

Claim: $-\epsilon \frac{\nabla f(x_k)}{\|\nabla f(x_k)\|} = \operatorname{argmin}_s \{ \nabla f(x_k)^T s \mid \|s\|_2 \leq \epsilon \}$

Proof. Suppose for contradiction that s' does better

$$\begin{aligned} \nabla f(x_k)^T s' &< \nabla f(x_k)^T \left(\frac{-\epsilon \nabla f(x_k)}{\|\nabla f(x_k)\|_2} \right) \\ &= -\epsilon \|\nabla f(x_k)\|_2 \end{aligned}$$

By Cauchy-Schwarz,

$$-\|\nabla f(x_k)\|_2 \|s'\|_2 < -\epsilon \|\nabla f(x_k)\|_2$$

$$\Rightarrow \|s'\|_2 > \epsilon.$$

$\Rightarrow s'$ is not feasible. \square

$$\Rightarrow x_{k+1} = x_k - \alpha_k \nabla f(x_k) \quad \text{for some } \alpha_k.$$

"Gradient Descent".

Starting 9/9

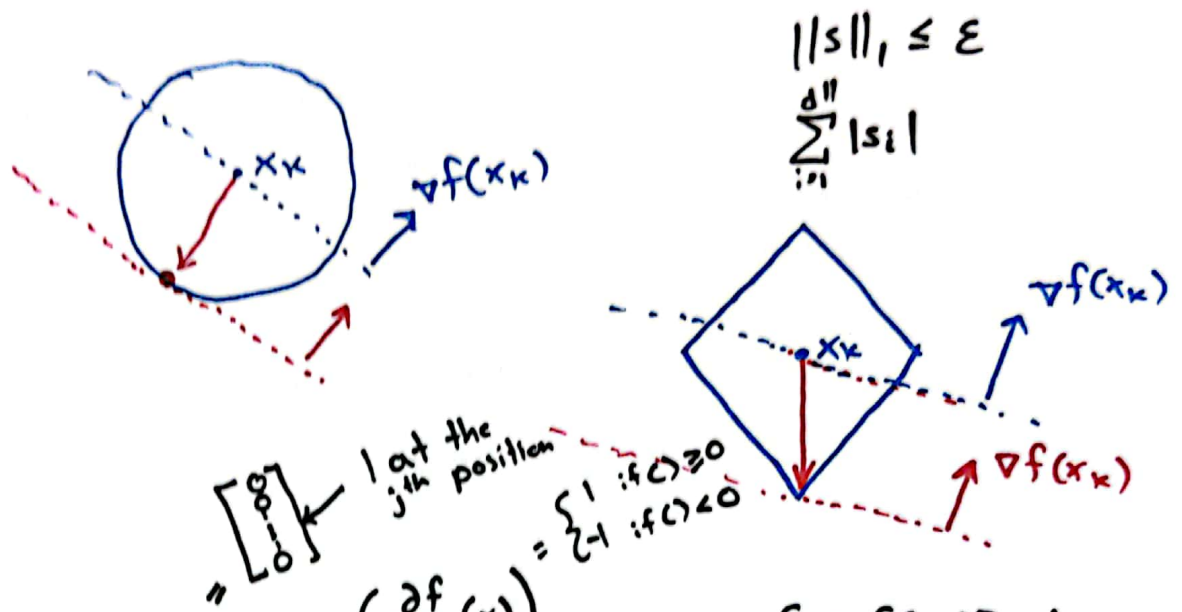
HW1 due before the next lecture.

Submit on time. I will release my solutions
and then cannot accept yours...

Approach 3

We could consider other norms.

We could look for best improvement with



Claim $s = -\epsilon e_j \operatorname{sign}\left(\frac{\partial f}{\partial x_j}(x_k)\right) \in \operatorname{argmin} \{ \nabla f(x_k)^T s \mid \|s\|_1 \leq \epsilon \}$
where $j \in \operatorname{argmax} \{ \left| \frac{\partial f}{\partial x_j}(x_k) \right| \}$

Proof. Suppose s' does better:

$$\nabla f(x_k)^T s' < \nabla f(x_k)^T s$$

$$= -\epsilon \frac{\partial f}{\partial x_j}(x_k) \operatorname{sign}\left(\frac{\partial f}{\partial x_j}(x_k)\right)$$

$$= -\epsilon \left| \frac{\partial f}{\partial x_j}(x_k) \right|$$

$$= -\epsilon \|\nabla f(x_k)\|_\infty \quad \left(\|y\|_\infty = \max |y_i| \right)$$

By Hölder's Inequality

$$-\|\nabla f(x_k)\|_\infty \|s'\|_1 < -\epsilon \|\nabla f(x_k)\|_\infty$$

$$\Rightarrow \|s'\|_1 > \epsilon.$$

$\Rightarrow s'$ is not feasible.

□

This motivates "Coordinate Descent"

$$x_{k+1} = x_k - \alpha_k e_j \operatorname{sign}\left(\frac{\partial f}{\partial x_j}(x_k)\right)$$

↑ the coordinate with largest partial derivative.

If $x_0 = 0$, x_k is k -sparse
("k nonzeros").

Second-Order Local Improvements

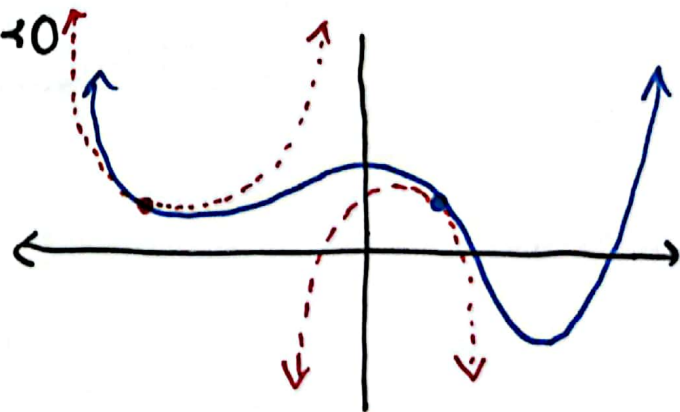
$$f(x_k + s) \approx f(x_k) + \nabla f(x_k)^T s + \frac{1}{2} s^T \nabla^2 f(x_k) s$$

Approach 4

$$\min_{s \in \mathbb{R}^d} \left(\right)$$

$$= \begin{cases} -\infty & \text{if } \nabla^2 f(x_k) \prec 0 \\ ??? & \end{cases}$$

↑ promising



If $\nabla^2 f(x_k) \succ 0$, then

$$x_{k+1} = x_k - \nabla^2 f(x_k)^{-1} \nabla f(x_k)$$

"Newton's Method"
for $\nabla f(x_k) = 0$

Approach 5 Trust - Region Methods

$$x_{k+1} = x_k + \underset{\|s\| \leq \epsilon}{\operatorname{argmin}} \left\{ f(x_k) + \nabla f(x_k)^T s + \frac{1}{2} s^T \nabla^2 f(x_k) s \right\}$$

How can we guarantee good quality?

Taylor Approximation Theorems.

Gradient Descent/Smooth Optimization

4 lectures

1. A Descent Lemma
2. Stepsizes/Line searches
3. Nonconvex Smooth Opt Guarantees
4. Shape of Smooth Convex Funcs
5. Better Guarantees
6. Complexity Lower bounds
7. Acceleration

1. A Descent Lemma

We defined gradient descent (GD) as

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k)$$


for some $x_0 \in \mathbb{R}^d$ and $\alpha_k \in \mathbb{R}$.

Lemma For any f with L -Lipschitz gradient, any $k \geq 0$ has


$$f(x_{k+1}) \leq f(x_k) - \left(\alpha_k - \frac{L\alpha_k^2}{2} \right) \|\nabla f(x_k)\|_2^2.$$

Consequences: Decrease whenever $\alpha_k - \frac{L\alpha_k^2}{2} > 0$
 $\Leftrightarrow \alpha_k < \frac{2}{L}$.

Better descent when L small
 $\|\nabla f(x_k)\|_2$ large.

Best decrease when $\alpha_k = \frac{1}{L}$ 
of $\frac{1}{2L} \|\nabla f(x_k)\|_2^2$

(Taylor Approx Theorem)

Proof. $|f(x_{k+1}) - (f(x_k) + \nabla f(x_k)^T (x_{k+1} - x_k))| \leq \frac{L}{2} \|x_{k+1} - x_k\|_2^2$ 

- dropping l.o.
- plug in
 $x_{k+1} - x_k$
 $= -\alpha_k \nabla f(x_k)$

$$\Rightarrow f(x_{k+1}) - f(x_k) - \alpha_k \|\nabla f(x_k)\|_2^2 \leq \frac{L\alpha_k^2}{2} \|\nabla f(x_k)\|_2^2$$

$$\Rightarrow f(x_{k+1}) \leq f(x_k) - \left(\alpha_k - \frac{L\alpha_k^2}{2} \right) \|\nabla f(x_k)\|_2^2.$$

□

2. Stepsize Choice / Linesearch

Based on our lemma, $\alpha_k = 1/L$ gives

$$\frac{1}{2L} \|\nabla f(x_k)\|_2^2.$$

Impractical. We do not often know L .

Exact Linesearch

Pick the best α_k

$$\alpha_k = \operatorname{argmin}_{\alpha \in \mathbb{R}} f(x_k - \alpha \nabla f(x_k))$$

Trivially outperforms $\alpha_k = 1/L$ since

$$\begin{aligned} \text{(namely, } f(x_{k+1}) &\leq f(x_k - \alpha \nabla f(x_k)) \quad \forall \alpha \\ &\leq f(x_k - \frac{1}{L} \nabla f(x_k)) \\ &\leq f(x_k) - \frac{1}{2L} \|\nabla f(x_k)\|_2^2 \quad \text{by Descent Lemma} \end{aligned}$$

Impractical. Slowdown our iteration.

Backtracking Linesearch

Pick $\alpha \in \mathbb{R}$, $\tau \in (0, 1)$

$$\alpha_k = \sup \left\{ \alpha \tau^n \mid n=0, 1, 2, \dots \right. \\ \left. f(x_k - \alpha \tau^n \nabla f(x_k)) < f(x_k) \right\}$$

"keep backing off exponentially until decrease is found."

Does this terminate? Yes!

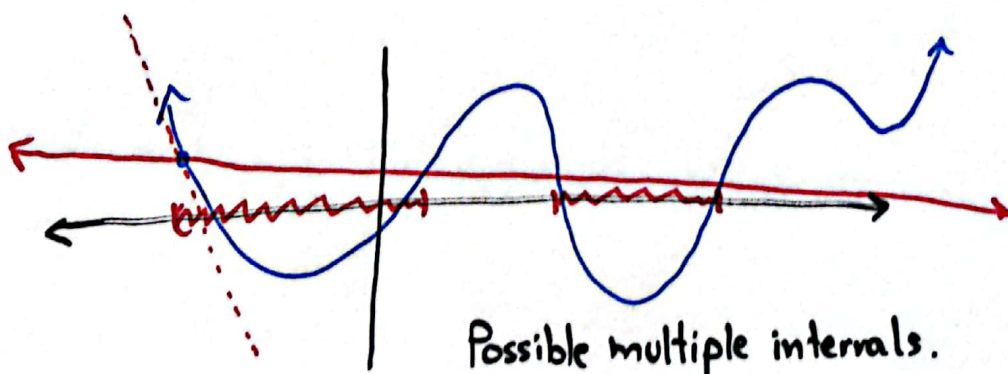
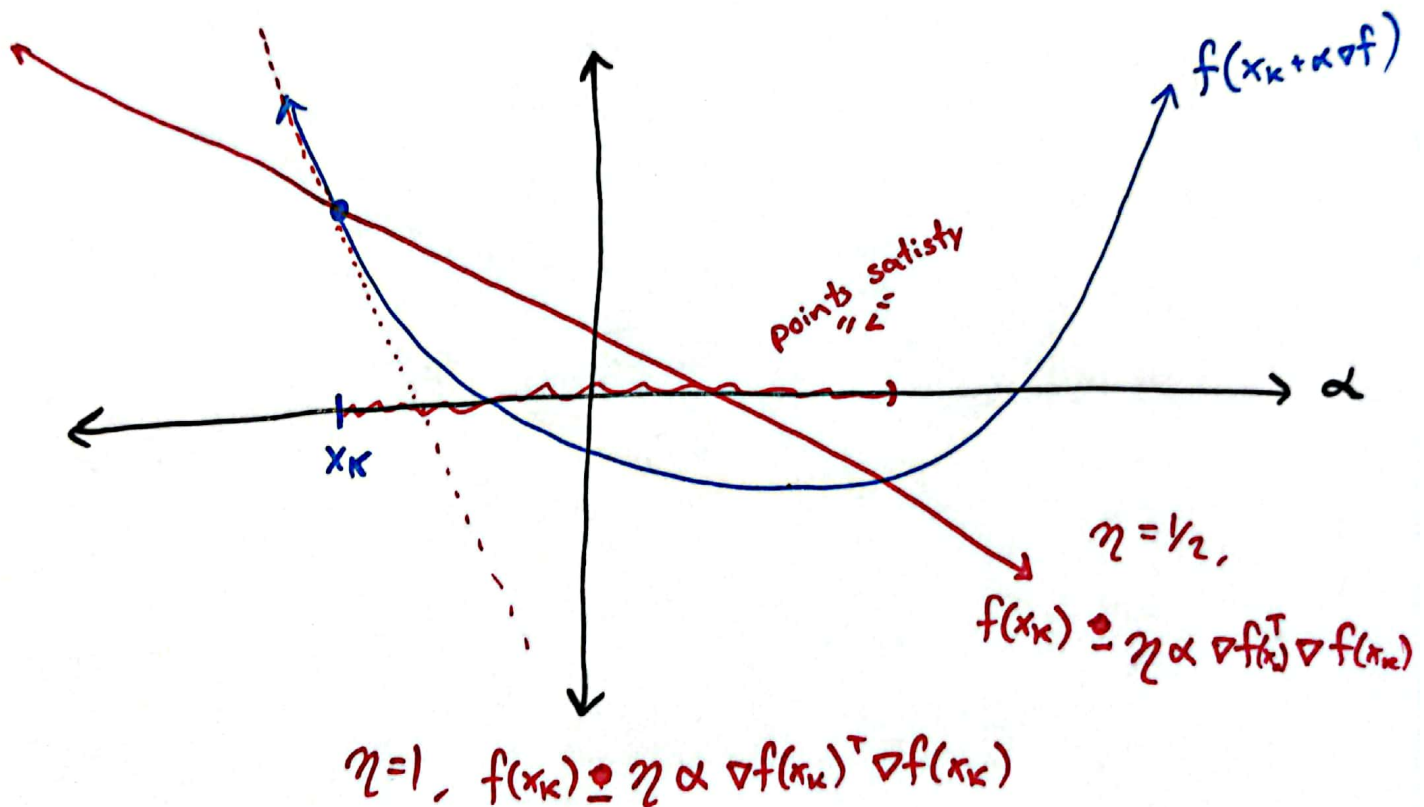
Descent Lemma ensures decrease $\alpha \in (0, \frac{2}{L})$.

Only $\lceil \log_{1/2} \left(\frac{\alpha}{2/L} \right) \rceil$ backtracking steps.

How should we measure descent? What is " \angle "?

Pick $\eta \in (0, 1)$, we want

"Armijo Condition" $f(x_k - \alpha_k \nabla f(x_k)) \leq f(x_k) - \eta \alpha_k \|\nabla f(x_k)\|_2^2$



Lemma The Armijo Condition holds for
 $\alpha \in \left[0, \frac{2(1-\eta)}{L}\right].$