

Amarth SIS, Jan 20

- * β -smoothness
- * upper bounds
- * Analysis of gradient descent for β -smooth fns

For If w:

For problem 1, I recommend
starting with simple facts:

$$\int c(x) (g(h(x)))' = g'(h(x)) h'(x)$$

for scalar funcs.

(2) Jacobian of $x \rightarrow Ax$
is A.

build gradient, hessian formulas

for $g(a^T x)$ first, then
 $g(Ax)$

$$\min_x \underbrace{\sum \log(1 + \exp(a_i^T x))}_{\frac{\partial}{\partial z} (Ax)}$$

$$g(z) = \sum_i \log(1 + \exp(z_i))$$

Easy to differentiate

$$\nabla g(z) = \begin{bmatrix} \frac{\exp(z_1)}{1 + \exp(z_1)} \\ \vdots \\ \frac{\exp(z_n)}{1 + \exp(z_n)} \end{bmatrix}$$

derivative

$$\nabla^2 g = \begin{bmatrix} \cdot & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & \ddots \end{bmatrix}$$

Last time: ended with rates

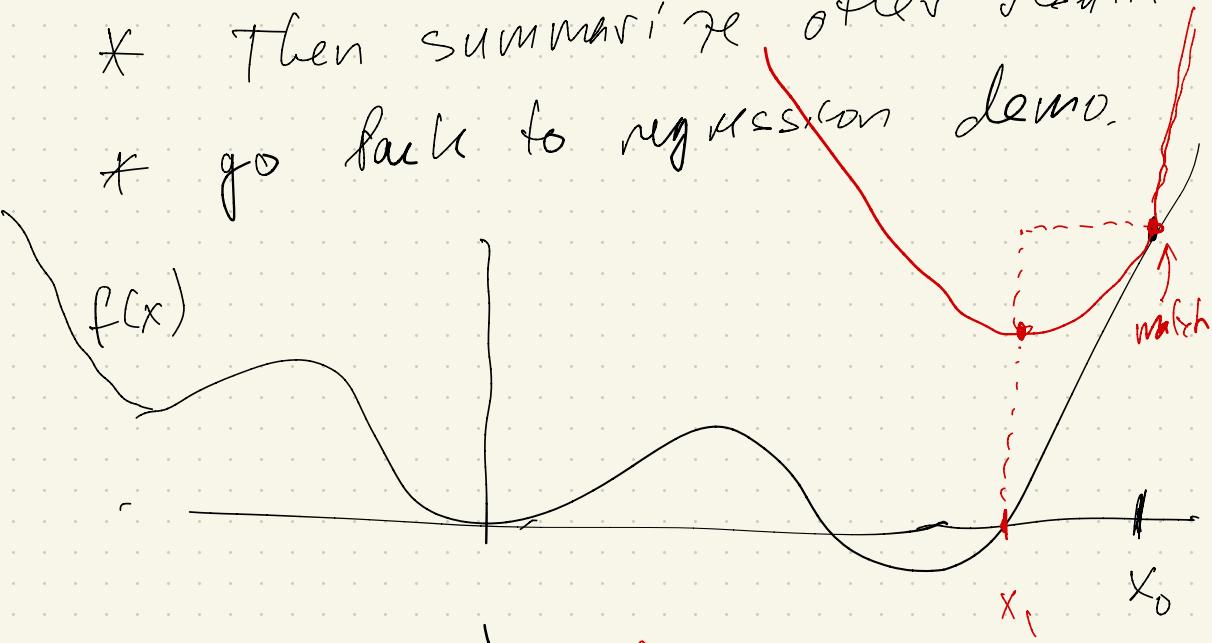
From sublinear to quadratic

Today: analyze gradient descent.

$$x_{n+1} = x_n - (\text{step}) \nabla f(x_n)$$

$$f(x) \leftarrow \min_x f(x)$$

- * Then summarize after result,
- * go back to regression demo.



$$\} \quad \text{get } f(x_1) < f(x_0)$$

for a global upper bound

Strategy:

- ① Build tight upper bounds
- match f and gradient at current point
 - above everywhere
- ② get iterates by minimizing these upper bounds.

Need a bit more notation!

Def: A function $F: \mathbb{R}^n \rightarrow \mathbb{R}^M$ is β -Lipschitz continuous if

$$\|F(x) - F(y)\| \leq \beta \|x - y\|$$

A function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is

β -smooth if $f \in C^1$ and

∇f is β -Lipschitz cont:

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

Thm 1: If f is β -smooth, then

$$\rightarrow f(y) \leq f(x) + \underbrace{\langle \nabla f(x), y-x \rangle}_{\text{tangent line approx at } x} + \frac{\beta}{2} \|y-x\|^2$$

\uparrow
quadratic

Thm 2: For $f \in C^2$, f is β -smooth

if and only if $\nabla^2 f \leq \beta I$

equivalently, largest eigenvalue of $\nabla^2 f \leq \beta$

Ex: $\min_x \frac{1}{2} \|Ax-b\|^2$

$$\nabla f(x) = \underbrace{A^T(Ax-b)}_{\text{linear in } x}$$

So f is β -smooth, with β the largest eigenvalue of $A^T A$

Since $\nabla^2 f(x) = A^T A$.

Main point: β -smoothness gives us very simple global tight upper bounds.

$$\text{Call } m_x(y) = f(x) + \langle \nabla f(x), y - x \rangle + \frac{\beta}{2} \|y - x\|^2$$

① Know $f(y) \leq m_x(y)$ everywhere, by theorem 1.

② $m_x(x) = f(x)$ (yay)

③ $\nabla m_x(y) = \nabla f(x) + \beta(y - x)$
at $y = x$, $\nabla m_x(x) = \nabla f(x)$

Strategy: minimize m_x at each iteration

$$x^+ = \underset{y}{\operatorname{argmin}} m_x(y) \quad \text{'next point'}$$

$$x_{k+1} = \underset{y}{\operatorname{argmin}} m_{x_k}(y) \quad \text{'k-th iterate'}$$

Let's minimize $m_x(y)$ by completing the square

$$\begin{aligned}
 m_x(y) &= f(x) + \underbrace{\langle \nabla f(x), y-x \rangle}_{\text{blue}} + \frac{\beta}{2} \|y-x\|^2 \\
 &= f(x) + \beta \left(\underbrace{\langle \frac{1}{\beta} \nabla f(x), y-x \rangle}_{\text{blue}} + \underbrace{\frac{1}{2} \|y-x\|^2}_{\text{blue}} \right) \\
 &= f(x) + \beta \left(\frac{1}{2} \left(\|y-x\| + \frac{1}{\beta} \|\nabla f(x)\| \right)^2 - \frac{1}{2\beta} \|\nabla f(x)\|^2 \right)
 \end{aligned}$$

$$\begin{aligned}
 \text{Check: } \beta &\left(\frac{1}{2} \|y-x\|^2 + \langle y-x, \frac{1}{\beta} \nabla f(x) \rangle \right. \\
 &\quad \left. + \frac{1}{2\beta} \|\nabla f(x)\|^2 \right) \\
 &= \frac{\beta}{2} \|y-x\|^2 + \langle y-x, \nabla f(x) \rangle \\
 &\quad + \frac{1}{2\beta} \|\nabla f(x)\|^2
 \end{aligned}$$

$$m_x(y) = f(x) + \beta \left(\frac{1}{2} \left\| y - \underbrace{\left(x - \frac{1}{\beta} \nabla f(x) \right)}_{x^+} \right\|^2 \right)$$

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

$m_x(y)$

minimizer: $x - \frac{1}{\beta} \nabla f(x)$

Let $x^+ = x - \frac{1}{\beta} \nabla f(x)$

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

gradient descent.

Some nice observations:

$$f(x_+) \leq m_x(x_+) = f(x) - \frac{1}{2\beta} \|\nabla f(x)\|^2$$

global upper bound

* We have $f(x^+)$ lower

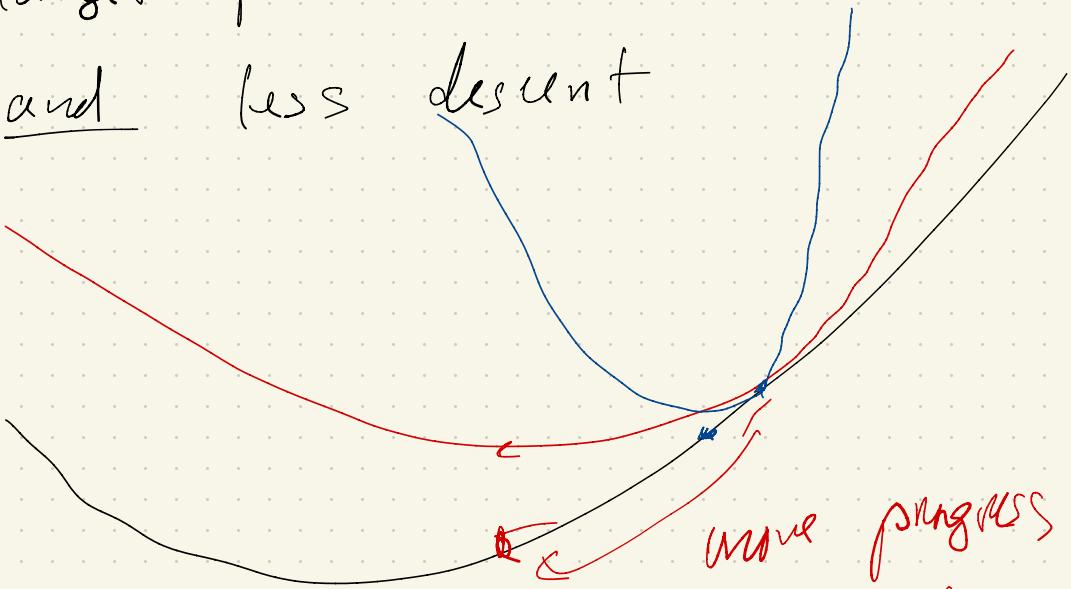
than $f(x)$ unless $\nabla f(x) = 0$,

in which case we are at a stationary point.

Step in our algorithm is $\frac{f}{\beta}$

larger $\beta \Rightarrow$ smaller step.

and less descent



more progress
if I get tighter
bounds,

$$\text{We have } f(x_{k+1}) \leq f(x_k) - \frac{1}{2\beta} \|\nabla f(x_k)\|^2$$

★

$$\text{for } x_{k+1} = x_k - \frac{1}{\beta} \nabla f(x_k)$$

Rearrange ★, we end up with

$$\frac{1}{2\beta} \|\nabla f(x_k)\|^2 \leq f(x_k) - f(x_{k+1})$$

+

$$\frac{1}{2\beta} \|\nabla f(x_{k-1})\|^2 \leq f(x_{k-1}) - f(x_k)$$

telescoping

$$+ \dots \leq f(x_1) - f(x_2)$$

$$\frac{1}{2\beta} \|\nabla f(x_0)\|^2 \leq f(x_0) - f(x_1)$$

$$\frac{1}{2\beta} \sum_{j=0}^k \|\nabla f(x_j)\|^2 \leq f(x_0) - f(x_{k+1})$$

$f(x_0) - f(x_{k+1})$

$$\underbrace{\frac{1}{K} \sum_{j=0}^K \| \nabla f(x_j) \|^2}_{\text{average value of observed}} \leq \underbrace{2\beta}_{K} (f(x_0) - f^*)$$

average value of observed
quantities,

$$\min_{j=0, \dots, K} \| \nabla f(x_j) \|^2 \leq \underbrace{2\beta}_{K} (f(x_0) - f^*)$$



$$\min \leq \text{avg}.$$

Result: $\min_{j=0, \dots, K} \| \nabla f(x_j) \|^2$

$$\leq \underbrace{\frac{1}{\sqrt{K}}}_{\text{Sublinear rate}} \underbrace{\int 2\beta (f(x_0) - f^*)}_{O(\sqrt{n})}$$

Sublinear rate, $O(\sqrt{n})$

Our first analysis!

These results get stronger as we add more assumptions.

- ① We can talk about progress in terms of fun values, distance to unique soln'
- ② Rates get better

Thm 2.25: If f is β -smooth and convex, x^* any minimizer we have

$$\min_{j=0, \dots, k} \|\nabla f(x_j)\| \leq \frac{2\beta \|x_0 - x^*\|}{K}$$

not $\sum K$

$$\underbrace{f(x_k) - f^*}_{\downarrow} \leq \frac{\beta \|x_0 - x^*\|^2}{2k}$$

fun values, when f convex.

Thm 2.26: If α -convex + β -smooth

$$\|x_k - x^*\|^2 \leq (\gamma)^k \|x_0 - x^*\|^2$$

↓ Linear rate

