

# Lecture 1 Gradient Descent

- idea behind gradient descent
- smoothness
- strong convexity
- convergence analysis of GD under smoothness and strong convexity / smoothness assumptions

convex optimization problem

$$\min_{x \in C} f(x)$$

convex function

convex sets

Let's assume for simplicity that  $C = \mathbb{R}^n$

We start looking at a greedy approach

- ① start at  $x_0$
  - ② for  $i = 0$  to  $\dots$ 
    - find  $y$  s.t.  $f(y) < f(x_i)$
    - $x_{i+1} = y$ $\Rightarrow$
- Question 1  
Does it work?

## Answer

As long as we can find  $y$  s.t.  $f(y) < f(x_i)$   
yes it works. Because all local minima  
are also global minima. ( $f$  is convex)

## Question 2

How do we find  $y$  s.t.  $f(y) < f(x_i)$

(let's assume that everything is  
differentiable and fine)

since  $f(y) \approx f(x) + \nabla f(x)(y-x)$  if  
 $y$  is close to  $x$

take  $y-x = \epsilon \cdot (-\nabla f(x)) \Rightarrow$

$$\Rightarrow f(y) \approx f(x) - \epsilon \|\nabla f(x)\|_2^2 < f(x)$$

$y = x + \epsilon (-\nabla f(x))$

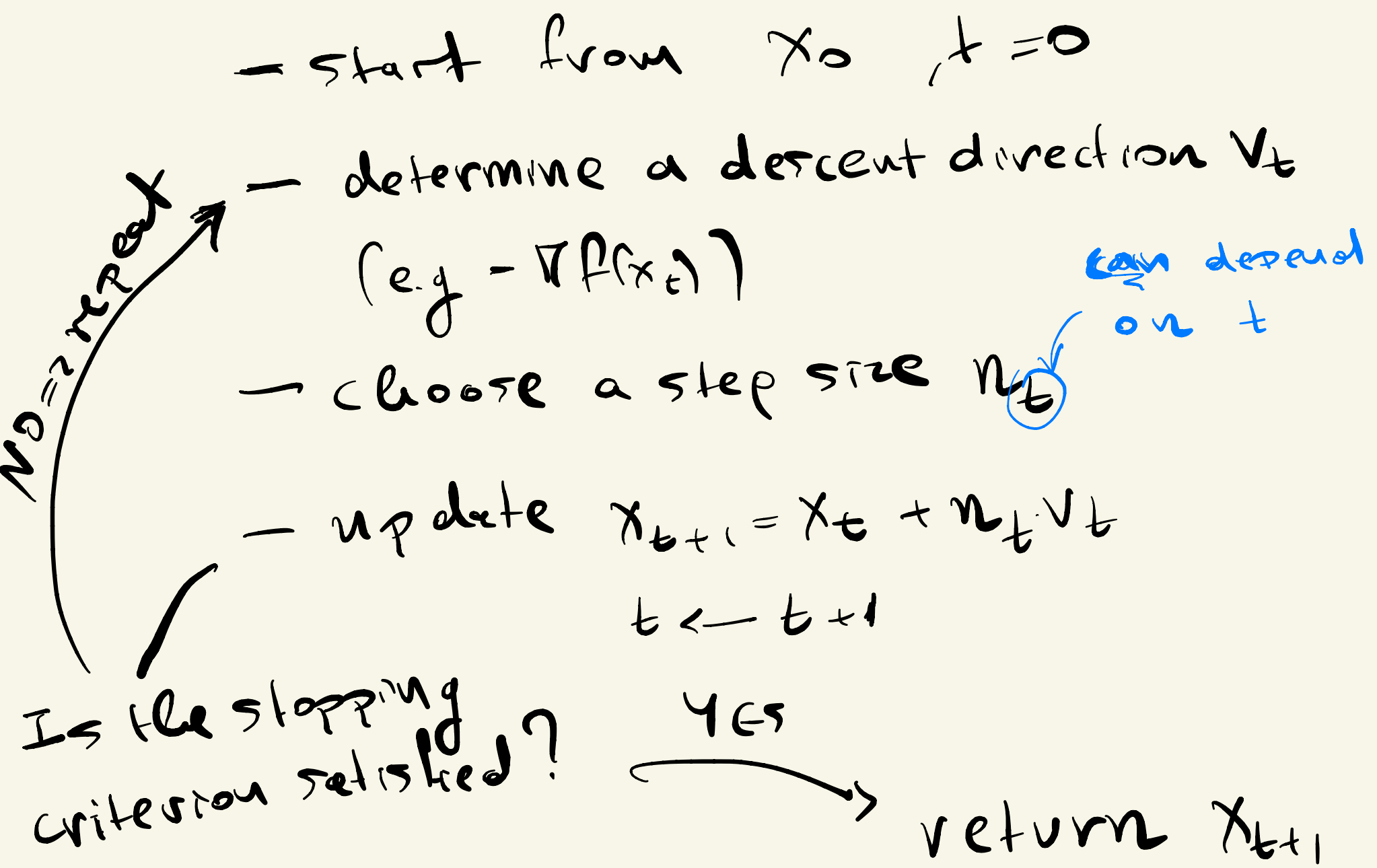
step size  $\epsilon$

descending direction

(some algorithms use different directions  $v$  as long as  $\langle v, \nabla f(x) \rangle < 0$ )



# Gradient Descent method (GD)

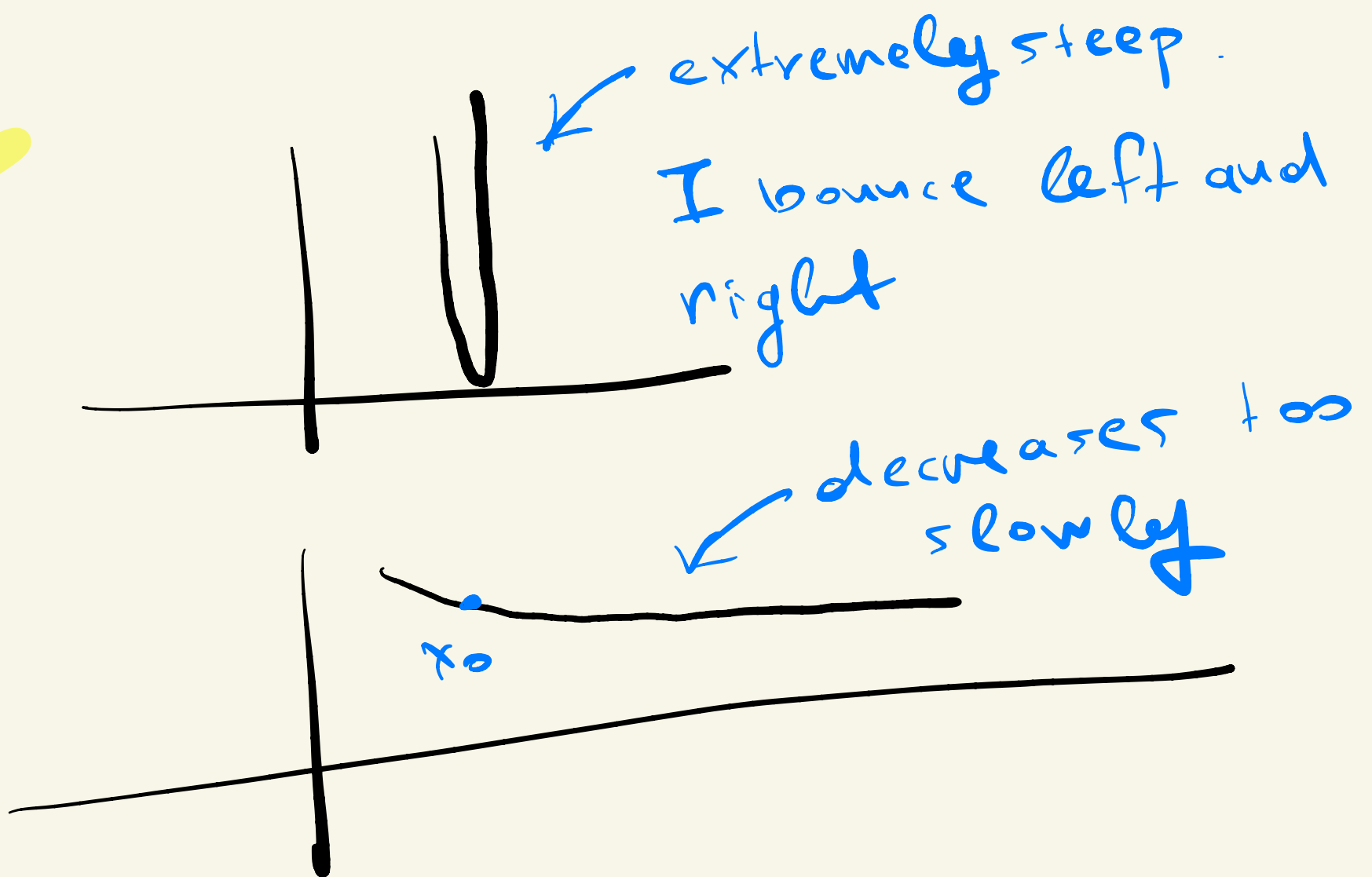


line search = choose the step size that minimizes the objective function along some fixed descent direction

why? Because it may be the case that  $\min_{\eta \in \mathbb{R}} f(x_t + \eta v_t)$  is easy!!!

Can we say something about the convergence rate without any additional assumptions?

No



## Smoothness

Def.

we say that a continuously differentiable function  $f$  is  $b$ -smooth if

for twice differentiable functions this is equivalent to  $\nabla^2 f \leq b \cdot I$

$$\|\nabla f(x) - \nabla f(y)\| \leq b \|x - y\|$$

↳ the gradient does not change too quickly

## Question.

when is the first order Taylor approximation  $f(x) + \nabla f(x)(y-x)$  of  $f(y)$  a good approximation?

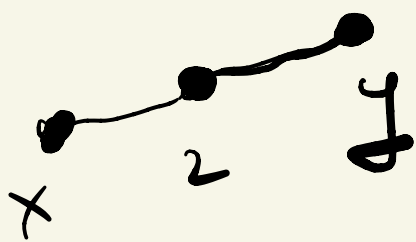
without  
any  
assumption

only if  
 $y$  is very  
close to  $x$

when  $f$  is  $L$ -smooth  
and  $L$  is small

$y$  can be more  
distant from  $x$ , because  
we know that the  
gradient changes  
slowly

elaborate  
more



is  $f(x) + \nabla f(x)(y-x)$  a good  
approximation for  $y$

$x-z$  are close

$$f(z) \approx f(x) + \nabla f(x)(z-x)$$

$z-y$  are close

$$f(y) \approx f(z) + \nabla f(z)(y-z)$$

smoothness

$$\nabla f(x) \approx \nabla f(z)$$

$$f(y) \approx f(x) + \nabla f(x)(y-x)$$

Intuitive take away point

If  $\|\nabla f(x)\|$  is large we can be aggressive in our optimization step  
(we are searching for  $x^*$  s.t.  $\|\nabla f(x^*)\| = 0$  and  $\nabla f(x)$  changes slowly)  
unconstrained case

Some claims about  $L$ -smooth functions

$$\textcircled{1} f(y) \leq f(x) + \nabla f(x)(y-x) + \frac{L}{2} \|y-x\|_2^2$$

good linear approximation

proof

$\exists z$  in the line connecting  $x$  and  $y$  s.t.

$$\begin{aligned} f(y) &= f(x) + \nabla f(x)(y-x) + \frac{1}{2} (y-x)^\top \nabla^2 f(z) (y-x) \\ &\leq f(x) + \nabla f(x)(y-x) + \frac{1}{2} (y-x)^\top L I (y-x) \\ &= f(x) + \nabla f(x)(y-x) + \frac{L}{2} \|y-x\|_2^2 \end{aligned}$$

$$\textcircled{2} f(x - \frac{1}{\textcircled{b}} \nabla f(x)) - f(x) \leq - \frac{1}{2\textcircled{b}} \|\nabla f(x)\|^2$$

taking the descent direction  
of the negative gradient  
and a step size of  $1/b$

I am sure I am gonna  
decrease!!!

small  $b$

① I can be  
very  
aggressive

② I decrease  
a lot

proof

$$\textcircled{1} \Rightarrow f(y) \leq f(x) + \nabla f(x)(y-x) + \frac{b}{2} \|y-x\|_2^2 \Rightarrow$$

$$\Rightarrow f(y) - f(x) \leq \underbrace{\nabla f(x)(y-x) + \frac{b}{2} \|y-x\|^2}_{\text{minimize the UB}}$$

minimize the UB

$$\nabla_y (UB) = \nabla f(x) + b(y-x) \Rightarrow$$

$$\Rightarrow \nabla_y (UB) = 0 \Rightarrow y = x + \left(\frac{1}{b}\right)(-\nabla f(x))$$

by replacing we get

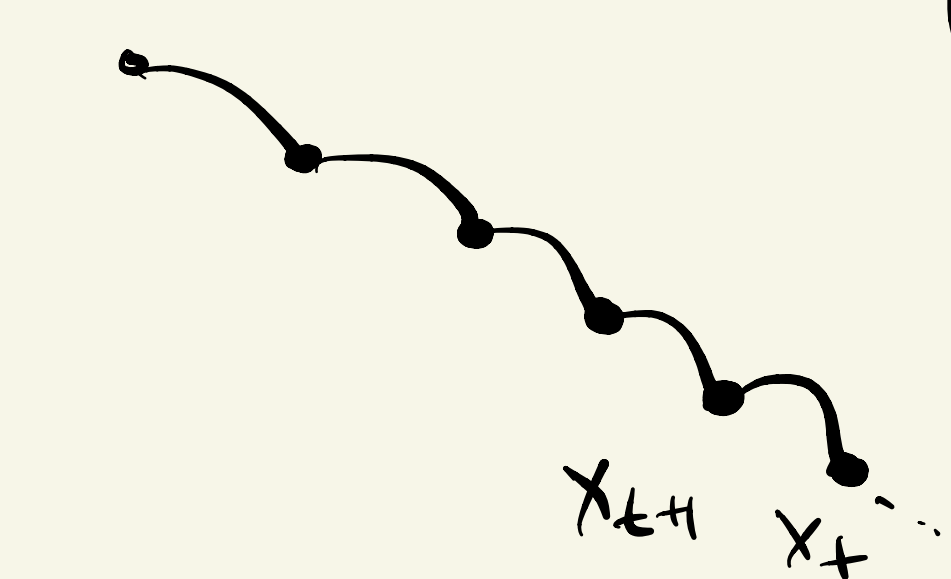
$$\begin{aligned} & \nabla f(x) \left(-\frac{\nabla f(x)}{b}\right) + \frac{b}{2} \left\| \frac{1}{b} (-\nabla f(x)) \right\|^2 = \\ & = -\frac{\|\nabla f(x)\|^2}{b} + \frac{1}{2b} \|\nabla f(x)\|^2 = -\frac{1}{2b} \|\nabla f(x)\|^2 \end{aligned}$$

$$\textcircled{3} f(y) \geq f(x) + \nabla f(x)(y-x) + \frac{1}{2L} \|\nabla f(x) - \nabla f(y)\|_2^2$$

a "little bit" better convexity.

Intuition: I increase by consecutive linear approximation by thinking that  $\nabla f(x)$  is fixed

$f(y)$



$$\nabla f(x_t)(x_{t+1} - x_t)$$

$$\approx \nabla f(x_0)(x_{t+1} - x_0)$$

$x_0$  fixed

$$f(y) - f(x) \geq \left( \nabla f(x_t)(x_{t+1} - x_t) \right) - f(x_0)$$

$$\approx \nabla f(y)\left(y - \frac{x+y}{2}\right) + \nabla f(x)\left(\frac{x+y}{2} - x\right) =$$

$$= \nabla f(y)\left(y - \frac{x}{2}\right) - \nabla f(x)\left(y - \frac{x}{2}\right) =$$

$$= (\nabla f(y) - \nabla f(x))\left(y - \frac{x}{2}\right) \approx$$

$$\|\nabla f(y) - \nabla f(x)\| \|y - x\| \geq \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|_2^2$$

let's assume that property of convexity

proof

for any  $z$  we have

$$f(x) + \nabla f(x)(z-x) \leq f(z) \leq f(y) + \nabla f(y)(z-y) + \frac{b}{2} \|z-y\|^2$$

$$\Rightarrow f(y) \geq f(x) + \nabla f(x)(z-x) - \nabla f(y)(z-y) - \frac{b}{2} \|z-y\|^2$$

(\*)

$$f(y) \geq f(x) + \nabla f(x)(y-x) + \underbrace{\nabla f(x)(z-y) - \nabla f(y)(z-y)}_{=0} - \frac{b}{2} \|z-y\|^2 =$$

$$\text{maximize } (\nabla f(y) - \nabla f(x))(z-y) + \frac{b}{2} \|z-y\|^2$$

$$\nabla_z (\swarrow) = 0 \Rightarrow (\nabla f(y) - \nabla f(x)) + bz - by = 0$$

$$\Rightarrow z = y - \frac{1}{b} (\nabla f(y) - \nabla f(x))$$

$$\text{replace } z-y = -\frac{1}{b} (\nabla f(y) - \nabla f(x)) + 0$$

(\*)



Corollary (coercivity of the gradient)

$$(\nabla f(y) - \nabla f(x))(y - x) \geq \frac{1}{2\beta} \|\nabla f(y) - \nabla f(x)\|_2^2$$

proof

from a "little bit" better convexity.

$$f(y) - f(x) \leq \nabla f(x)(y - x) - \frac{1}{2\beta} \|\nabla f(y) - \nabla f(x)\|_2^2$$

$$f(x) - f(y) \leq \nabla f(y)(y - x) - \frac{1}{2\beta} \|\nabla f(y) - \nabla f(x)\|_2^2$$

$\oplus \Rightarrow \dots$

Strong convexity

Def

we say that  $f$  is  $\alpha$ -strongly convex

$$\text{if } f(y) \geq f(x) + \nabla f(x)(y - x) + \frac{\alpha}{2} \|y - x\|_2^2$$

Intuitive take-away point

when  $\|\nabla f(x)\|$  is small

stop. (because  $\|y - x\|^2$  dominates the RHS and consequently slightly move)

$f(y) \geq f(x)$  if we move optimal solution is closer

$$\nabla^2 f(x) \succeq \alpha \cdot I$$

$f$  twice differentiable

in 0-smoothness we have

$$\frac{1}{2\beta} \|\nabla f(y) - \nabla f(x)\|_2^2$$

## claims

$$\textcircled{1} \quad \frac{1}{2\alpha} \|\nabla f(x)\|_2^2 \geq f(x) - f(x^*)$$

If gradient is small we  
are good!!!

## proof

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

$$\Rightarrow f(x) - f(y) \leq \underbrace{-\frac{\alpha}{2} \|y-x\|^2 - \nabla f(x)(y-x)}_{\text{maximize}}$$

$$\nabla_y (u^*) = 0 \Rightarrow -\alpha(y^*-x) - \nabla f(x) = 0$$

$$\Rightarrow y^* = x - \frac{1}{\alpha} \nabla f(x)$$

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

holds  $\forall y \Rightarrow y = x^*$  ← optimal solution

$$\textcircled{2} \quad \|x - x^*\|_2 \leq \frac{2}{\alpha} \|\nabla f(x)\|_2^2$$

If the gradient is small  
not only  $f(x) \approx f(x^*)$  but  
 $x$  is actually close to  $x^*$

proof

$$f(x^*) \geq f(x) + \nabla f(x)(x^* - x) + \frac{\alpha}{2} \|x - x^*\|^2$$

$$\geq f(x) - \|\nabla f(x)\| \|x^* - x\| + \frac{\alpha}{2} \|x^* - x\|^2$$

↑  
Cauchy-Schwarz

$$f(x^*) - f(x) \leq 0 \quad (x^* \text{ is a minimizer})$$

$\Rightarrow \textcircled{2}$

Other common assumptions

L-Lipschitz

$$|f(x) - f(y)| \leq L \|x - y\|$$

↑  
my function does not change  
very fast

Diameter

$$\forall x, y \in \mathcal{C}$$

$$\|x - y\| \leq R$$

↑  
I am always at most  $R$  far  
away from  $x^*$

# convergence analysis of GD

①  $a$ -strongly convex  $(b > a > 0)$   
 $b$ -smooth

$$x_{t+1} = x_t - \frac{1}{b} \nabla f(x_t)$$

$$\text{then } f(x_t) - f(x^*) \leq \left(1 - \frac{a}{b}\right)^t (f(x_0) - f(x^*))$$

$$\stackrel{e^{x_2, 1+x_1, \dots, x_n}}{\leq} e^{-\frac{a}{b}t} (f(x_0) - f(x^*)) \leq$$

To ensure

$$f(x_t) - f(x^*) < \varepsilon$$

$$\stackrel{\leq}{\leq} e^{-\frac{a}{b}t} \left(\frac{b}{2}\right) \|x_0 - x^*\|_2^2$$

$\hookrightarrow b$ -smoothness  
definition and

$$\nabla f(x^*) = 0$$

$$T = \left(\frac{b}{a}\right) \ln \left(\frac{f(x_0) - f(x^*)}{\varepsilon}\right)$$

$\downarrow$   
condition  
numbers  
(bigger when  
it is small)

$$f(x_0) \leq f(x^*) + \nabla f(x^*)^T (x^* - x_0) + \frac{b}{2} \|x^* - x_0\|_2^2$$

$$O(\log 1/\varepsilon)$$

Proof  $f(x_t) - f(x^*) \leq \left(1 - \frac{\alpha}{L}\right)^t (f(x_0) - f(x^*))$

Idea

In every iteration  $x_t \rightarrow x_{t+1}$

① we decrease our objective by

$$\frac{1}{2L} \|\nabla f(x_t)\|^2$$

② we are at most  $\frac{1}{2\alpha} \|\nabla f(x_{t+1})\|^2$  far away from the optimal solution

$$f(x_{t+1}) \leq f(x_t) - \frac{1}{2L} \|\nabla f(x_t)\|^2 =$$

$$= f(x_{t+1}) - f(x^*) \leq f(x_t) - f(x^*) - \frac{1}{2L} \underbrace{\|\nabla f(x_t)\|^2}_{\geq 2\alpha(f(x_t) - f(x^*))}$$

$$\Rightarrow f(x_{t+1}) - f(x^*) \leq \left(1 - \frac{\alpha}{L}\right) (f(x_t) - f(x^*))$$

## ② only $L$ -smoothness

$$x_{t+1} = x_t - \frac{1}{L} \nabla f(x_t)$$

then  $f(x_t) - f(x^*) \leq \frac{16L \|x_0 - x^*\|^2}{T}$

proof (helpful notation)  $\delta_t = f(x_t) - f(x^*)$

$$\delta_{t+1} - \delta_t \leq -\frac{1}{2L} \|\nabla f(x_t)\|^2$$

$$f(x_{t+1}) - f(x_t) \leq -\frac{1}{2L} \|\nabla f(x_t)\|^2 \quad (\text{by the smoothness assumption, claim ② for smooth functions})$$

by first order condition

$$\begin{aligned} f(x_t) - f(x^*) &\leq \nabla f(x_t) (x_t - x^*) \leq \\ &\leq \|\nabla f(x_t)\| \|x_t - x^*\| \Rightarrow \end{aligned}$$

$$\|\nabla f(x_t)\| \geq \frac{f(x_t) - f(x^*)}{\|x_t - x^*\|}$$

decreases  
let's be optimistic and assume that decreases

then  $\|\nabla f(x_t)\| \geq \frac{f(x_t) - f(x^*)}{\|x_0 - x^*\|} = c \cdot \delta_t$

$$\delta_{t+1} - \delta_t \leq -\frac{c^2 \delta_t^2}{2L}$$

$$(c = \frac{1}{\|x_0 - x^*\|})$$

$$\delta_{t+1} - \delta_t \leq -\frac{c^2}{2b} \delta_t^2$$

how much time I do need to  
half the distance  $\delta_t \approx \frac{\delta_+}{2}$ ?

Before I get  $\delta_t' < \frac{\delta_+}{2}$

I have a decrease of at least

$$\frac{c^2}{2b} \frac{\delta_+^2}{4} = \frac{c^2 \delta_+^2}{8b} \text{ at each step}$$

$$\Rightarrow T \cdot \frac{c^2 \delta_+^2}{8b} < \frac{\delta_+}{2} \Rightarrow T < \frac{4b}{c^2 \delta_+}$$

If I need to half my error  $\delta$  times  $\Rightarrow$

$$\Rightarrow T_{\text{total}} < \frac{4b}{c^2 \delta_0} + \frac{4b}{c^2} \left(\frac{2}{\delta_0}\right) + \dots + \frac{4b}{c^2} \cdot \frac{2^J}{\delta_0} =$$

$$= O\left(\frac{4b}{c^2} \cdot \frac{2^{J+1}}{\delta_0}\right) \Rightarrow \frac{c}{\|x^* - x_0\|^2}$$

$$\Rightarrow \frac{\delta_0}{2^J} = O\left(\frac{b \cdot \|x^* - x_0\|^2}{T_{\text{total}}}\right) \Rightarrow f(x_{T_{\text{total}}}) - f(x^*) = O\left(\frac{b \cdot \|x^* - x_0\|^2}{T_{\text{total}}}\right)$$



Last claim to prove

$\|x_t - x^*\|^2$  is non decreasing

proof

$$\|x_{t+1} - x^*\|^2 = \|x_t - \frac{1}{b} \nabla f(x_t) - x^*\|^2$$

$$= \|x_t - x^*\|^2 + \frac{1}{b^2} \|\nabla f(x_t)\|^2$$

$$- \frac{2}{b} \nabla f(x_t) (x_t - x^*) \leq$$

$$\leq \|x_t - x^*\|^2 - \frac{1}{b^2} \|\nabla f(x_t)\|^2$$

by coercivity and  $\nabla f(x^*) = 0$   
( $\nabla f(y) - \nabla f(x) \geq \frac{1}{b} \|\nabla f(y) - \nabla f(x)\|_2^2$ )

Now we summarize the performance of GD for different properties of the convex  $f$  that we are optimizing over, and we compare it to Nesterov accelerated method together with lower bounds on the performance of any optimization algorithm.

# iterations to get an $\varepsilon$ -additive error	GD	Nesterov's accelerated method	LB
$L$ -Lipschitz convex $f$	$O(\frac{L^2}{\varepsilon^2})$		$O(\frac{L^2}{\varepsilon^2})$
$b$ -smooth convex $f$	$O(\frac{b}{\varepsilon})$	$O(\sqrt{\frac{b}{\varepsilon}})$	$O(\sqrt{\frac{b}{\varepsilon}})$
$b$ -smooth and $\alpha$ -strongly convex $f$	$O(\frac{b}{\alpha} \log \frac{1}{\varepsilon})$	$O(\sqrt{\frac{b}{\alpha}} \log \frac{1}{\varepsilon})$	$O(\sqrt{\frac{b}{\alpha}} \log \frac{1}{\varepsilon})$

we hide the dependence on the initial distance from the optimum ( $\|x_0 - x^*\|$ ) in the  $O(\cdot)$  notation. Here we focus

mainly on the dependence with  $\varepsilon$ .

Although, as we will see for the min-cut problem, getting a better dependence w.r.t.  $\|x^* - x_0\|$  is important, and Nesterov's accelerated method also achieves that.

For a table summarizing many more convergence rates of various optimization algorithms. Please see the end of chapter 1 in Bebeck's book (Convex Optimization: Algorithms and Complexity)