

- \* Note that projected GD can be modified in the same way for  $L$ -smooth functions to obtain the same convergence rate when solving

$$\min_{x \in \Omega} f(x) \quad \text{where } \Omega \subset \mathbb{R}^n \text{ is convex.}$$



### Picking $\mu$ in Practice:

- The convergence theorems that we've seen so far require knowing  $L$ , the Lipschitz constant (or the smoothness constant) in order to set  $\mu$ , but we may not know  $L$  in practice.
- Idea to get around this issue:

Use the best  $\mu^{(t)}$  at every iteration.

$$x^{(t+1)} = x^{(t)} - \mu^{(t)} \nabla f(x^{(t)})$$

$\Rightarrow$  Pick  $\mu^{(t)}$  to minimize:

$$f(x^{(t+1)}) = f\left(\underbrace{x^{(t)}}_{\text{known}} - \underbrace{\mu^{(t)}}_{\text{the variable we're optimizing}} \underbrace{\nabla f(x^{(t)})}_{\text{known}}\right)$$

Issue: Solving for  $\mu^{(t)}$  exactly is often hard, so we often settle for an approximation

Possible sol'n: Use, at every iteration  $t$ , a "backtracking Line Search"

How does this work?

Note that for any descent direction  $\vec{p}$ :

$$f(x) - \nabla f(x)^T p \leq \underline{f(x-p)} \leq \underbrace{f(x) - \gamma \nabla f(x)^T p}_{(*)}$$

for some small  $\gamma > 0$ .

Pick  $\mathcal{P} = \mu \nabla f(x)$ , as we do in GD

and plug it into (\*):

$$\underbrace{f(x - \mu \nabla f(x)) \leq f(x) - \delta \mu \|\nabla f(x)\|^2}_{(**)}$$

So we expect that for  $\mu$  that is small enough, (\*\*) should hold.

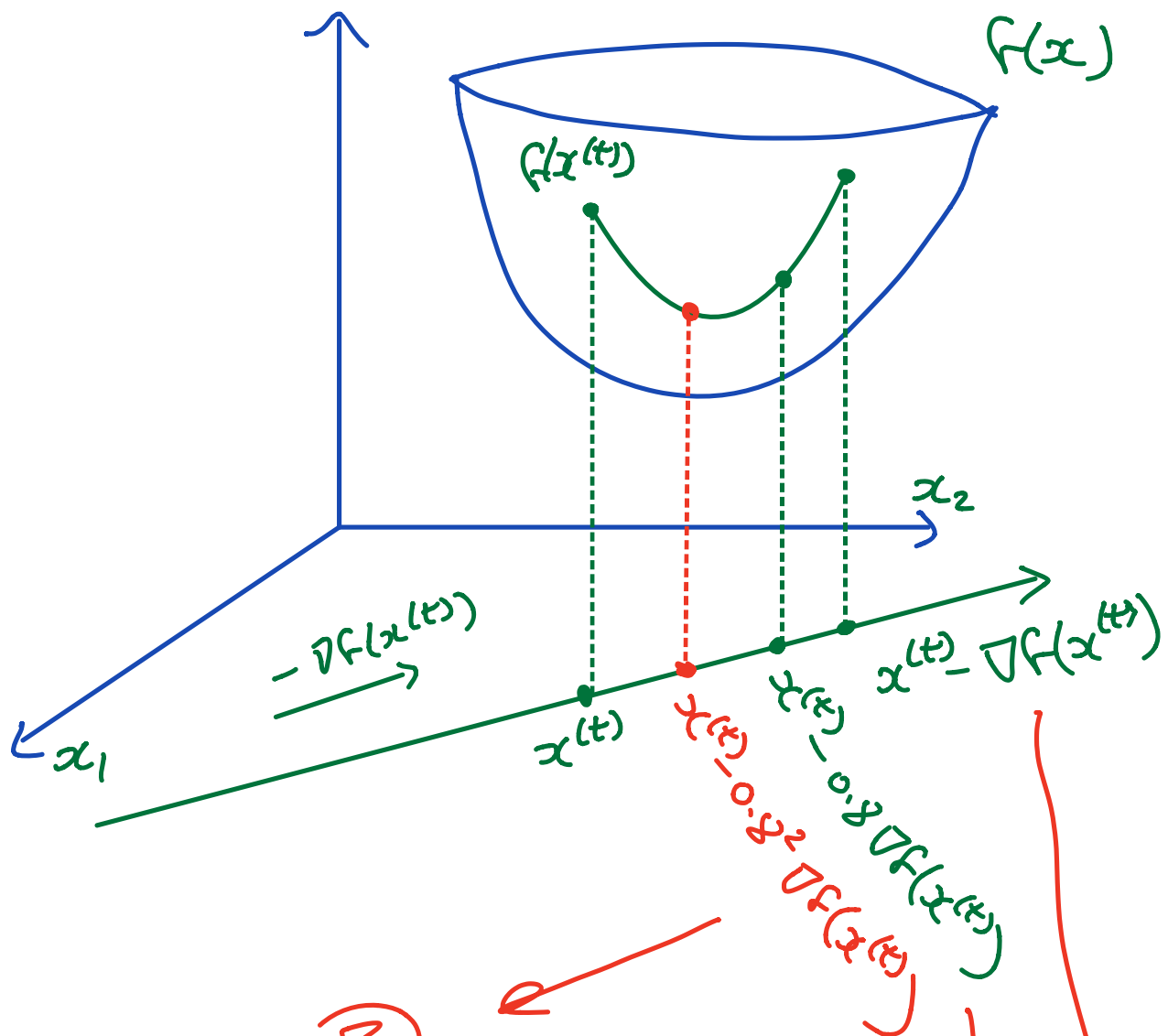
Idea: • Fix  $\delta$ , e.g.  $\delta = 0.5$ ,

- Start with  $\mu = 1$  (for example)
- decrease  $\mu$  iteratively till

$$\begin{aligned} f(x^{(t+1)}) &= f(x^{(t)} - \mu^{(t)} \nabla f(x^{(t)})) \\ &\leq \underbrace{f(x^{(t)}) - \mu^{(t)} \delta \|\nabla f(x^{(t)})\|^2}_{\text{Armijo Condition}} \end{aligned}$$

For example

$$\begin{array}{ccccc} \mu^{(t)} = 1 & \rightarrow & \mu^{(t)} = 0.8 & \rightarrow & \mu^{(t)} = 0.8^2 \\ \times & & \times & & \checkmark \end{array}$$



now, we're good

③  $f(x^{(t)} - 0.8 \nabla f(x^{(t)}))$  is too big  
↓  
still

②  $f(x^{(t)} - 1 \times \nabla f(x^{(t)}))$  is too big

①

Then we discussed Backtracking Line search for selecting  $\mu$ .

→ Pick  $\beta < 1, \delta < 1$

At each GD step  $t$ :

① Set  $v = -\nabla f(x^{(t)})$

② Set  $\mu^{(t)} = 1$

③ If  $f(x^{(t)} - \mu^{(t)} \nabla f(x^{(t)})) \leq f(x^{(t)}) - \mu \delta \|\nabla f(x^{(t)})\|^2$

then: keep  $\mu^{(t)}$

Else: Set  $\mu^{(t)} \leftarrow \beta \mu^{(t)}$   
and repeat ③

Example : Consider  
 $f(x_1, x_2) = (x_1 - 1)^4 + (x_1 + x_2 - 1)^2$   
(and note that  $x^* = (1, 0)$ )

$$\nabla f = (4(x_1 - 1)^3 + 2(x_1 + x_2 - 1), 2(x_1 + x_2 - 1))$$

Suppose that  $x^{(0)} = (0, 0)$

$$\begin{aligned} x^{(1)} &= x^{(0)} - \mu \nabla f(x^{(0)}) \\ &= (0, 0) - \mu (-6, -2) \end{aligned}$$

Ideally, want to pick  $\mu$  to minimize

$$\begin{aligned} f((0, 0) - \mu(-6, -2)) &= f(6\mu, 2\mu) \\ &= (6\mu - 1)^4 + (8\mu - 1)^2 \end{aligned}$$

Picking  $\mu$  this way entails solving a nonlinear opt. problem, which may be hard.

So let's use backtracking line search :

We start with  $\mu = 1$  :

$$f(x^{(1)}) = 674 > \underbrace{f(x^{(0)})}_{2} - \underbrace{\mu}_{1} \underbrace{\delta}_{0.5} \underbrace{\| \nabla f(x^{(0)}) \|^2}_{40}$$

-18

$x^{(0)} - \mu \nabla f(x^{(0)})$

So we try  $\mu = 1 \times 0.8 = 0.8$

$$f(x^{(1)}) = 237.6736 > \underbrace{f(x^{(0)}) - 0.8 \times 0.5 \| \nabla f(x^{(0)}) \|^2}_{-14}$$

So we try  $\mu = 0.8^2$   
(also fails)

So we keep going until  $\mu = 0.8''$

gives

$$f(x^{(1)}) = 0.1530 \leq f(x^{(0)}) - 0.5 \times 0.8'' \| \nabla f(x^{(0)}) \|^2$$

So we choose this  $\mu$  ( $\mu = 0.8''$ )

$$\text{and set } x^{(1)} = \underbrace{x^{(0)}}_{(0,0)} - \underbrace{\mu}_{0.8''} \underbrace{\nabla f(x^{(0)})}_{\nabla f(0,0)}$$

and we repeat this process for  
 $t = 2, 3, \dots$

(for example  $x^{(1000)} = (0.979, 0.021)$ )

Theorem: For an  $L$ -smooth convex function with  $\mu^{(t)}$  set by backtracking line search, GD gives

$$f(x^{(t)}) - f(x^*) \leq \frac{1}{2t \min_{s=1, \dots, t} \mu^{(s)}}$$

$$\text{and } \min_{s=1, \dots, t} \mu^{(s)} \geq \min\left(1, \frac{\beta}{L}\right)$$

↳ this guarantees

$$f(x^{(t)}) - f(x^*) \leq \frac{L}{2t \beta}$$