

Convex functions, part 3: Lipschitz gradient, etc.

Sunday, January 31, 2021 3:49 PM

mostly details not in BV'04

Recall a function $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is **Lipschitz continuous** if $\exists L > 0$ st. $\forall x, y \in \text{dom}(F)$,
 $\|F(x) - F(y)\| \leq L \cdot \|x - y\|$.

If F' exists, then $\|F'\| \leq L \Rightarrow F$ is Lipschitz continuous
 appropriate operator norm, usually 1-norm if 1D or spectral norm

What do we mean by this?

The Jacobian, where if $F(x) = \begin{bmatrix} F_1(x) \\ \vdots \\ F_m(x) \end{bmatrix}$, $F_i: \mathbb{R}^n \rightarrow \mathbb{R}$,

$$J_{ij} = \frac{\partial F_i(x)}{\partial x_j} \quad \left(\text{or } \frac{\partial F_j(x)}{\partial x_i} \right), \quad \text{I can never remember, and conventions aren't consistent anyhow}$$

$m \times n$ (or $n \times m$) matrix.

In optimization, "Jacobian" is often confusing, since it's unclear what "F" is.

Ex: $f: \mathbb{R}^n \rightarrow \mathbb{R}$, $f(x)$ a scalar

$\nabla f: \mathbb{R}^n \rightarrow \mathbb{R}^n$, "gradient", $\nabla f(x)$ a vector, "operates" on directions d
 like $\langle \nabla f(x), d \rangle$ (so a linear operator)

$\nabla^2 f: \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$, "Hessian", $\nabla^2 f(x)$ a matrix, "operates" on d
 (symmetric matrix) like $\langle d, \nabla^2 f(x) d \rangle$

$$(\nabla f(x))_i = \frac{\partial f}{\partial x_i}$$

$$(\nabla^2 f(x))_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

... so Jacobian of $F=f$ is the gradient (transposed)

... but Jacobian of $F=\nabla f$ is the Hessian

Fact Suppose $f \in C^2(U)$ for some open set $U \subseteq \mathbb{R}^n$, then

$(\nabla f \text{ is } L\text{-Lipschitz cts. on } U)$ iff $(\forall x \in U, \nabla^2 f(x) \leq L \cdot I)$

i.e., all eigenvalues $(\nabla^2 f(x)) \leq L$

$\Rightarrow \|\nabla^2 f(x)\| \leq L$

Fact Suppose... (same as above), then

$(\nabla f \text{ is } \mu\text{-strongly convex on } U)$ iff $(\forall x \in U, \mu I \leq \nabla^2 f(x))$
 (w/ respect to $\|\cdot\|_2$)

(need $\mu > 0$. If $\mu = 0$ this is plain old convexity)

So... one of our common assumptions will be $\forall f$ is L Lipschitz ($\nabla^2 f \leq L\mathbf{I}$)

And, a bit less often, also assume strong convexity ($\nabla^2 f \geq \mu\mathbf{I}$)

Q1 Is $f(x) = e^{-x}$ convex, strictly convex, strongly convex?
Is f' Lipschitz?

Q2 Is $f(x) = \begin{cases} -\log(x) & x > 0 \\ +\infty & x \leq 0 \end{cases}$ convex....?
 f' Lipschitz?

A1 e^{-x} on \mathbb{R} is strictly convex (hence convex) but not strongly conv
($f'(x) = e^{-x} \rightarrow 0$ as $x \rightarrow \infty$). It is strongly convex on the domain $(-\infty, R]$ for any $R < \infty$.



Similarly, f' isn't Lipschitz on \mathbb{R} but it is Lipschitz on $[-R, \infty)$ $\forall R < \infty$

A2 $-\log(x)$ is convex and strictly convex



Since $f' = \frac{-1}{x}$ is monotone
... but not strongly convex unless we look at $(-\infty, R]$ again.

For f' Lipschitz, it isn't on $(0, \infty)$ but it is on $[\delta, \infty)$ $\forall \delta > 0$.
} Can create problems for some algorithms if they converge to $x=0$

Ex $f(x) = \frac{1}{2} \|x\|_2^2$, $\nabla f(x) = x$, $\nabla^2 f(x) = \mathbf{I}$
 $\Rightarrow L=1$, $\mu=1$. Only function w/ this property.

NICEST FUNCTION EVER:

Calculus: $f(x) = e^x$

Statistics: $f(x) = e^{-x^2/2}$ or multivar. version

Optimization: $f(x) = \frac{1}{2} x^2$ (= negative log-likelihood of Gaussian!)

Def The condition # of f is $K_f = \frac{L}{\mu}$.
 $K_f \approx 1$ good
 $K_f \approx \infty$ bad

Why these assumptions?

Taylor's Thm: $f(y) = f(x) + f'(x) \cdot (y-x) + \frac{1}{2} f''(\xi) (y-x)^2$
 for some $\xi \in [x, y]$ (or in $[y, x]$). Similar in higher-dim.

if $f''(\xi) \leq L \quad \forall \xi$, then

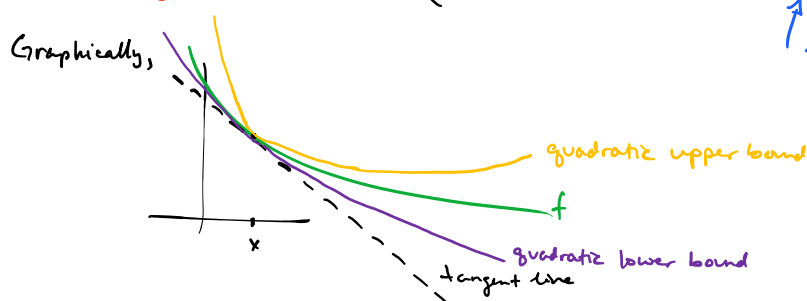
$$f(y) = \dots \leq f(x) + f'(x)(y-x) + \frac{1}{2} L (y-x)^2.$$

So...

Theorem If ∇f is L -Lipschitz and f is μ -strongly convex (and take $\mu=0$ if just convex)

then $\forall x, y \in \text{dom}(f)$,

$$\mu/2 \|y-x\|^2 \leq f(y) - \left(f(x) + \langle \nabla f(x), y-x \rangle \right) \leq L/2 \|y-x\|^2$$



This inequality is sometimes called "The Descent Lemma"

Usually f is complicated, but now we can "sandwich" it between

a quadratic upper bound and a quadratic lower bound (if strongly conv, $\mu > 0$)
 or a linear lower bound (if just conv, $\mu = 0$)

and quadratics are easy to work with, eg. easy to minimize in closed form, etc.

More properties

eg. f convex $\Rightarrow \nabla f$ monotone, meaning $\langle x-y, \nabla f(x) - \nabla f(y) \rangle \geq 0$

These can be strengthened w/ our $\mu > 0$ and L assumptions

See github class website

Handouts / Strong Convexity Lipschitz.pdf