

DSA3102: Lecture 12 (Readings Sec 9.4, 9.5 & my notes on Newton's method)

↳ very fast under the right conditions

Last time: An optimization algorithm produces a sequence $\{x^{(k)}\}_{k \in \mathbb{N}} \subset \mathbb{R}^n$ with the hope that $x^{(k)} \rightarrow x^*$
 (infinite collection of vectors) $x^* = \underset{x}{\operatorname{argmin}} f(x)$. If $f(x)$ strictly convex,
 x^* is unique

start at x^0
 Gradient descent: $x^{(k+1)} = x^{(k)} - t \nabla f(x^{(k)})$ move in direction of negative gradient
 Step size choice: $t = \underset{s \geq 0}{\operatorname{argmin}} f(x^{(k)} - s \nabla f(x^{(k)}))$

Backtracking line search.

only make use of first derivative

Conditions on f :

- $m < M$
- ① m -strongly convex: $\nabla^2 f(x) \geq mI \quad \forall x \in S$.
 - ② M -smoothness: $\nabla^2 f(x) \leq MI \quad \forall x \in S$.

Thm: GD with stepsize chosen according to exact line search on f that satisfies m -s.c. & M -smoothness yields.

linear convergence
 k th iteration $f(x^{(k)}) - p^* \leq \underbrace{(1 - \frac{m}{M})^k}_{\text{sub-optimality at the start}} (f(x^{(0)}) - p^*)$ geometric series
 $\|x^{(k)} - x^*\| \leq c \|x^{(k-1)} - x^*\| \xrightarrow[0.8]{10^{-6}} 0.8 \times 10^{-6} \rightarrow 0.64 \times 10^{-6}$ fast

Thm: GD with stepsize chosen according to exact line search on f that satisfies M -smoothness yields
 (no m strong convex)

cannot guarantee you converge
so quickly

$$f(x^{(k)}) - p^* = O(\gamma_k)$$

goes down
Polynomially fast

Snapshot of today main theoretical result:

$$\|x^{(k+1)} - x^*\| \leq c \|x^{(k)} - x^*\|^2.$$

quadratic
convergence

Recall Steepest Descent: motivation for Newton's method

$$f(x+v) \approx f(x) + \nabla f(x)^T v$$

Minimize $\nabla f(x)^T v$ over all directions v . (want to push f to be as small as possible)

Def: Normalized steepest descent direction.

$$\Delta x_{nsd} = \underset{\text{norm 1}}{\arg \min} \{ \nabla f(x)^T v : \|v\| = 1 \}$$

Def: Unnormalized steepest descent direction.

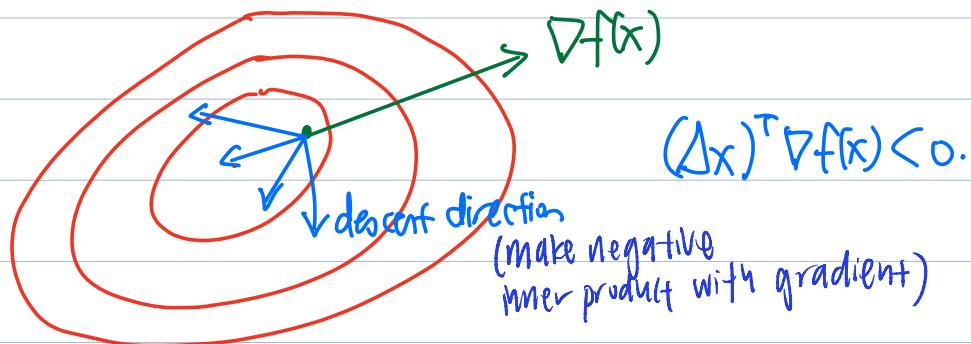
$$\Delta x_{sd} = \|\nabla f(x)\| \Delta x_{nsd}$$

Dual norm: $\|x\|_*: \sup \{ y^T x : \|y\| \leq 1 \}$
 $- \|x\|_* := \inf \{ y^T x : \|y\| \leq 1 \}$

Rmk: For the steepest descent step,

consider inner product between gradient and steepest descent direction

$$\begin{aligned}\nabla f(x)^T \Delta x_{\text{sd}} &= \nabla f(x)^T [\|\nabla f(x)\|_* \Delta x_{\text{nsd}}] \\ &= \|\nabla f(x)\|_* \nabla f(x)^T \Delta x_{\text{nsd}} \\ &= \|\nabla f(x)\|_* (-\|\nabla f(x)\|_*) \\ &= -\|\nabla f(x)\|_*^2 < 0\end{aligned}$$



The steepest descent direction is a descent direction.

$$-\|a\|_2 \|b\|_2 \leq a^T b \leq \|a\|_2 \|b\|_2 \quad \text{Cauchy}$$

Ex: Euclidean norm.

$$\Delta x_{\text{nsd}} = \arg \min \{ \nabla f(x)^T v : \|v\|_2 = 1 \}.$$

different norm get different algorithm

$$= -\nabla f(x) \quad \begin{matrix} \text{use lower bound} \\ \text{of Cauchy inequality} \end{matrix}$$

Rmk: If we instantiate steepest descent to the Euclidean norm,
we get gradient descent.

$$\frac{1}{p} + \frac{1}{q} = 1$$

Ex: ℓ_1 -norm

$$-\|a\|_1 \|b\|_\infty \leq a^T b \leq \|a\|_1 \|b\|_\infty$$

$$\Delta x_{\text{nsd}} = \arg \min \{ \underbrace{\nabla f(x)^T v}_{-\|\nabla f(x)\|_\infty} : \|v\|_1 = 1 \}.$$

differentiate wrt one coordinate

$$-\|\nabla f(x)\|_\infty$$

By Holder inequality, one vector constrained to l_1 norm, the other should be constrained by l_∞ norm

$$\|y\|_\infty = \max_{1 \leq i \leq n} |y_i| \quad \|y\|_1 = \sum_{i=1}^n |y_i|$$

partial derivative in i th coordinate

Let $i \in [n]$ be s.t. $|(\nabla f(x))_i| = \|\nabla f(x)\|_\infty$

$$\Delta x_{nsd} = -(\text{sign}) \left(\frac{\partial f(x)}{\partial x_i} \right) e_i.$$

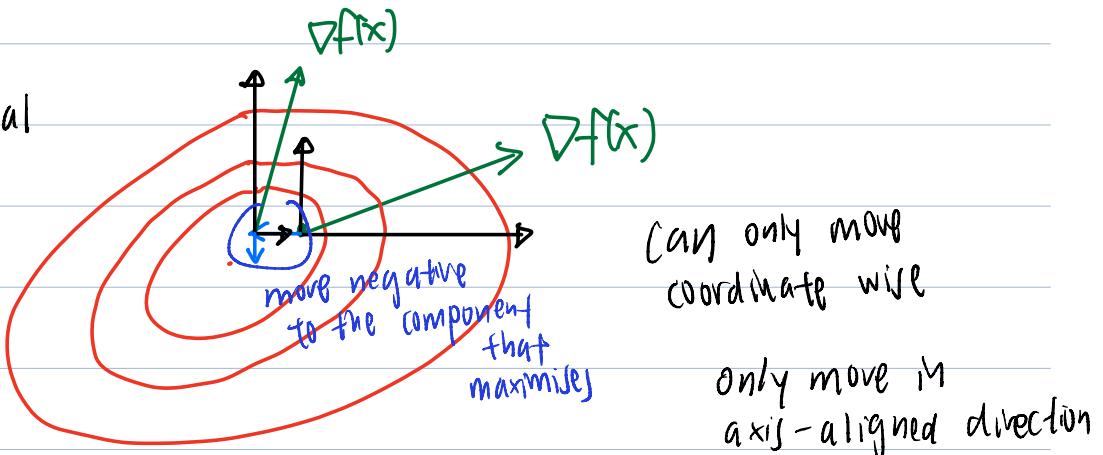
either up/down

look for coordinate which has maximum abs value

$$e_i = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 0 \end{pmatrix} \leftarrow \begin{array}{l} \text{i-th position} \\ (\text{move in only} \\ \text{1 particluar} \\ \text{direction}) \end{array}$$

Rmk: If we initialize steepest descent to the ℓ_1 norm, we get coordinate descent, i.e., move in the direction of maximal decrease.

find direction which yields maximal decrease



Newton's Method

hessian

Newton direction: $\Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x)$.

Descent direction: $(\Delta x_{nt})^\top \nabla f(x) = (-\nabla^2 f(x)^{-1} \nabla f(x))^\top \nabla f(x)$

take inner product of presumed desc direction with gradient and check that it is $-ve$

must be
positive
definite

$$= -\nabla f(x)^T \boxed{\nabla^2 f(x)^{-1}} \nabla f(x) < 0.$$

The Newton direction is a descent direction.

Rmk: We assume $\nabla^2 f(x)$ & $\nabla^2 f(x)^{-1}$ are positive definite. so that it is invertible

Ex: $f(x) = \frac{1}{2} x^T Q x$ $Q = \begin{bmatrix} 0.01 & 0 \\ 0 & 1 \end{bmatrix}$.

$$= \frac{1}{2} [0.01x_1^2 + x_2^2].$$

$$\underline{x^* = 0} \quad -\nabla f(x) = -Qx = \begin{bmatrix} 0.01x_1 \\ x_2 \end{bmatrix}.$$

If we start GD at $x^{(0)}$, the update is

$$x^{(1)} = x^{(0)} + t \begin{bmatrix} -0.01x_1^{(0)} \\ -x_2^{(0)} \end{bmatrix}.$$

Is this a good direction?

What is the best direction if we start at $x^{(0)} \neq 0$?

If is $-x^{(0)} = \begin{bmatrix} -x_1^{(0)} \\ -x_2^{(0)} \end{bmatrix}$ \leftarrow a direction

Compare the GD direction $\begin{bmatrix} -0.01x_1^{(0)} \\ -x_2^{(0)} \end{bmatrix}$ with the best direction $\begin{bmatrix} -x_1^{(0)} \\ -x_2^{(0)} \end{bmatrix}$. \rightarrow get to the optimum straight away

They can be very diff.
due to coefficient

Newton direction $\Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x)$ converge M
 $= -Q^{-1} Q x = -x.$ one single step

$$f(x) = \frac{1}{2} x^T Q x$$

$$\nabla f(x) = Q x$$

$$\nabla^2 f(x) = Q \text{ invertible}$$

\uparrow
gets best direction
(the minus of the point)

Newton direction: $\Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x)$

Steepest descent: $\Delta x_{sd} = \arg \min \left\{ \nabla f(x)^T v : \|v\|_{\infty} = 1 \right\}$

$$\Delta x_{sd} = -\|\nabla f(x)\|_* \Delta x_{sd}.$$

$$\Delta x_{nt} = \arg \min \left\{ \nabla f(x)^T v : \underbrace{\|v\|_{\nabla^2 f(x)}}_{\text{Hessian norm}} = 1 \right\}$$

steepest desc with hessian norm

$$\|v\|_P = \sqrt{v^T P v}$$

Verify that Δx_{sd} when instantiated to the Hessian norm gives Newton's direction.

$$C = \nabla f(x)$$

$$P = \nabla^2 f(x) > 0.$$

$$\min_v \left\{ C^T v : v^T P v = 1 \right\}, \quad \min_v \left\{ \nabla f(x)^T v : v^T \nabla^2 f(x) v = 1 \right\}$$

KKT condition: $L(v, \mu) = C^T v + \mu (v^T P v - 1)$

nodal
feasibility as no
inequality constraints

derivative
(wrt v)

Stationarity: $\nabla_v L(v, \mu) = c + 2\mu P v = 0$

$$v = -\frac{1}{2\mu} P^{-1} c.$$

Primal feasibility: $v^T P v = 1$

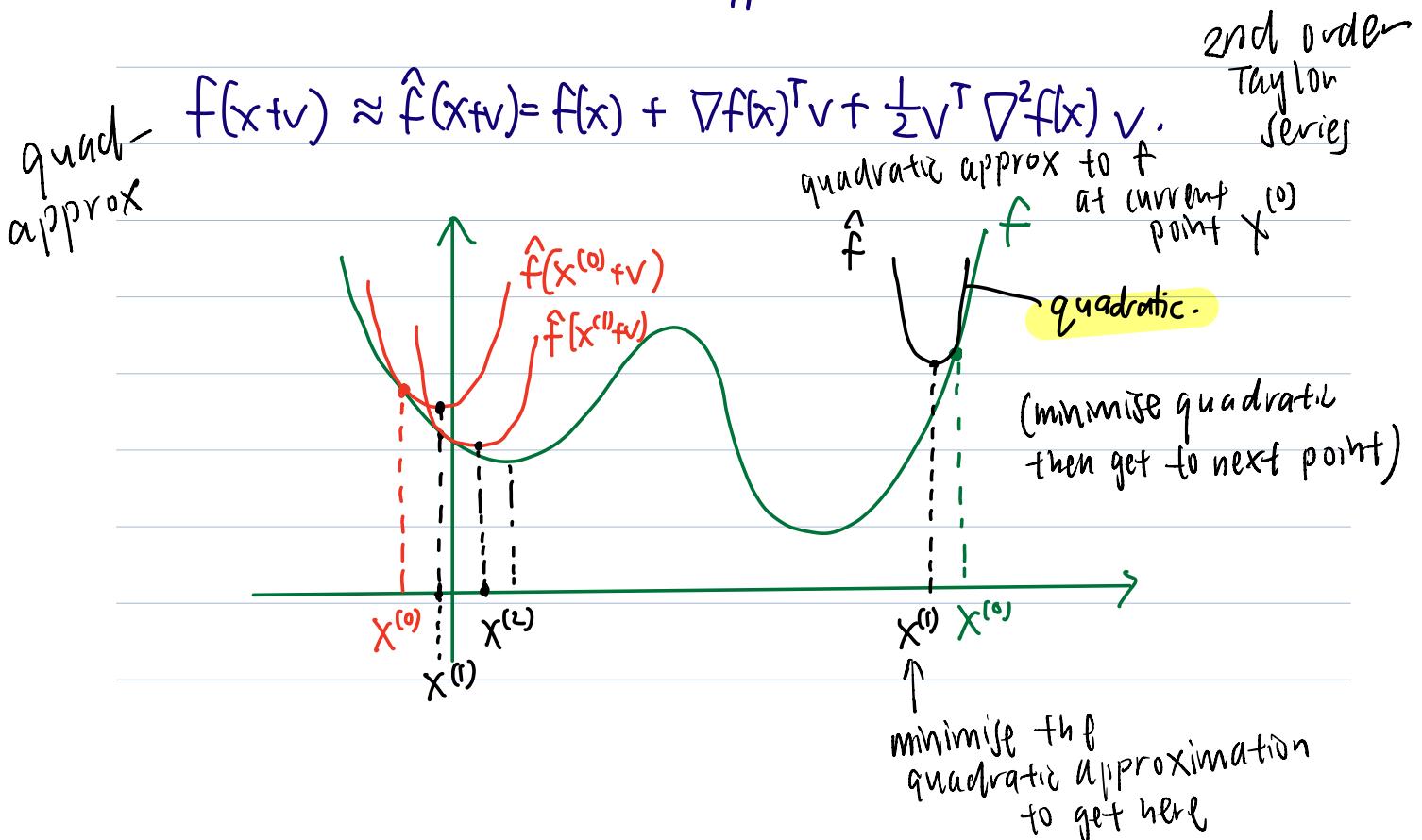
$$\left(-\frac{1}{2\mu} P^{-1} c\right)^T P \left(-\frac{1}{2\mu} P^{-1} c\right) = 1$$

$$\mu^2 = \frac{1}{4} c^T P^{-1} c \Rightarrow v = -\frac{1}{(c^T P^{-1} c)^{1/2}} P^{-1} c.$$

$$\Delta x_{nsd} = -((\nabla f(x))^T \nabla^2 f(x)^{-1} \nabla f(x))^{-1/2} \nabla^2 f(x)^{-1} \nabla f(x),$$

$$\Delta x_{sd} = \|\nabla f(x)\|_* \Delta x_{nsd} = -\nabla^2 f(x)^{-1} \nabla f(x)$$

Minimizer of the second-order approximation to f



Newton direction \rightarrow best direction according to 2nd order Taylor approx of f at current point

Differentiate $f(x) + \nabla f(x)^T v + \frac{1}{2} v^T \nabla^2 f(x) v$ w.r.t. v .

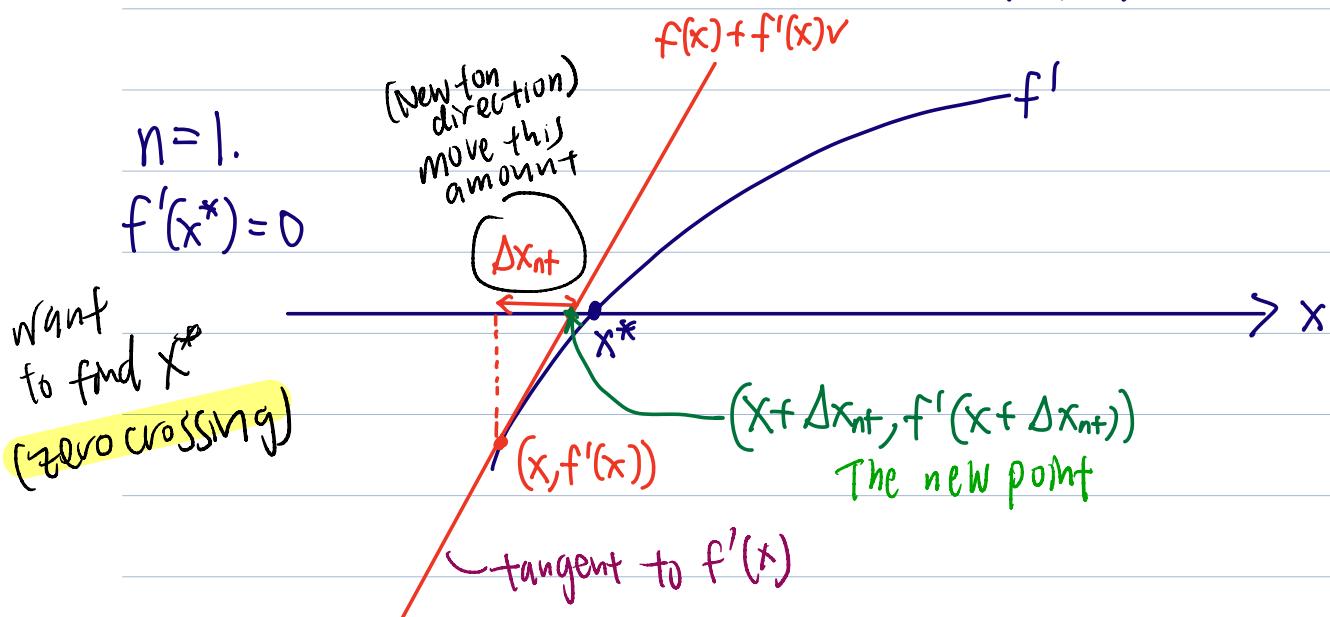
$$\text{direction to move } \nabla f(x) + \nabla^2 f(x) v = 0 \\ v = -\nabla^2 f(x)^{-1} \nabla f(x).$$

Solution of linearized optimality condition.

$$g = \nabla f \quad \text{gradient of } f \quad g(x^*) = 0 \quad \therefore \nabla f(x^*) = 0.$$

Linearizing this near x : $\nabla f(x) + \nabla^2 f(x) v = 0$

$$v = -\nabla^2 f(x)^{-1} \nabla f(x)$$



Affine invariance of Newton step:

Newton step is independent of linear or affine changes of the coordinates.

$T \in \mathbb{R}^{mn}$: invertible
non-singular.

T : matrix
 x, y : vector

$$g(y) = f(Ty)$$

$$x = Ty.$$

$$x, y \in \mathbb{R}^n$$

$$\nabla g(y) = T^T \nabla f(x), \quad \nabla^2 g(y) = T^T \nabla^2 f(x) T.$$

$$\begin{aligned}\Delta y_{nt} &= -\nabla^2 g(y)^{-1} \nabla g(y) \\ &= -(T^T \nabla^2 f(x) T)^{-1} T^T \nabla f(x) \\ &= -T^{-1} \nabla^2 f(x)^{-1} \cancel{T^{-1}} \cancel{T^T} \nabla f(x) \\ &= -T^{-1} \underbrace{\nabla^2 f(x)^{-1} \nabla f(x)}_{\text{Newton step in the}} \\ &= T^{-1} \Delta x_{nt} \quad \text{original coordinate system}\end{aligned}$$

\Rightarrow Newton steps of f and g are related by the same linear transformation.

$$x + \Delta x_{nt} = T(y + \Delta y_{nt})$$

Newton's Method:

Tolerance

Given a starting point $x^{(0)} \in \text{dom } f$, $\text{Tol } \varepsilon > 0$

loop

repeat

1. Compute $\Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x)$ newton direction computationally costly.

Newton decrement: $\lambda(x) = (\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x))^{\frac{1}{2}}$

↳ smaller this is,
the closer to the
optimum

Problem of Newton method:
inverting a matrix takes (step 1)
long time (solution: Approximate
the hessian (diagonal matrix))

(not
scalable) 2. Quit if $\lambda(x)^2/2 < \varepsilon$ $\|\nabla f(x)\|/\|\nabla^2 f(x)^{-1}\|$

3. Line search: Choose t by backtracking line search
or other methods.

4. $x^+ = x + t \underbrace{\Delta x_{nt}}_{\text{step size, newton direction}}$

Analysis of Newton's method near the optimum. \rightarrow converge very quickly
(pure Newton phase)

$$x^{(k+1)} = x^{(k)} - t \nabla^2 f(x^{(k)})^{-1} \nabla f(x)$$

Near the optimum, backtracking line search results in no reduction
of step size $\Rightarrow t=1$.
step size always 1

Can be seen as a solver for the system of n equations.

$$g(x) = 0 \quad g: \mathbb{R}^n \rightarrow \mathbb{R}^n$$

find x
such
that
derivative
vanishes
(find $x^{(k)} \rightarrow x^*$)

$$g = \nabla f$$

$$\nabla g = \nabla^2 f. \quad x^{(k+1)} = x^{(k)} - \nabla g(x^{(k)})^{-1} g(x^{(k)})$$

Intuitive argument to show why Newton converges very fast to x^* satisfying $g(x^*)=0$.

Taylor series approx

$$0 = g(x^*) = g(x^{(k)}) + \nabla g(x^{(k)})^\top (x^{(k)} - x^*) + o(\|x^{(k)} - x^*\|)$$

smaller terms

Multiply this with $\nabla g(x^{(k)})^{-1}$

$$x^{(k)} - x^* - \nabla g(x^{(k)})^{-1} g(x^{(k)}) = o(\|x^{(k)} - x^*\|).$$

$$x^{(k+1)} - x^* = o(\|x^{(k)} - x^*\|) \quad \text{small things}$$

the dist from next iterate to optimum divided by dist from current iterate to optimum is 0

$$\text{Compare this to GD} \quad \|x^{(k+1)} - x^*\| \leq 0.8 \|x^{(k)} - x^*\|$$

$$\limsup_{k \rightarrow \infty} \frac{\|x^{(k+1)} - x^*\|}{\|x^{(k)} - x^*\|} \leq 0.8 \quad \leftarrow \frac{M}{M}$$

⇒ linear convergence

Thm: $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$ ∇g is symmetric Jacobian matrix.
(not examinable)

For $\delta > 0$, let $S_\delta = \{x \in \mathbb{R}^n : \|x - x^*\| \leq \delta\}$.
↳ neighbourhood of optimum (very close to x^*)

Assume g is continuously differentiable in S_δ and that $\nabla g(x^*) \in \mathbb{R}^{n \times n}$ is invertible.

initialise very close to optimum

1. $\exists \delta > 0$ s.t. if $x^{(0)} \in S_\delta$, the sequence $\{x^{(k)}\}$ converges to x^* and the convergence of

$\{\|x^{(k)} - x^*\|\}$ is superlinear.

invertible at x^*
 ↳ in a neighbourhood
 also invertible if
 gradient is continuous

$$\lim_{k \rightarrow \infty} \frac{\|x^{(k+1)} - x^*\|}{\|x^{(k)} - x^*\|} = 0.$$

2. If $\exists L, M > 0$ s.t.

$$\begin{aligned}\|\nabla g(x) - \nabla g(y)\| &\leq L\|x - y\| \\ \|\nabla g(x)^{-1}\| &\leq M\end{aligned}\quad \forall x, y$$

then if $x^{(0)} \in S_S$, we have

\leftarrow

$$\|x^{(k+1)} - x^*\| \leq \frac{LM}{2} \|x^{(k)} - x^*\|^2$$

Convergence is quadratic.

Pf: Choose $\delta > 0$ s.t. $(\nabla g(x))^{-1}$ exists for all $x \in S_S$.

Claim: $g(x^{(k)}) = \int_0^1 \nabla g(x^* + t(x^{(k)} - x^*))^T (x^{(k)} - x^*) dt$.

Pf: Let $h(t) = g(x^* + t(x^{(k)} - x^*))$

$$h(1) - h(0) = \int_0^1 \frac{dh}{dt} dt \quad - (1)$$

$$\begin{array}{c} \parallel \\ g(x^{(k)}) \\ \parallel \\ g(x^*) \end{array}$$

By the chain rule applied to h ,

$$\frac{dh}{dt} = \nabla g(x^* + t(x^{(k)} - x^*))^T (x^{(k)} - x^*) \quad - (2)$$

$$g(x^{(k)}) = \int_0^1 \nabla g(x^* + t(x^{(k)} - x^*))^\top (x^{(k)} - x^*) dt.$$

as desired. provide upper bound

Now we will estimate $\|x^{(k+1)} - x^*\|$ in terms of $\|x^{(k)} - x^*\|$.

$$\|x^{(k+1)} - x^*\| = \|x^{(k)} - \nabla g(x^{(k)})^{-1} g(x^{(k)}) - x^*\|$$

$$= \|x^{(k)} - x^* - \nabla g(x^{(k)})^{-1} g(x^{(k)})\|$$

$$= \|\nabla g(x^{(k)})^{-1} (\nabla g(x^{(k)}) (x^{(k)} - x^*) - g(x^{(k)}))\|$$

$$\stackrel{\text{Claim}}{=} \|\nabla g(x^{(k)})^{-1} \left[\nabla g(x^{(k)}) - \int_0^1 \nabla g(x^* + t(x^{(k)} - x^*))^\top dt \right] (x^{(k)} - x^*)\|$$

$$= \|\nabla g(x^{(k)})^{-1} \left[\int_0^1 \nabla g(x^{(k)}) - \nabla g(x^* + t(x^{(k)} - x^*))^\top dt \right] (x^{(k)} - x^*)\|$$

$$\leq \underbrace{\|\nabla g(x^{(k)})^{-1}\|}_{\leq M} \left\| \int_0^1 \nabla g(x^{(k)}) - \nabla g(x^* + t(x^{(k)} - x^*))^\top dt \right\| \cdot \|x^{(k)} - x^*\|$$

$$\leq M \underbrace{\left\| \int_0^1 \nabla g(x^{(k)}) - \nabla g(x^* + t(x^{(k)} - x^*))^\top dt \right\|}_{A < \varepsilon} \|x^{(k)} - x^*\|$$

By the continuity of ∇g , we may take δ suff. small to make A arbitrarily small.

$\Rightarrow x^{(k)}$ converges x^*

& $\{\|x^{(k)} - x^*\|\}$ converges superlinearly fast

$\therefore A$ can be made arbitrarily small.

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

Part 2: $\|x^{(k+1)} - x^*\| \leq M A \|x^{(k)} - x^*\|$

$$A = \left\| \int_0^1 \nabla g(x^{(k)}) - \nabla g(x^* + t(x^{(k)} - x^*)) dt \right\|$$

$$\leq \int_0^1 \left\| \nabla g(x^{(k)}) - \nabla g(x^* + t(x^{(k)} - x^*)) \right\| dt.$$

$$\leq L \int_0^1 \left\| x^{(k)} - (x^* + t(x^{(k)} - x^*)) \right\| dt.$$

$$= L \int_0^1 (1-t) \|x^{(k)} - x^*\| dt.$$

$$= L \|x^{(k)} - x^*\| \underbrace{\int_0^1 1-t dt}_{\frac{1}{2}}$$

$$A \leq \frac{L}{2} \|x^{(k)} - x^*\|$$

$$\|x^{(k+1)} - x^*\| \leq M A \|x^{(k)} - x^*\|$$

$$\leq \frac{LM}{2} \|x^{(k)} - x^*\|^2 \quad (\text{qed}) \quad //.$$