

In multivariate case we similarly use the Hessian to understand the curvature

$\nabla^2 f(x)$ is a square matrix with eigenvalues λ_i at the point x .

- ① if $\lambda_i \geq 0$ for all $i \Rightarrow$ locally convex (local minima) (strong if strict)
- ② if $\lambda_i \leq 0$ for all $i \Rightarrow$ locally concave (local maxima) (strong if strict)

otherwise we will have some sort of saddle point.

$$\mu I = \begin{bmatrix} \mu & & 0 \\ & \mu & \\ 0 & & \mu \end{bmatrix}$$

$\hookrightarrow \lambda_i = \mu$

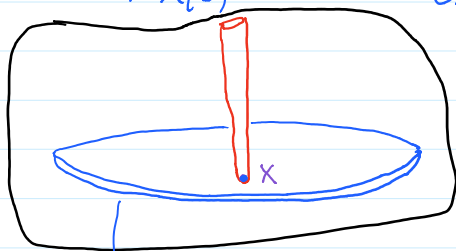
$$L I = \begin{bmatrix} L & & 0 \\ & L & \\ 0 & & L \end{bmatrix}$$

$\hookrightarrow \lambda_i = L$

$$\mu I \preceq \nabla^2 f(x) \preceq L I$$

\hookrightarrow Comparison of the eigenvalues.

if true at all x , then we are certain our function does not have a flat surface ($\lambda_i \geq \mu > 0$) hence strongly convex. When upper bounded by L , it means it can't change too quickly hence strongly smooth.



an alternative definition is to say $\mu \|v\|^2 \leq v^T \nabla^2 f(x) v \leq L \|v\|^2$ for all v !
equivalent ($\lambda_i = \operatorname{argmin} \frac{v^T \nabla^2 f(x) v}{\|v\|^2}$)

\hookrightarrow roughly saying that the function will be above the blue hemisphere (μ) and below the red hemisphere (L) around x .

We can approximate a function f around the point x using the derivatives of the function at x .

Let's say $x, f(x), \frac{df(x)}{dx}$, and $\frac{d^2f(x)}{dx^2}$ are given, we then will have

$$f(y) \approx g(y) = \underbrace{f(x)}_{\text{Const}} + \underbrace{\frac{df(x)}{dx}}_{\text{Const}} (y-x) + \underbrace{\frac{d^2f(x)}{dx^2}}_{\text{Const}} (y-x)^2 \frac{1}{2}$$

$$f(x) \approx f(x) + \underbrace{\dots}_{\text{Const}} dx + \underbrace{\dots}_{\text{Const}} dx^2 + \dots$$

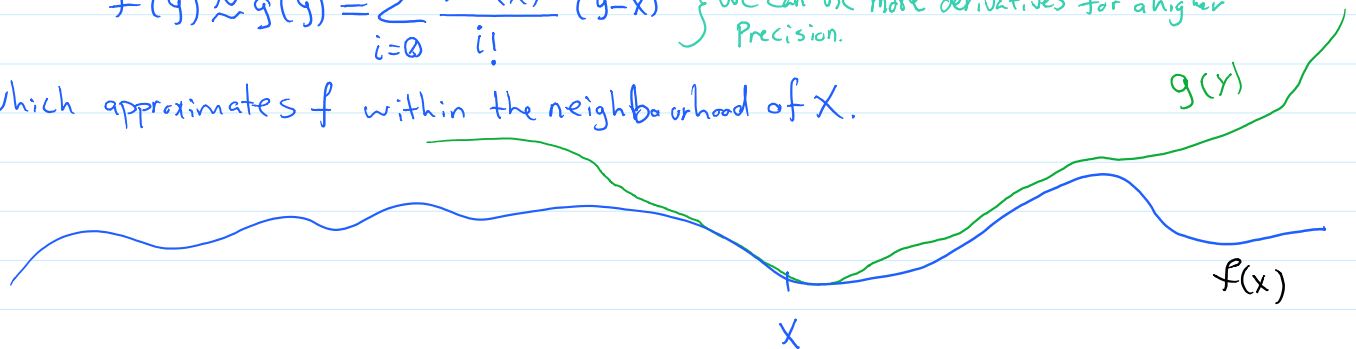
↳ g is a function of y !

more generally a polynomial approximation $g(y)$ can be written as (Taylor expansion)

$$f(y) \approx g(y) = \sum_{i=0}^{\infty} \frac{f^{(i)}(x)}{i!} (y-x)^i$$

We can use more derivatives for a higher precision.

Which approximates f within the neighbourhood of x .



We use a generalization of this in the multivariate case

$$g(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(x) (y-x) + \underbrace{O(\|y-x\|^2)}_{\text{The discarded terms can be bounded by this.}}$$

or this particular variant that is applicable to our problem (convex optimization)

$$\exists z: f(y) = g(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x)$$

To prove the convergence rate of a gradient step $x^{t+1} = x^t - \frac{1}{L} \nabla f(x^t)$ we need to

- ① Guarantee progress
- ② Guarantee termination

↳ not enough by itself because I can make progress, but small enough that I only converge to a fixed point, e.g., $\sum_{i=0}^{\infty} \frac{1}{2^i} = 2$

↳ as a measure of $f(x^t) - f(x^*)$?
 $\|x^t - x^*\|$?

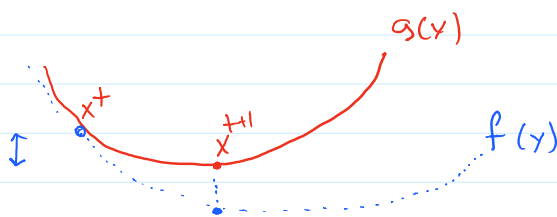
We assume $\mu I \preceq \nabla^2 f(x) \preceq L I$ for all x .

① Progress. We start with the Taylor expansion

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x) \quad \text{for some } z.$$

$< L \|y-x\|^2$

$$\Rightarrow f(y) \leq \underbrace{f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|^2}_{g(y)}$$



g is a quadratic upper bound on f!

Since g is always above f , it is a safe bet if I minimize g instead of f !

Set $\nabla g(y) = 0$ to find the minimizer of g . note that $\nabla^2 g(y) > 0 \Rightarrow$ Convex

$$\Rightarrow \nabla g(y) = \nabla f(x) + L(y-x) = 0 \Rightarrow \boxed{y = x - \frac{1}{L} \nabla f(x)}$$

if I pick the next point x^{++} , I will be jumping to the minimizer of g , and we know the real function will be below g , so we're safe. So, how much progress have we made? plugging the value of y back, we'd get

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

$$\Rightarrow f(x) - \frac{1}{L} \|\nabla f(x)\|^2 + \frac{1}{2L} \|\nabla f(x)\|^2 = f(x) - \frac{1}{2L} \|\nabla f(x)\|^2$$

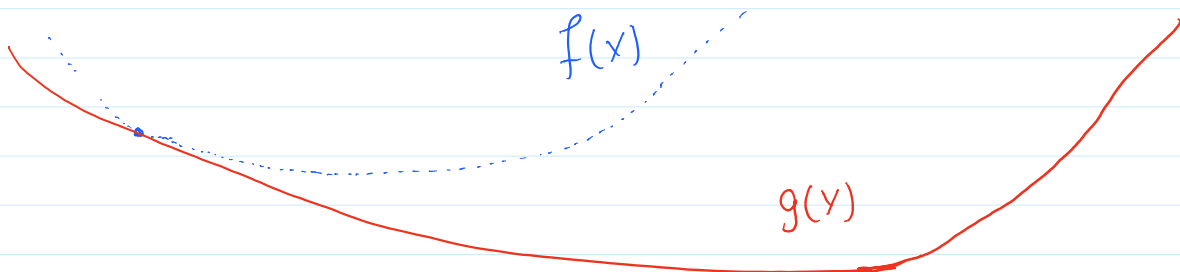
So we have decreased $f(x^*)$ by jumping to $f(x^{++})$ by at least $\frac{1}{2L} \|\nabla f(x^*)\|^2$

as long as $\|\nabla f(x^*)\|$ is not zero, progress is guaranteed!

now from the other side we have

$$f(x) \geq \underbrace{f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2}_{g(y)}$$

g is always below f



Now we have a quadratic lower bound on f , which is useful to bound the maximum distance to the real minimizer because we can minimize $g(x)$ exactly.

if we minimize $g(x)$ we'd get $f(x) \geq \underbrace{f(x^*)}_{\text{function}} - \underbrace{\frac{1}{2\mu} \|\nabla f(x)\|^2}_{\text{a fixed value}} *$

if $f(x)$ is always above $*$, then $f(x^*)$ is also above $*$

$$\Rightarrow f(x^*) \geq f(x^t) - \frac{1}{2\mu} \|\nabla f(x^t)\|^2 \quad \textcircled{\text{I}}$$

we previously had $f(x^{t+1}) \leq f(x^t) - \frac{1}{2L} \|\nabla f(x^t)\|^2 \quad \textcircled{\text{II}}$

and we wish to show $f(x^*) - f(x^t)$ goes to 0 as t goes to ∞ .

$$\textcircled{\text{I}} \Rightarrow \mu (f(x^*) - f(x^t)) \geq -\frac{1}{2} \|\nabla f(x^t)\|^2 + \textcircled{\text{II}}$$

$$\hookrightarrow f(x^{t+1}) \leq f(x^t) + \frac{\mu}{L} (f(x^*) - f(x^t)) \quad \text{subtract } f(x^*)$$

$$f(x^{t+1}) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right) (f(x^t) - f(x^*)) \quad \textcircled{\text{III}}$$

\Rightarrow This means we get closer to $f(x^*)$ by $\rho = \left(1 - \frac{\mu}{L}\right)$ with each step.

by applying $\textcircled{\text{III}}$ to itself repeatedly, we'll have

$$f(x^t) - f(x^*) \leq \rho^t (f(x^0) - f(x^*)) \quad \rho < 1$$

which means our distance to the optimal objective is at most ρ^t times our initial distance.

if we wish to get at least ϵ close to the optimal objective, we can say:

$$f(x^t) - f(x^*) \leq \rho^t (f(x^0) - f(x^*)) \leq \epsilon$$

$$\hookrightarrow \rho^t \leq \epsilon \quad \Rightarrow \left(\frac{1}{\rho}\right)^t \leq \frac{1}{\epsilon}$$

$$\Rightarrow -t \log \frac{1}{\rho} \leq \log \epsilon \quad \Rightarrow t \geq \log \frac{1}{\epsilon} \cdot C \quad \Rightarrow t = o(\log \frac{1}{\epsilon})$$