# Newton's Method: Convergence Guarantees

## Assumptions

① Function $f \in C^2(\mathbb{R}^n)$ and $m$-strongly convex

$\Rightarrow$ on set $S = \{ x : f(x) \leq f(x^{(0)}) \}$

$$mI \preceq \nabla^2 f(x) \preceq MI \quad \forall \ x \in S$$

② We have $L$-Lipschitz Hessians

$$\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq L \|x-y\|_2$$

If $3^{rd}$ derivative exists $\Leftrightarrow$ Equivalent condition on $3^{rd}$ derivative that it is bounded.

⇓ Assump. 2 can be replaced by working with self-concordant functions (optional reading in §9.6).

e.g, $f(x) = x^T P x$ ; $P \succeq mI$

$$\nabla^2 f(x) = P$$

$\Rightarrow$ All quadratics Satisfy Assump. 2.

## Convergence behavior of Newton's Method

It has two phases of convergence:

① Phase I $\Rightarrow$ Damped phase $\Rightarrow$ It has linear convergence.

In this phase, backtracking provides a step size

that satisfies:

$$t \geq \min \left\{ \beta \frac{m}{M}, 1 \right\}$$

② <u>Phase II</u> $\Rightarrow$ Quadratic Convergent phase

$\Rightarrow$ Full Newton step phase $\Rightarrow t \approx 1 \; \forall \; l \geq k$
(some fixed $k$)

In this phase, we have convergence behavior is

$$f\left(x^{(l)}\right) - p^* = O\left(c^{2^l}\right) \quad ; \quad c = 0.5$$

## Summary of discussion

① Rapid (super linear / Quadratic) convergence eventually. Once in Quadratic convergent phase, we need only six to eight more iterations to reach optimal value.

② Newton's method is also not affected by change of coordinates

 ↳ It is much less sensitive to the condition number of a problem.

③ Performance scales well with the number of dimensions.

④ Back tracking parameters also do not affect the performance that much.

**<u>Drawback</u> :** Memory and Computation

**<u>Damped Phase (Linearly Convergent Phase)</u>**

Newton's method gives us linear convergence from $k=0$ upto some finite as long as

$$\| \nabla f(x^{(k)}) \|_2 \geq \eta \quad \text{for some}$$

$$0 < \eta < \frac{m^2}{L}$$

and $t \geq \min \left\{ \beta \frac{m}{M}, 1 \right\}$ for these iterations.

Basically, in each iteration, we will reduce the objective function by a constant $\gamma > 0$

$$\gamma = \alpha \beta \eta^2 \frac{m}{M^2}.$$

<u>Reminder.</u>  $\lambda(x) = \left( \nabla f(x)^T (\nabla^2 f(x))^{-1} \nabla f(x) \right)^{1/2}$

$\underset{\Downarrow}{^2}$

$\lambda(x)^2 = \Delta x_{nt}^T \nabla^2 f(x) \Delta x_{nt} \Longleftrightarrow \nabla f(x)^T \Delta x_{nt} = -\lambda(x)^2$

b/c $\Delta x_{nt} = -(\nabla^2 f(x))^{-1} \nabla f(x)$

<u>Quadratic upper bound</u>

$$f(x + t \Delta x_{nt}) \leq f(x) + t \underbrace{\nabla f(x)^T \Delta x_{nt}}_{-\lambda(x)^2} + \frac{M t^2}{2} \| \Delta x_{nt} \|_2^2$$

$$\leq f(x) - t \lambda(x)^2 + \frac{M t^2}{2} \| \Delta x_{nt} \|_2^2$$

b/c $\lambda^2(x) = \Delta x_{nt}^T \nabla^2 f(x) \Delta x_{nt} \geq m \| \Delta x_{nt} \|_2^2$

Remember: $\lambda_{min}(A) \leq \dfrac{V^T A V}{\|V\|_2^2} \leq \lambda_{max}(A)$ ; A is diagonalizable

Note: In class, we upperbounded $\lambda(x)^2$; we should have upperbounded $\|\Delta x_{nt}\|_2^2 \Rightarrow \|\Delta x_{nt}\|_2^2 \leq \dfrac{\lambda^2(x)}{m}$. Below are the original steps in class, that led to an extra $m$ factor. I am correcting them below, so you can tell what things were corrected in relation to class notes. Corrections will be in red.

$\Rightarrow f(x + t\Delta x_{nt}) \leq f(x) - t\,m\,\|\Delta x_{nt}\|_2^2 + \dfrac{m t^2}{2}\|\Delta x_{nt}\|_2^2$

$\quad\quad\quad$ $t\lambda(x)^2 \quad\quad\quad \dfrac{M}{2m}t^2\lambda(x)^2$

Pick $\boxed{\hat{t} = \dfrac{m}{M}}$

$\quad = f(x) - \dfrac{m}{M}\|\Delta x_{nt}\|_2^2 + \dfrac{M}{2m}\times\dfrac{m^2}{M^2}\cdot\|\Delta x_{nt}\|_2^2$

$\quad\quad\quad\quad\quad\quad \lambda(x)^2 \quad\quad\quad\quad\quad\quad \lambda(x)^2$

$\quad = f(x) - \left(\dfrac{m}{M} - \dfrac{m}{2M}\right)\|\Delta x_{nt}\|_2^2$

$\quad\quad\quad\quad\quad\quad\quad\quad \lambda(x)^2$

$\quad\quad\quad\quad\quad\quad \underbrace{\quad\quad}_{\frac{1}{2}\frac{m}{M} \to m}$

$f(x + t\Delta x_{nt}) \leq f(x) - \dfrac{1}{2}\cdot\dfrac{m}{M}\|\Delta x_{nt}\|_2^2 \quad \lambda(x)^2$

$f(x + t\Delta x_{nt}) \leq f(x) - \dfrac{1}{2}\dfrac{m}{M}\|\Delta x_{nt}\|_2^2\,\lambda(x)^2$

$\quad\quad\quad\quad$ If $\alpha \in \left(0, \dfrac{1}{2}\right]$

$\quad\quad\quad\quad \Rightarrow \leq f(x) - \alpha t\|\Delta x_{nt}\|_2^2\,\lambda(x)^2$

$\Rightarrow \hat{t} = \dfrac{m}{M}$ satisfies the back tracking condition

(remember, back tracking in Newton's method uses $f(x) - \alpha t\lambda(x)^2$

(remember, back tracking in
Newton's method uses $f(x) - \alpha t \lambda(x)^2$
condition)

we are bound to accept $t \geq \frac{\beta m}{M}$

$$f(x + t\Delta x_{nt}) \leq f(x) - \alpha t \underbrace{\|\Delta x_{nt}\|_2^2}_{\geq \alpha \beta \eta^2 \boxed{\frac{m}{M^2}}} \quad \lambda(x)^2$$

$\Rightarrow \alpha t \lambda(x)^2$
$\geq \alpha \beta \frac{m}{M} \times \lambda(x)^2$

Also;

$\lambda(x)^2 \geq \dfrac{1}{\lambda_{max}(\nabla^2 f(x))} \times \|\nabla f(x)\|_2^2 \geq \dfrac{\eta^2}{M}$

## Quadratic Convergence Phase

Once $\|\nabla f(x^{(u)})\|_2$ goes below $\eta$, then we

always have $t^{(u)} = 1$ and

$$\frac{L}{2m^2} \|\nabla f(x^{(u+1)})\|_2 \leq \left( \frac{L}{2m^2} \|\nabla f(x^{(u)})\|_2 \right)^2$$

## How many iterations for Newton's method?

Linear phase $\Rightarrow f(x)$ decreases by at least $\gamma$ in
each iteration

$$\Rightarrow \# \text{ of iterations} = \frac{f(x^{(0)}) - p^*}{\gamma}$$

Say we want final accuracy to be $\epsilon$

$$f(x^{(l)}) - p^* \leq \epsilon$$

$\forall \, l \geq k$

$\hookrightarrow = O(0.5^{2^l}) \leq \epsilon$

$\epsilon_0 = 2m^3/L^2$

$\log_2 \log_2 \left\lceil \frac{\epsilon_0}{\epsilon} \right\rceil = \#$ of iterations.

$$\log_2 \log_2 \left( \frac{\epsilon_0}{\epsilon} \right) \xrightarrow{\epsilon_0 = 2m^3/L^2} = \text{\# of iterations.}$$

$$\text{\# of iterations} \propto \log_2 \log_2 (\epsilon^{-\gamma})$$

$$\text{vs.} \quad GD \propto \log(\epsilon^{-1})$$

$$\nabla^2 f(x) =$$