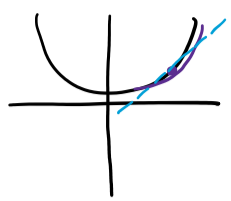# Strongly Convex

## Strong Convexity

A function $f: \mathbb{R}^n \to \mathbb{R}$ is **strongly convex** if $\exists \, c > 0$ s.t. $x, y \in \mathbb{R}^n$

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{c}{2} \| y - x \|^2$$



Instead of **tangent line**, we actually have a **tangent curve**

(arrow) This derive into this

Size of the gradient tells you how close you are to the $x^*$ when your condition is **under strongly convex**, and the relation is **squared**

$$\Rightarrow \quad f(x) - f(x^*) \leq \frac{\| \nabla f(x) \|^2}{2c} \quad \text{(proved in hw)}$$

← This is a property of the function, it is given to you

## Exponential Convergence Rate

We will show that strong convexity converges super quick.

1. GD: $x^{(t+1)} = x^{(t)} - \mu_t \nabla f(x^{(t)})$
2. Assuming $L$-smooth, bounds gradient $( \| \nabla^2 f(x) \| \leq L )$

$$f(x^{(t+1)}) \leq f(x^{(t)} - \underbrace{\mu_t f(x^{(t)})}_{h} )$$

$\leq L \| h \|^2$

Taylor theory → $f(x^{(t)}) - h^T \nabla f(x^{(t)}) + \frac{1}{2} h^T \nabla^2 f(z) h$

Cauchy-Schwarz: $\langle h, \nabla^2 f(z) h \rangle \leq \| h \| \, \| \nabla^2 f(z) h \|$
$\leq \| h \| \, \| \nabla^2 f(z) \| \, \| h \|$

↑ Biggest the hessian can be: $L$

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

$$\downarrow$$

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

Now nice property show up: $f(x^{(t+1)}) \leq f(x^{(t)}) - \mu_t \left( 1 - \frac{L\mu_t}{2} \right) \| \nabla f(x^{(t)}) \|^2$

Pick $\mu_t \leq \frac{1}{L}$ ($L$-smooth)

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

2 Directions now

### A: No Assumption of strong Convexity

Telescoping Series: $\frac{\mu_t}{2} \sum_{t=0}^{T-1} \| \nabla f(x^{(t)}) \|^2 \leq \sum_{t=0}^{T-1} \left[ f(x^{(t)}) - f(x^{(t+1)}) \right]$

As $T \to \infty$, $\| \nabla f(x^{(t)}) \|^2$ must $\to 0$
Done this before

$= f(x^{(0)}) - f(x^T)$

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

### B: Assuming strong Convexity

$$\| \nabla f(z) \|^2 \geq 2c \left[ f(z) - f(x^*) \right]$$

Then $f(x^{(t+1)}) \leq f(x^{(t)}) - \frac{\mu_t}{2} \left[ 2c \left[ f(x^{(t)}) - f(x^*) \right] \right]$   Do $- f(x^*)$ on both side

$f(x^{(t+1)}) - f(x^*) \leq (1 - \mu_t c) \left( f(x^{(t)}) - f(x^*) \right)$

$\downarrow$

Assume $\mu_t \leq \frac{1}{L}$  $\leq (1 - \frac{c}{L}) (f(x^{(t)}) - f(x^*))$

$f(x^{(t+1)}) - f(x^*) \leq (1 - \frac{c}{L}) (f(x^{(t)}) - f(x^*))$

$\leq (1 - \frac{c}{L})^t (f(x^{(0)}) - f(x^*))$  Chase all the way back to $x^{(0)}$

This is true with GD + ① Strong Convexity
② $L$-Smooth $\| \nabla^2 f(x) \| \leq L$
③ $\mu_t = 1/L$

This is a really really fast exponential convergence rate

## Generalization by Squeeze

When our function is quadratic $f(x) = \frac{1}{2} x^T A x$

$$\frac{c}{L} = \frac{1}{K} \leftarrow \text{Condition \#}$$

↑
Largest Eigen value of $\nabla^2 f(x)$ (A), can't stretch it more than $L$
direction of max convexity

Smallest Eigenvalue of $A$, direction of least convexity

If $\frac{c}{L} = 1$, It's saying GD of the quadratic with identity matrix.
Closer $c$ to $L$, easier

Condition # for general strongly convex function:

### 1. $L$-Smooth

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

$f(x+h) = f(x) + \nabla f(x)^T h + \frac{1}{2} h^T \nabla^2 f(z) h$
$\leq f(x) + \nabla f(x)^T h + \frac{L}{2} \| h \|^2$

$\rightarrow Q_L(h)$

Can't be "too convex"
Bounded above, can't be too bended

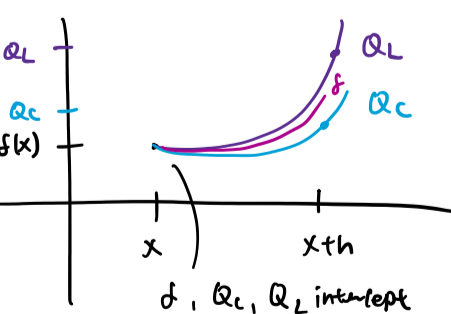Construct Linear Quadratic in $h$
(This is similar notion to Condition # $K$)

### 2. Strongly Convex Part

$f(x+h) \geq f(x) + \nabla f(x)^T h + \frac{c}{2} \| h \|^2$

$Q_c(h)$

Must be "some what Convex"
Can't be too flat



function is stuck between $Q_c$ and $Q_L$
→ Squeeze theorem

Note: $f(x)$ doesn't need to Quadratic any more, as long as it satisfy ① and ②, it doesn't matter what $f$ is, we have auto Quadratic bounding

At this sweet spot, GD Converges exponentially fast

## Under Strongly Convex, previous Method all generalizes well

When $\frac{c}{L} \equiv K$   GD w/momentum have Conv rate of $\frac{\sqrt{K} - 1}{\sqrt{K} + 1}$
for Strong Convexity and $\| \nabla^2 f(x) \| \leq L$

GD w/ N.A. have Convergence rate of $\sqrt{\frac{\sqrt{K} - 1}{\sqrt{K}}}$

Turns out the things we talked about Quadratic function is true for function with Strongly Convexity
When removed Strongly Convexity, momentum would fail, but N.A. would still work

What Ridge regression is doing is to artificially create a $c$ and make a function strongly Convex