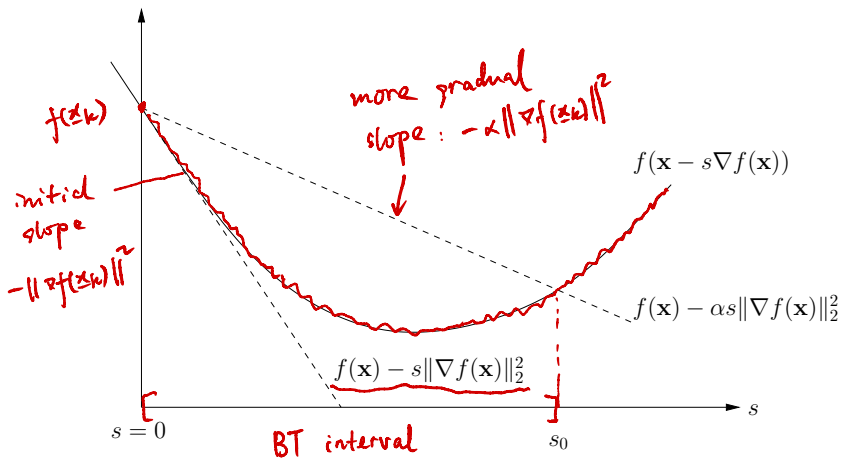
*shrink by  
a factor  $\beta$ .*

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

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

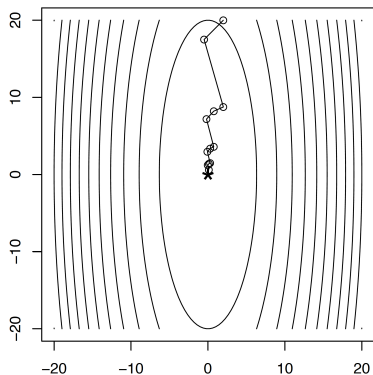
# Backtracking Interpretation





# Backtracking Example

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



# Convergence Rate Analysis: Fixed Step Size

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

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

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

(Nesterov notation:  $f \in \mathcal{F}_L^{1,1}$ )

## Theorem 1

Gradient descent with fixed step size  $s \leq 1/L$  satisfies

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

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.

# Convergence Rate Analysis: Fixed Step Size

*Proof Sketch.*

- $\nabla f$  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: 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

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

□

# GD Convergence: Fixed Step Size under Strong Convexity

Assume that  $f$  is convex is differentiable,  $\nabla f$   $L$ -Lipschitz, and  $\mu$ -strongly convex, i.e.,  $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2}\mu\|\mathbf{y} - \mathbf{x}\|_2^2$ . (Nesterov notation:  $\mathcal{S}_{\mu,L}^{1,1}$ )

## Theorem 2

Gradient descent with fixed step size  $0 < s \leq 2/(L + \mu)$  satisfies

$$\|\mathbf{x}_k - \mathbf{x}^*\| \leq \left( \sqrt{1 - \frac{2s\mu L}{\mu + L}} \right)^k \|\mathbf{x}_0 - \mathbf{x}^*\|$$

i.e., GD has *linear* convergence rate. If  $s = \frac{2}{\mu + L}$ , then

$$\|\mathbf{x}_k - \mathbf{x}^*\| \leq \left( \frac{\kappa - 1}{\kappa + 1} \right)^k \|\mathbf{x}_0 - \mathbf{x}^*\|$$

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{L}{2} \left[ \left( \frac{\kappa - 1}{\kappa + 1} \right)^2 \right]^k \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$$

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}^*)^\top (\mathbf{x}_k - \mathbf{x}^*) + \frac{L}{2} \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \leq \frac{L}{2} \left( \frac{\kappa - 1}{\kappa + 1} \right)^{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$$

$$\begin{aligned} & \sqrt{1 - \frac{2}{\mu+L} \times \frac{2}{\mu+L}} \\ &= \sqrt{1 - \frac{4}{(\mu+L)^2}} = \sqrt{\frac{(L-\mu)^2}{(L+\mu)^2}} \\ &= \frac{L-\mu}{L+\mu} = \frac{\kappa-1}{\kappa+1} \end{aligned}$$

# GD Convergence: Fixed Step Size under Strong Convexity

*Proof Sketch.* For notational convenience, let  $r_k$  denote the residual  $\|\mathbf{x}_k - \mathbf{x}^*\|$ .

- Consider  $r_{k+1}^2$ , we have

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

- According to [Nesterov, Thm 2.1.12], if  $f \in \mathcal{S}_{\mu,L}^{1,1}$ , we have

$$\nabla f(\mathbf{x}_k)^\top (\mathbf{x}_k - \mathbf{x}^*) \geq \frac{\mu L}{\mu + L} r_k^2 + \frac{1}{\mu + L} \|\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}^*)\|_2^2$$

- Plugging in (1) and using the fact that  $\nabla f(\mathbf{x}^*) = 0$ , we have:

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

- The last inequality in Thm 2 follows from the  $L$ -Lipschitz gradient assumption. □

# Convergence Rate Analysis: Backtracking

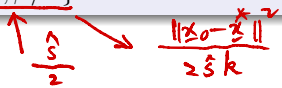
Same assumption:  $f \in \mathcal{F}_L^{1,1}$ .

## Theorem 3

Gradient descent with backtracking line search satisfies:

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{4\alpha \min\{1, \beta/L\}k} = O(1/k)$$

- Same **sublinear** rate as fixed step size
- If  $\beta$  is not too small, then we don't lose much compared to fixed step size ( $\beta/L$  vs  $1/L$ )



# Convergence Rate Analysis: Backtracking

Proof.

- Recall BT exit condition:  $f(\mathbf{x}_k - s\nabla f(\mathbf{x}_k)) \leq f(\mathbf{x}_k) - \alpha s \|\nabla f(\mathbf{x}_k)\|_2^2$ . This is satisfied when  $s \leq 1/L$ , because  $s \leq 1/L \Rightarrow -s + \frac{Ls^2}{2} \leq -\frac{s}{2}$

- Using this and the  $L$ -Lipschitz assumption, we have

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

In const. step:

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

- Hence, BTLS terminates either with  $s = 1$  or with  $s \geq \beta/L$ . Thus, we have

$$f(\mathbf{x}_{k+1}) \leq \underline{f(\mathbf{x}_k)} - \min\{\alpha, \beta\alpha/L\} \|\nabla f(\mathbf{x}_k)\|_2^2$$

$$\leq f(\mathbf{x}^*) + \nabla f(\mathbf{x}_k)^T (\mathbf{x}_k - \mathbf{x}^*) - \min\{\alpha, \frac{\beta\alpha}{L}\} \|\nabla f(\mathbf{x}_k)\|_2^2 \quad (\Delta)$$

- The rest of the proof follows essentially the same line of arguments as in the proof for the fixed step size case. □

Call:  $\frac{s}{2}$

(starting from  $(\Delta)$ , the rest just follows from fixed step-size case.)



# GD Convergence: Backtracking under Strong Convexity

Assume that  $f \in \mathcal{S}_{\mu,L}^{1,1}$

## Theorem 4

*Gradient descent with backtracking line search satisfies*

$$\|f(\mathbf{x}_k) - f(\mathbf{x}^*)\| \leq (1 - \min\{2\mu\alpha, 2\beta\alpha\mu/L\})^k \|f(\mathbf{x}_0) - f(\mathbf{x}^*)\|$$

*i.e., GD has a **linear** convergence rate.*

# GD Convergence: Backtracking under Strong Convexity

Proof.

- From the proof of the weakly convex case, we have obtained:

$$f(\mathbf{x}_{k+1}) \leq \underbrace{f(\mathbf{x}_k)}_{-f(\mathbf{x}^*)} - \min\{\alpha, \beta\alpha/L\} \underbrace{\|\nabla f(\mathbf{x}_k)\|_2^2}_{-f(\mathbf{x}^*)}$$

- Noting  $\|\nabla f(\mathbf{x}_k)\|_2^2 \geq 2\mu(f(\mathbf{x}_k) - f(\mathbf{x}^*))$  & subtracting  $f(\mathbf{x}^*)$  on both sides:

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq (1 - \min\{2\mu\alpha, 2\beta\alpha\mu/L\}) (f(\mathbf{x}_k) - f(\mathbf{x}^*)).$$

Proof. If  $f \in \mathcal{S}_{\mu/L}^{bl}$ , recall we've shown:

$$f(y) \geq \underbrace{f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2}_{\text{RHS is convex \& quad fn. of } y \text{ (given fixed } x\text{)}} \quad (\text{def. of strong convexity}).$$

Minimize RHS w.r.t.  $y$  (find least LB)  $\Rightarrow \tilde{y} = x - \frac{1}{\mu} \nabla f(x)$ . (5)

Plug (5) into RHS yields:  $f(y) \geq f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2$ . Multiply both sides by  $2\mu$  rearranging.  $\square$

# Convergence Rate Analysis: Exact LS

Assume that  $f \in \mathcal{F}_L^{1,1}$

## Theorem 5

*Gradient descent with exact line search satisfies*

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{L \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2k} = O(1/k)$$

*i.e., GD has a **sublinear** convergence rate.*

# Convergence Rate Analysis: Exact LS

Proof.

- From  $L$ -Lipschitz and  $\mathbf{x}_{k+1} = \mathbf{x}_k - s \nabla f(\mathbf{x}_k)$ , we have

optimal step-size is also used by the min of the RMS

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

- Minimize over  $s$  on both sides yields:

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

$\Rightarrow s_e = \frac{1}{L}$

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

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

- The rest of the proof follows exactly the same arguments as in the proof of the fixed step size case □

(DIY)

Σ

# GD Convergence: Exact LS under Strong Convexity

Assume that  $f \in \mathcal{S}_{\mu,L}^{2,1}$

## Theorem 6

*Gradient descent with exact line search satisfies*

$$\|f(\mathbf{x}_k) - f(\mathbf{x}^*)\| \leq (1 - \mu/L)^k \|f(\mathbf{x}_0) - f(\mathbf{x}^*)\|$$

*i.e., GD has a **linear** convergence rate.*

## Observation

No improvement in the linear rate over fixed step size!

# GD Convergence: Exact LS under Strong Convexity

*Proof.*

- From  $L$ -Lipschitz and  $\mathbf{x}_{k+1} = \mathbf{x}_k - s\nabla f(\mathbf{x}_k)$ , we have

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

- Minimize over  $s$  on both sides yields:

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

- Following the same step in BTLS, subtracting  $f(\mathbf{x}^*)$  from both sides and noting  $\|\nabla f(\mathbf{x}_k)\|_2^2 \geq 2\mu(f(\mathbf{x}_k) - f(\mathbf{x}^*))$ , we have

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \underbrace{\left(1 - \frac{\mu}{L}\right)}_{\frac{1}{\kappa}} (f(\mathbf{x}_k) - f(\mathbf{x}^*)).$$

□

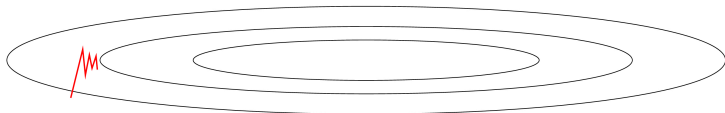
# In Summary

	Convex	Strongly Convex
Fixed Step Size	$O(1/k)$	$O\left[\left(\frac{\kappa-1}{\kappa+1}\right)^{2k}\right]$
BTLS	$O(1/k)$	$O((1 - \min\{2\mu\alpha, 2\beta\alpha/\kappa\})^k)$
Exact LS	$O(1/k)$	$O\left[\left(\frac{\kappa-1}{\kappa}\right)^k\right]$

only slightly better.

# The Slow Linear Rate Is Typical

Not just pessimistic bound – It really is quite slow!





## Next Class

# Accelerated First-Order Methods

①

Thm 1: If  $f \in \mathcal{F}_L^{1,1}(\mathbb{R}^n)$ , then GD with const. step-size  $s \in (0, \frac{1}{L})$  satisfies:

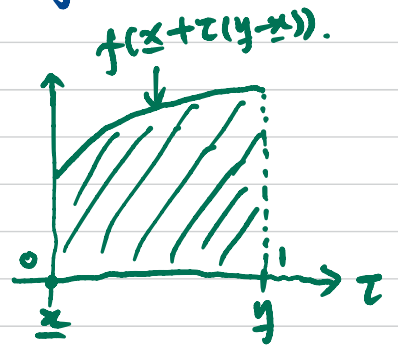
$$f(x_k) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2sk}, \quad \forall k \quad \left( \begin{array}{l} \mathcal{O}(\frac{1}{k}) \text{ sublinear} \\ \text{convergence rate} \end{array} \right)$$

Proof. step ① Claim: If  $\nabla f$  is Lipschitz, then:

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

To show  $(\star)$ , we start from the following:

$$f(y) = f(x) + \int_0^1 \underbrace{f'(x + \tau(y-x))}_{\text{dir. der.}} d\tau$$



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

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 taking absolute value on both sides:

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

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

$$\leq \int_0^1 \underbrace{\|\nabla f(x + \tau(y-x)) - \nabla f(x)\|}_{L\tau\|y-x\|} \cdot \|y-x\| d\tau \quad \left( \begin{array}{l} \text{Cauchy-Schwartz Ineq.} \\ |a^T b| \leq \|a\| \cdot \|b\| \end{array} \right)$$

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

②

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

$\mathbf{z}_{k+1} = \mathbf{z}_k - s_k \nabla f(\mathbf{z}_k)$ . Plugg this in (\*):

$$\begin{aligned} f(\mathbf{z}_{k+1}) &\leq f(\mathbf{z}_k) + \nabla f(\mathbf{z}_k)^T \underbrace{(\mathbf{z}_{k+1} - \mathbf{z}_k)}^{-s \nabla f(\mathbf{z}_k)} + \frac{L}{2} \underbrace{\|\mathbf{z}_{k+1} - \mathbf{z}_k\|_2^2}_{-s \nabla f(\mathbf{z}_k)} \\ &= f(\mathbf{z}_k) - s \|\nabla f(\mathbf{z}_k)\|_2^2 + \frac{Ls^2}{2} \|\nabla f(\mathbf{z}_k)\|_2^2 \\ &= f(\mathbf{z}_k) - s \left(1 - \frac{Ls}{2}\right) \|\nabla f(\mathbf{z}_k)\|_2^2 \quad (\Delta) \end{aligned}$$

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

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

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

Plugging (\*) into (Δ) yields:

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

Now, let's take step-size  $s \in (0, \frac{1}{L}]$ , then

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

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

Using above in (\*\*\*)  $\Rightarrow$

$$f(\mathbf{z}_{k+1}) - f(\mathbf{z}^*) \leq \underbrace{\nabla f(\mathbf{z}_k)^T (\mathbf{z}_k - \mathbf{z}^*) - \frac{s}{2} \|\nabla f(\mathbf{z}_k)\|_2^2}_{\text{green underline}} \quad (\Delta\Delta).$$

③

Step ④: Consider RHS:  $\nabla f(x_k)^T (x_k - x^*) - \frac{s}{2} \|\nabla f(x_k)\|_2^2$

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

Therefore,  $(\Delta\Delta) \Rightarrow$

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

Step ⑤: Summing (12) from 1 to k (telescoping):

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

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

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

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

Thm 2: If  $f \in S_{\mu, L}^{1,1}$  and  $s \leq \frac{2}{L+\mu}$ , then

$$\|z_k - z^*\| \leq \left( \sqrt{1 - \frac{2\mu L}{L+\mu}} \right)^k \|z_0 - z^*\|.$$

Proof: Consider:  $\|z_{k+1} - z^*\|^2$ :

$$\begin{aligned} \|z_{k+1} - z^*\|^2 &= \|z_k - z^* - s \nabla f(z_k)\|^2 \\ &= \|z_k - z^*\|^2 + s^2 \|\nabla f(z_k)\|^2 - 2s \nabla f(z_k)^T (z_k - z^*). \end{aligned} \quad (\Delta)$$

Lemma ([Nesterov, Thm 2.1.2]): If  $f \in S_{\mu, L}^{1,1}$ , then

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

← (\*)

Using (\*), with  $x = z_k$  and  $y = z^*$ , noting  $\nabla f(z^*) = 0$ , we have:

$$\nabla f(z_k)^T (z_k - z^*) \geq \frac{1}{\mu + L} \|\nabla f(z_k)\|^2 + \frac{\mu L}{\mu + L} \|z_k - z^*\|^2.$$

$$\text{Then: } (\Delta) \leq \left(1 - \frac{2\mu L}{\mu + L}\right) \|z_k - z^*\|^2 + s \underbrace{\left(s - \frac{2}{\mu + L}\right)}_{\leq 0} \|\nabla f(z_k)\|^2$$

$$\Rightarrow \|z_{k+1} - z^*\|^2 \leq \left(1 - \frac{2\mu L}{\mu + L}\right) \|z_k - z^*\|^2.$$

Taking square root on both sides. done!

Now, it remains to show (\*) is true.

Since  $f \in S_{\mu, L}^{1,1}$ , we have  $\forall x, y \in \mathbb{R}^n$ :

$$\subset \mathcal{F}_{L}^{1,1}$$

$$1^\circ. f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|^2. \quad (\text{from last lecture}) \quad \textcircled{5}$$

$$2^\circ. f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2. \quad (\mu\text{-strongly convex}).$$

From 1°: interchange  $x$  &  $y$ :  $f(x) \leq f(y) + \nabla f(y)^T (x-y) + \frac{L}{2} \|y-x\|^2$ .

Adding two copies:  $(\nabla f(x) - \nabla f(y))^T (x-y) \leq L \|y-x\|^2$

By the same token on 2°, we can show,

$(\nabla f(x) - \nabla f(y))^T (x-y) \geq \mu \|x-y\|^2$ .

Now, let's pick  $x_0 \in \mathbb{R}^n$  and consider the auxiliary fn:

$$\phi(y) \triangleq f(y) - \nabla f(x_0)^T y. \Rightarrow \nabla \phi(y) = \nabla f(y) - \nabla f(x_0)$$

It's clear that  $\phi(\cdot) \in \mathcal{F}_L^{\text{li}}$ , and its opt. pt. is  $y^* = x_0$ .

(take the grad of  $\phi(\cdot)$  & set it to 0  $\Rightarrow \nabla \phi(y) = \nabla f(y) - \nabla f(x_0) = 0$   
 $\Rightarrow y^* = x_0$ ).

Thus,  $\phi(x_0) = \phi(y^*) \leq \phi(y - \frac{1}{L} \nabla \phi(y))$ .

From 1°:  $\phi(y - \frac{1}{L} \nabla \phi(y)) - \phi(y) - \frac{\nabla \phi(y)^T (y - \frac{1}{L} \nabla \phi(y) - y)}{-\frac{1}{L} \|\nabla \phi(y)\|^2} \leq \frac{\frac{1}{2} \|\nabla \phi(y)\|^2}{\frac{1}{2} \|\nabla \phi(y)\|^2}$

$\Rightarrow \underbrace{\phi(x_0)}_{f(x_0) - \nabla f(x_0)^T x_0} + \frac{1}{2L} \|\nabla f(y) - \nabla f(x_0)\|^2 \leq \phi(y) = f(y) - \nabla f(x_0)^T y$

$f(x_0) + \nabla f(x_0)^T (y - x_0) + \frac{1}{2L} \|\nabla f(y) - \nabla f(x_0)\|^2 \leq f(y)$ .

Interchange  $x_0$  &  $y$ :

$$f(y) + \nabla f(y)^T (x_0 - y) + \frac{1}{2L} \|\nabla f(y) - \nabla f(x_0)\|^2 \leq f(x_0).$$

Adding two copies (rename  $x_0$  as  $x$ ):

$$\Rightarrow \frac{1}{L} \|\nabla f(y) - \nabla f(x)\|^2 \leq (\nabla f(x) - \nabla f(y))^T (x - y). \quad \left( \begin{array}{l} \text{coercivity of} \\ \text{gradients} \end{array} \right)$$

In summary: we have

$$(1) \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2 \leq (\nabla f(x) - \nabla f(y))^T (x - y)$$

$$(2) \mu \|x - y\|^2 \leq (\nabla f(x) - \nabla f(y))^T (x - y).$$

$$\text{Thus, (1)} \times \frac{L}{\mu + L} + (2) \times \frac{L}{\mu + L} \Rightarrow$$

$$\frac{\mu L}{\mu + L} \|x - y\|^2 + \frac{L}{\mu + L} \|\nabla f(x) - \nabla f(y)\|^2 \leq \frac{2L}{\mu + L} (\nabla f(x) - \nabla f(y))^T (x - y).$$

Clearly, if  $L = \mu$ , then we are done.

So, it remains to show the case w/  $L > \mu$ .

Let's consider  $d(x) \triangleq f(x) - \frac{1}{2}\mu \|x\|^2$ . Then,  $d(x) \in \mathcal{F}_{L-\mu}^{1,1}$ .

Note:  $\nabla d(x) = \nabla f(x) - \mu x$ . (\*\*).

$$\text{Using (1): } \underbrace{(\nabla d(x) - \nabla d(y))^T}_{\downarrow} (x - y) \geq \frac{1}{L - \mu} \|\nabla d(x) - \nabla d(y)\|^2 \quad (*)$$

$$\text{Using (**): } \nabla d(x) - \nabla d(y) = [\nabla f(x) - \nabla f(y)] - \mu(x - y).$$

Plugging this into (\*) yields:

⑦

$$(L-\mu) [(\nabla f(x) - \nabla f(y)) - \mu(x-y)]^T (x-y)$$

$$\geq \|\nabla f(x) - \nabla f(y)\|^2 + \mu \|x-y\|^2 - 2\mu (\nabla f(x) - \nabla f(y))^T (x-y)$$

$$\Rightarrow -\mu \stackrel{\mu^2 = \mu L}{(L-\mu)} \|x-y\|^2 + (L-\mu) (\nabla f(x) - \nabla f(y))^T (x-y)$$

$$\geq \|\nabla f(x) - \nabla f(y)\|^2 + \cancel{\mu^2} \|x-y\|^2 - 2\mu (\nabla f(x) - \nabla f(y))^T (x-y)$$

$$\Rightarrow \cancel{(L-\mu)} (\nabla f(x) - \nabla f(y))^T (x-y) \geq \underbrace{\|\nabla f(x) - \nabla f(y)\|^2}_{L+\mu} + \mu L \|x-y\|^2$$

Dividing  $(L+\mu)$  on both sides, we're done!

□



## Useful Ineq. in Convex Analysis:

$$\begin{cases} f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|x-y\|^2, & \text{if } f \in \mathcal{F}_L^{1,1} \\ f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|x-y\|^2, & \text{if } f \in \mathcal{S}_\mu^1 \end{cases}$$

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x-y\|, \quad \text{if } f \in \mathcal{F}_L^{1,1}.$$

$$f(x) + \nabla f(x)^T (y-x) + \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|^2 \leq f(y), \quad \text{if } f \in \mathcal{F}_L^{1,1}$$

$$f(x) + \nabla f(x)^T (y-x) + \frac{1}{2\mu} \|\nabla f(y) - \nabla f(x)\|^2 \geq f(y), \quad \text{if } f \in \mathcal{S}_\mu^1$$

Interchanging & Adding:

$$\begin{cases} \frac{1}{L} \|\nabla f(y) - \nabla f(x)\|^2 \leq (\nabla f(x) - \nabla f(y))^T (x-y) \leq L \|x-y\|^2, & \text{if } f \in \mathcal{F}_L^{1,1} \\ \mu \|x-y\|^2 \leq (\nabla f(x) - \nabla f(y))^T (x-y) \leq \frac{1}{\mu} \|\nabla f(x) - \nabla f(y)\|^2, & \text{if } f \in \mathcal{S}_\mu^1. \end{cases}$$

Strongest Result: If  $f \in \mathcal{S}_{\mu,L}^{1,1}$

$$(\nabla f(x) - \nabla f(y))^T (x-y) \geq \frac{\mu L}{\mu+L} \|x-y\|^2 + \frac{1}{\mu+L} \|\nabla f(x) - \nabla f(y)\|^2.$$

Convex Combinations:  $\alpha \in [0, 1]$

$$f(\alpha x + (1-\alpha)y) + \frac{\alpha(1-\alpha)}{2L} \|\nabla f(x) - \nabla f(y)\|^2 \leq \alpha f(x) + (1-\alpha)f(y)$$

$$\leq f(\alpha x + (1-\alpha)y) + \frac{\alpha(1-\alpha)L}{2} \|x-y\|^2.$$