

L-Smooth: No Convexity Needed

2024年10月22日 星期二 15:43

Can we do something better

We made assumptions that

1. f is L -Lip
2. $\|x_0 - x^*\| \leq R$
3. f convex, diff



No real Hessian, The ball doesn't pick up momentum

$$\Rightarrow f\left(\frac{1}{T} \sum_{s=0}^{T-1} x_s\right) - f(x^*) \leq \frac{RL}{\mu T} \quad \text{where } \mu = \frac{R}{LJT}$$

L-Smooth Definition

Def: a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth (stronger than L -Lip)

if gradient is L -Lip, bounding the gradient instead of bounding function

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

Much nicer statement is given now

1. Gradient L -Lip
2. Hessian cannot be crazy

L-Smooth Bounds on Gradient

Thm: If f is L -smooth and twice differentiable (Hessian exist) then $V^T \nabla^2 f(x) V \leq L$, $\forall x \in \mathcal{R}, V \in \mathbb{R}^n, \|V\|=1$ (upperbounded by L)

$$0 \leq V^T \nabla^2 f(x) V \leq L$$

Min and max both bounded



$$\|\nabla^2 f(x)\|_2$$

matrix norm: How much it can be stretched

L-Smooth gives Stronger Convergence

Thm: If $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth and convex, and $0 \leq \mu \leq \frac{1}{L}$ (bound μ)

This binding is saying that when hessian big, L is big, take smaller upward steps

Thm GD satisfies $f(x^{(t)}) - f(x^*) \leq \frac{1}{2t\mu} \|x^{(0)} - x^*\|^2$

Remark: $\|x^{(0)} - x^*\| \leq R, \mu = \frac{1}{L} \Rightarrow f(x^{(t)}) - f(x^*) \leq \frac{RL}{2t}$ Much stronger change at each step

Example

$$\|x^{(0)} - x^*\| \leq 10$$

$$L=2, \mu=1/2 \Rightarrow t \geq 10,000$$

$$f(x^{(t)}) - f(x^*) \leq \frac{1}{1000} = 0.001 \text{ way better bound}$$

Remark: Didn't assume twice differentiable

Pf strategy: $f(y) \leq f(x) + \nabla f^T(x)(y-x) + \frac{1}{2}(y-x)^T(LI)(y-x)$

This shows that even once differentiable, we can prove this theorem

$$\frac{1}{2}\|y-x\|^2$$

How do we converge

Does gradient blow up, very small, blows up, ... How smooth does the convergence come

Thm: Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$, L -smooth, For any $0 \leq \mu \leq \frac{1}{L}$, each step of GD gives

$$f(x^{(t+1)}) \leq f(x^{(t)}) - \frac{\mu}{2} \|\nabla f(x^{(t)})\|^2$$

It definitely goes down, at least at the previous gradient or smaller (converges)

Proof

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

$$\text{GD tells: } x^{(t+1)} = x^{(t)} - \mu \nabla f(x^{(t)}) \rightarrow x^{(t+1)} - x^{(t)} = -\mu \nabla f(x^{(t)})$$

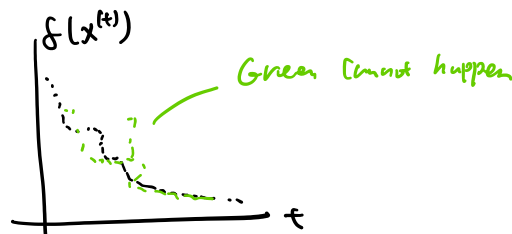
$$\Rightarrow f(x^{(t+1)}) \leq f(x^{(t)}) - \mu \nabla f(x^{(t)})^T \nabla f(x^{(t)}) + \frac{L}{2} \mu^2 \|\nabla f(x^{(t)})\|^2$$

$$\Rightarrow f(x^{(t+1)}) \leq f(x^{(t)}) - \mu \left(1 - \frac{L}{2}\mu\right) \|\nabla f(x^{(t)})\|^2$$

$$\text{When } \mu \leq \frac{1}{L} \Rightarrow \geq \frac{1}{2}$$

$$\Rightarrow f(x^{(t+1)}) \leq f(x^{(t)}) - \frac{\mu}{2} \|\nabla f(x^{(t)})\|^2 \quad \square$$

If we think about



How Big is the number gonna be

Thm: $f: \mathbb{R}^n \rightarrow \mathbb{R}$ L -smooth, $0 \leq \mu \leq \frac{1}{L}$

Then GD for T iterations means at least at one x_t must satisfy

$$\|\nabla f(x^{(t)})\| \leq \sqrt{\frac{2[f(x^{(0)}) - f(x^*)]}{\mu T}}$$

At some point, we must reach a small gradient, gradient always get smaller

Proof

Assume f is lower bounded (not $-\infty$ $f(x^*)$)

$$f(x^{(T)}) - f(x^{(0)}) = \sum_{t=0}^{T-1} [f(x^{(t+1)}) - f(x^{(t)})] \quad \text{Telescoping Theorem}$$

$$-\frac{\mu}{2} \|\nabla f(x^{(t)})\|^2 \quad \text{proved earlier}$$

$$f(x^{(T)}) - f(x^{(0)}) \leq -\sum_{t=0}^{T-1} \frac{\mu}{2} \|\nabla f(x^{(t)})\|^2$$

$$\sum_{t=0}^{T-1} \|\nabla f(x^{(t)})\|^2 \leq \frac{2}{\mu} (f(x^{(0)}) - f(x^{(T)}))$$

Since $f(x^{(T)}) \geq f(x^*)$, so can change to $f(x^*)$ maintains the relationship

$$\sum_{t=0}^{T-1} \|\nabla f(x^{(t)})\|^2 \leq \frac{2}{\mu} (f(x^{(0)}) - f(x^*))$$

We say that there exists at least one $x^{(t)}$ such that

$$\|\nabla f(x^{(t)})\|^2 \leq \frac{2}{\mu T} (f(x^{(0)}) - f(x^*))$$

$\sum_{t=0}^{T-1} \|\nabla f(x^{(t)})\|^2$ is a converging series and for a

converging series, averaging T terms, one of them must be smaller than or equal to T

$$\Rightarrow \|\nabla f(x^{(t)})\| \leq \sqrt{\frac{2(f(x^{(0)}) - f(x^*))}{\mu T}} \quad \square$$

NO Convexity Required Remarks

Last 2 theorems of how we converge and how big we converge

Never assumed convexity!!! Things we are about convex function will transfer to non-convex

Just can't say global minimum for sure, but it will work, function value will go down, and it will stop!!!