

Numerical Optimization

Unconstrained Optimization

Shirish Shevade

Computer Science and Automation
Indian Institute of Science
Bangalore 560 012, India.

NPTEL Course on Numerical Optimization

Consider the problem to minimize

$$\min f(\mathbf{x}) \triangleq \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{c}^T \mathbf{x}$$

where \mathbf{H} is a symmetric positive definite matrix.

- *Condition number* of the Hessian matrix controls the convergence rate of steepest descent method.
- Faster convergence if the Hessian matrix is \mathbf{I}
- Let $\mathbf{H} = \mathbf{L}\mathbf{L}^T$ be the Cholesky decomposition of \mathbf{H}
- Define $\mathbf{y} = \mathbf{L}^T \mathbf{x}$. Therefore, the function $f(\mathbf{x})$ is transformed to the function $h(\mathbf{y})$.

$$h(\mathbf{y}) \triangleq f(\mathbf{L}^{-T} \mathbf{y})$$

$$\begin{aligned}
h(\mathbf{y}) &= f(\mathbf{L}^{-T}\mathbf{y}) \\
&= \frac{1}{2}\mathbf{y}^T\mathbf{L}^{-1}\mathbf{H}\mathbf{L}^{-T}\mathbf{y} - \mathbf{c}^T\mathbf{L}^{-T}\mathbf{y} \\
&= \frac{1}{2}\mathbf{y}^T\mathbf{L}^{-1}\mathbf{L}\mathbf{L}^T\mathbf{L}^{-T}\mathbf{y} - \mathbf{c}^T\mathbf{L}^{-T}\mathbf{y} \\
&= \frac{1}{2}\mathbf{y}^T\mathbf{y} - \mathbf{c}^T\mathbf{L}^{-T}\mathbf{y}
\end{aligned}$$

- The Hessian matrix of $h(\mathbf{y})$ is \mathbf{I}
- Let us apply steepest descent method in \mathbf{y} -space

$$\begin{aligned}
\mathbf{y}^{k+1} &= \mathbf{y}^k - \nabla h(\mathbf{y}^k) \\
&= \mathbf{y}^k - \mathbf{L}^{-1}\nabla f(\mathbf{L}^{-T}\mathbf{y}^k) \\
\therefore \mathbf{L}^{-T}\mathbf{y}^{k+1} &= \mathbf{L}^{-T}\mathbf{y}^k - \mathbf{L}^{-T}\mathbf{L}^{-1}\nabla f(\mathbf{L}^{-T}\mathbf{y}^k) \\
\therefore \mathbf{x}^{k+1} &= \mathbf{x}^k - \mathbf{H}^{-1}\nabla f(\mathbf{x}^k)
\end{aligned}$$

Newton Method

Consider the problem,

$$\min_{\mathbf{x}} f(\mathbf{x})$$

- Let $f \in \mathcal{C}^2$ and f be bounded below.
- Use second order information to find a descent direction
- At every iteration, use Taylor series to approximate f at \mathbf{x}^k by a quadratic function and find the minimum of this quadratic function to get \mathbf{x}^{k+1}

$$\begin{aligned} f(\mathbf{x}) \approx f_q(\mathbf{x}) &= f(\mathbf{x}^k) + \mathbf{g}^{kT}(\mathbf{x} - \mathbf{x}^k) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^k)^T \mathbf{H}^k (\mathbf{x} - \mathbf{x}^k) \\ \mathbf{x}^{k+1} &= \arg \min_{\mathbf{x}} f_q(\mathbf{x}) \end{aligned}$$

- $\nabla f_q(\mathbf{x}) = 0 \Rightarrow \mathbf{x}^{k+1} = \mathbf{x}^k - (\mathbf{H}^k)^{-1} \mathbf{g}^k$ (assuming \mathbf{H}^k is invertible)

$\mathbf{x}^{k+1} = \mathbf{x}^k - (\mathbf{H}^k)^{-1} \mathbf{g}^k$ is of the form, $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha^k \mathbf{d}^k$

- **Classical Newton Method:**

- Newton Direction: $\mathbf{d}_N^k = -(\mathbf{H}^k)^{-1} \mathbf{g}^k$
- Step Length: $\alpha^k = 1$

- Is \mathbf{d}_N^k a descent direction?

$\mathbf{g}^{kT} \mathbf{d}_N^k = -\mathbf{g}^{kT} (\mathbf{H}^k)^{-1} \mathbf{g}^k < 0$ if \mathbf{H}^k is positive definite.

\mathbf{d}_N^k is a descent direction if \mathbf{H}^k is positive definite

- Consider the problem to minimize, $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} - \mathbf{c}^T \mathbf{x}$ where \mathbf{H} is a symmetric positive definite matrix.

$\mathbf{g}(\mathbf{x}) = 0 \Rightarrow \mathbf{x}^* = \mathbf{H}^{-1} \mathbf{c}$ is a strict local minimum

Let $\mathbf{x}^0 \in \mathbb{R}^n$ be any point. $\mathbf{g}(\mathbf{x}^0) = \mathbf{H} \mathbf{x}^0 - \mathbf{c}$, $\mathbf{H}(\mathbf{x}^0) = \mathbf{H}$.

Using classical Newton method,

$$\mathbf{x}^1 = \mathbf{x}^0 - \mathbf{H}^{-1}(\mathbf{H} \mathbf{x}^0 - \mathbf{c}) = \mathbf{H}^{-1} \mathbf{c} = \mathbf{x}^*.$$

Using classical newton method, the minimum of a strictly convex quadratic function (with invertible Hessian matrix) is attained in **one** iteration from *any starting point*.

Classical Newton Algorithm

(1) Initialize \mathbf{x}^0 and ϵ , set $k := 0$.

(2) **while** $\|\mathbf{g}^k\| > \epsilon$

(a) $\mathbf{d}^k = -(\mathbf{H}^k)^{-1} \mathbf{g}^k$

(b) $\alpha^k = 1$

(c) $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha^k \mathbf{d}^k$

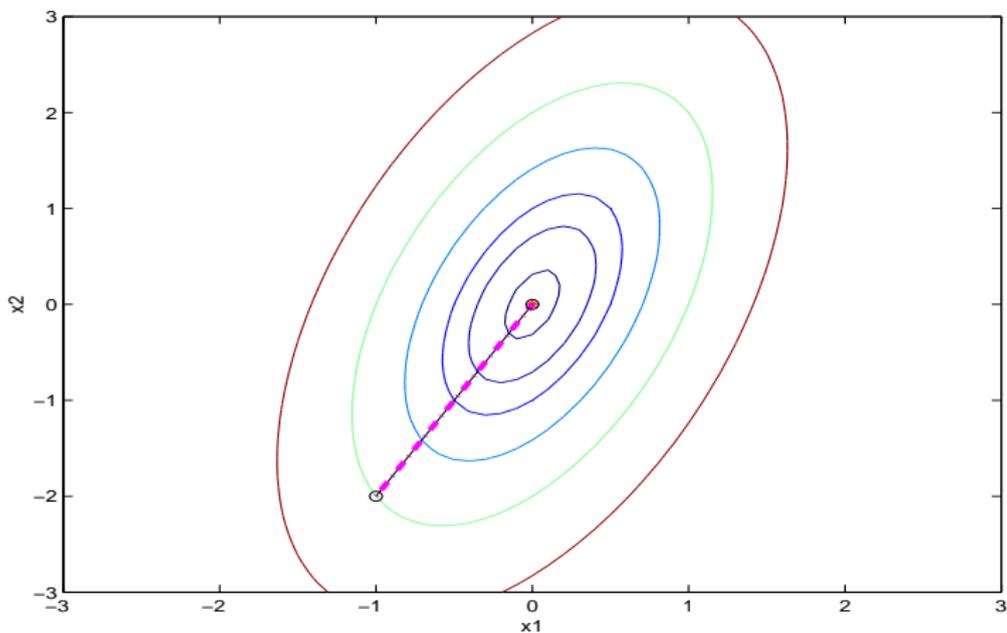
(d) $k := k + 1$

endwhile

Output : $\mathbf{x}^* = \mathbf{x}^k$, a stationary point of $f(\mathbf{x})$.

Example:

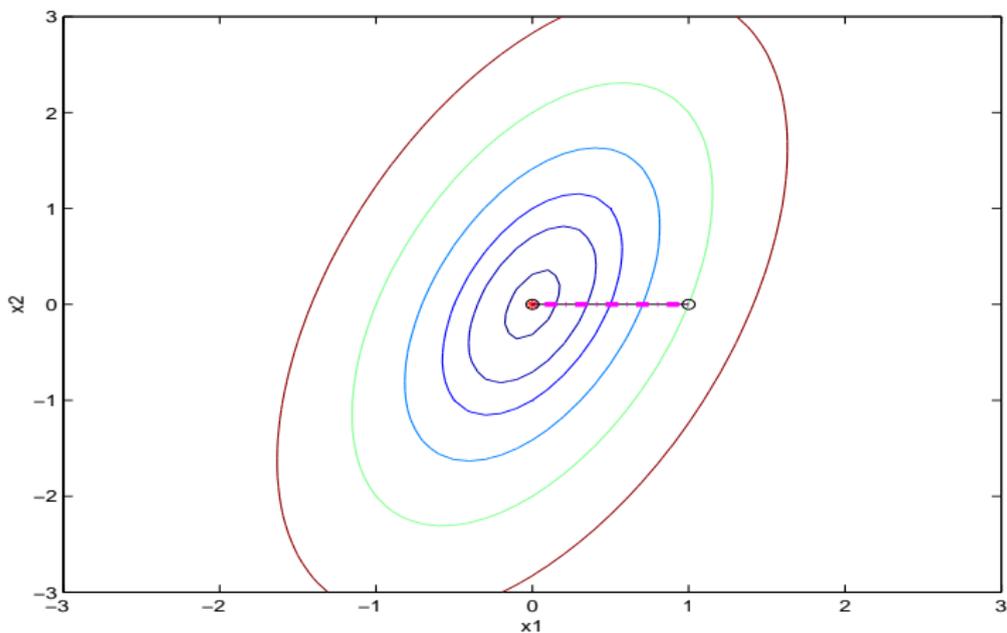
$$\min f(\mathbf{x}) \triangleq 4x_1^2 + x_2^2 - 2x_1x_2$$



Classical Newton algorithm applied to $f(\mathbf{x})$ converges to \mathbf{x}^* in **one** iteration from any starting point

Example:

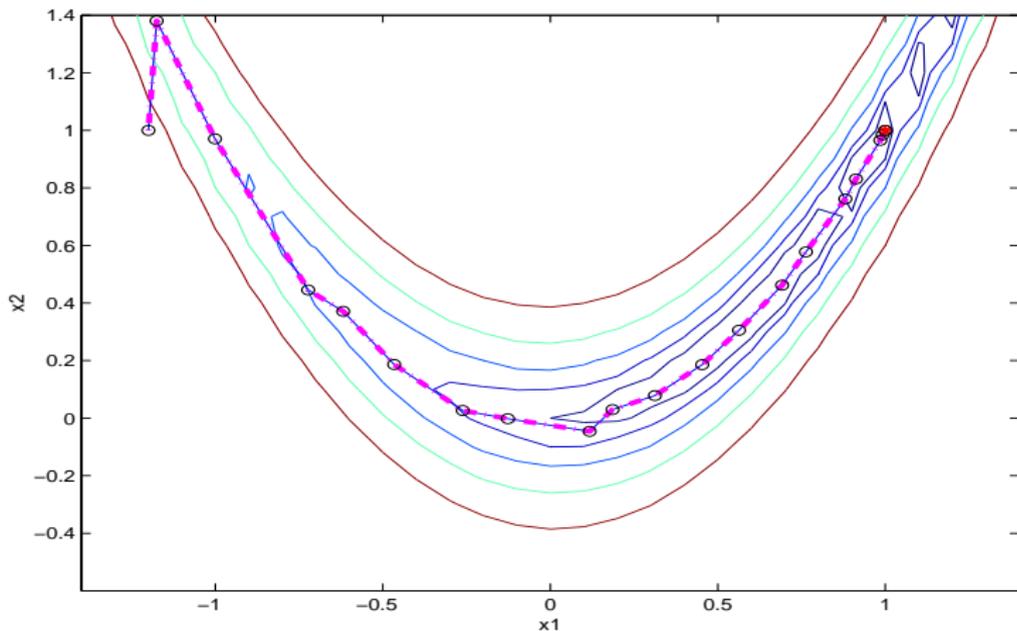
$$\min f(\mathbf{x}) \triangleq 4x_1^2 + x_2^2 - 2x_1x_2$$



Classical Newton algorithm applied to $f(\mathbf{x})$ converges to \mathbf{x}^* in **one** iteration from any starting point

Example (Rosenbrock function):

$$\min f(\mathbf{x}) \triangleq 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$



Behaviour of classical Newton algorithm (with backtracking line search) applied to $f(\mathbf{x})$ using $\mathbf{x}^0 = (-1.2, 1)^T$

Example (Rosenbrock function):

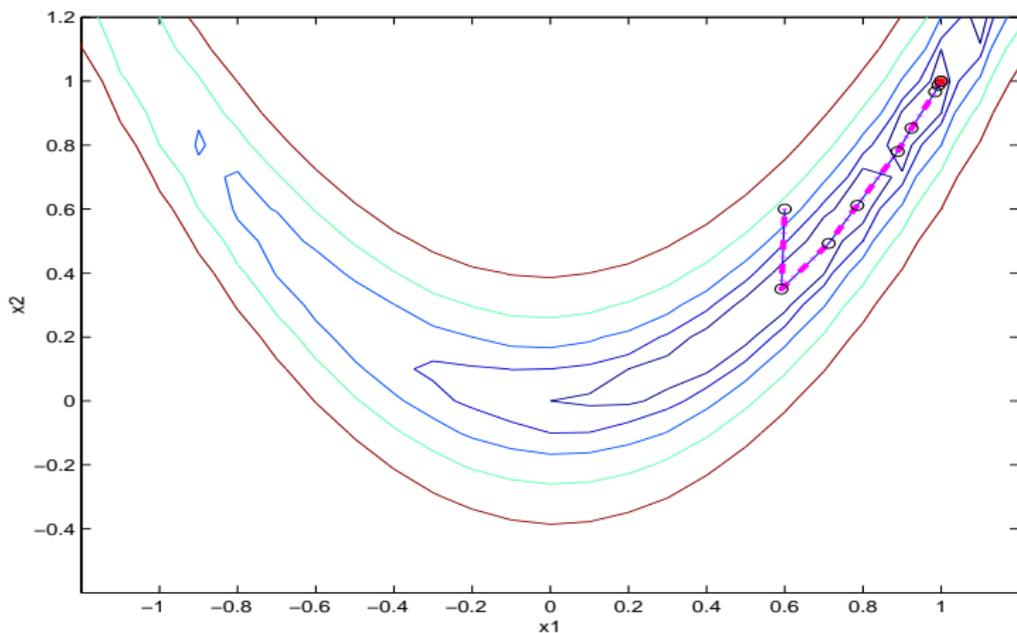
$$\min f(\mathbf{x}) \triangleq 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

k	x_1^k	x_2^k	$f(\mathbf{x}^k)$	$\ \mathbf{g}^k\ $	$\ \mathbf{x}^k - \mathbf{x}^*\ $
0	-1.2	1	24.2	232.86	2.2
1	-1.17	1.38	4.73	4.64	2.21
2	-1.00	0.97	4.01	17.54	2.00
3	-0.72	0.45	3.57	30.06	1.81
4	-0.62	0.37	2.63	6.34	1.74
5	-0.47	0.19	2.24	10.64	1.68
10	0.31	0.08	0.51	4.00	1.15
15	0.88	0.76	0.03	5.37	0.27
20	0.99	0.99	7.38×10^{-13}	1.3×10^{-6}	1.9×10^{-6}

Table: Classical Newton algorithm (with backtracking line search) applied to Rosenbrock function, using $\mathbf{x}^0 = (-1.2, 1.0)^T$, $\hat{\alpha} = 1$, $\rho = .3$ and $c_1 = 1.0 \times 10^{-4}$.

Example (Rosenbrock function):

$$\min f(\mathbf{x}) \triangleq 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$



Behaviour of classical Newton algorithm (with backtracking line search) applied to $f(\mathbf{x})$ using $\mathbf{x}^0 = (0.6, 0.6)^T$

Example (Rosenbrock function):

$$\min f(\mathbf{x}) \triangleq 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

k	x_1^k	x_2^k	$f(\mathbf{x}^k)$	$\ \mathbf{g}^k\ $	$\ \mathbf{x}^k - \mathbf{x}^*\ $
0	0.6	0.6	5.92	75.5947	0.57
1	0.59	0.35	0.17	0.80	0.77
2	0.71	0.49	0.10	4.64	0.58
3	0.79	0.61	0.05	1.65	0.44
4	0.89	0.78	0.02	4.18	0.25
5	0.92	0.85	0.01	0.40	0.17
9	0.99	0.99	5.76×10^{-13}	2.77×10^{-6}	1.69×10^{-6}

Table: Classical Newton algorithm (with backtracking line search) applied to Rosenbrock function, using $\mathbf{x}^0 = (0.6, 0.6)^T$, $\hat{\alpha} = 1$, $\rho = .3$ and $c_1 = 1.0 \times 10^{-4}$.

Classical Newton Algorithm

(1) Initialize \mathbf{x}^0 and ϵ , set $k := 0$.

(2) **while** $\|\mathbf{g}^k\| > \epsilon$

(a) $\mathbf{d}^k = -(\mathbf{H}^k)^{-1} \mathbf{g}^k$

(b) $\alpha^k = 1$

(c) $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha^k \mathbf{d}^k$

(d) $k := k + 1$

endwhile

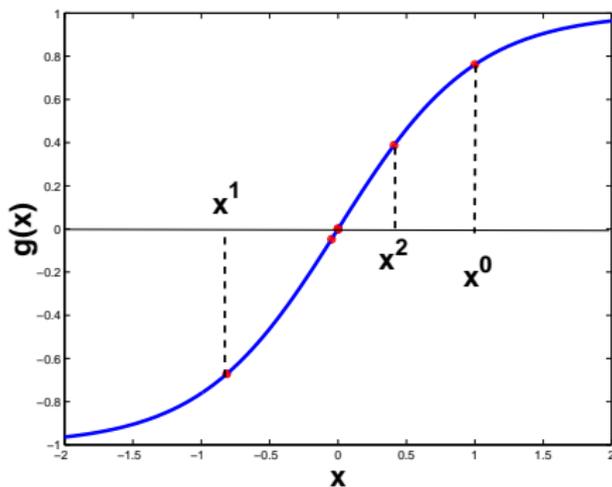
Output : $\mathbf{x}^* = \mathbf{x}^k$, a stationary point of $f(\mathbf{x})$.

- Requires $O(n^3)$ computational effort for every iteration (Step 2(a))
- *No guarantee* that \mathbf{d}^k is a descent direction
- *No guarantee* that $f(\mathbf{x}^{k+1}) < f(\mathbf{x}^k)$ (no line search)
- Sensitive to initial point (for non-quadratic functions)

Consider the problem,

$$\min_{x \in \mathbb{R}} \log(e^x + e^{-x})$$

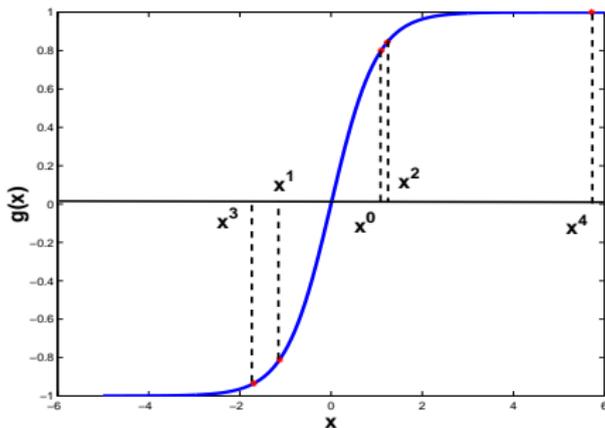
- $f(x) = \log(e^x + e^{-x})$
- $g(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$



Consider the problem,

$$\min_{x \in \mathbb{R}} \log(e^x + e^{-x})$$

- $f(x) = \log(e^x + e^{-x})$
- $g(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$



Classical Newton algorithm does not converge with this initialization of x^0

Definition

An iterative optimization algorithm is said to be **locally convergent** if for each solution \mathbf{x}^* , there exists $\delta > 0$ such that for any initial point $\mathbf{x}^0 \in B(\mathbf{x}^*, \delta)$, the algorithm produces a sequence $\{\mathbf{x}^k\}$ which converges to \mathbf{x}^* .

- **Classical Newton algorithm is locally convergent**

Let $f : \mathbb{R} \rightarrow \mathbb{R}$, $f \in \mathcal{C}^2$.

Consider the problem:

$$\min f(x)$$

Let $x^* \in \mathbb{R}$ be such that $g(x^*) = 0$ and $g'(x^*) > 0$.

Assume that x^0 is *sufficiently* close to x^* .

Suppose we apply classical Newton algorithm to minimize $f(x)$.

At k -th iteration,

$$\begin{aligned}x^{k+1} &= x^k - \frac{g(x^k)}{g'(x^k)} \\ \therefore x^{k+1} - x^* &= x^k - x^* - \frac{g(x^k) - g(x^*)}{g'(x^k)} \\ &= -\frac{(g(x^k) - g(x^*) + g'(x^k)(x^* - x^k))}{g'(x^k)}\end{aligned}$$

If we assume that $f \in \mathcal{C}^3$ (or $g \in \mathcal{C}^2$), then using truncated Taylor series,

$$g(x^*) = g(x^k) + g'(x^k)(x^* - x^k) + \frac{1}{2}g''(\bar{x}^k)(x^* - x^k)^2$$

where $\bar{x}^k \in LS(x^*, x^k)$.

Therefore,

$$x^{k+1} - x^* = \frac{1}{2} \frac{g''(\bar{x}^k)}{g'(x^k)} (x^k - x^*)^2$$

$$|x^{k+1} - x^*| = \frac{1}{2} \frac{|g''(\bar{x}^k)|}{|g'(x^k)|} |x^k - x^*|^2$$

Suppose there exist α_1 and α_2 such that

$$\begin{aligned} |g''(\bar{x}^k)| &< \alpha_1 \quad \forall \bar{x}^k \in LS(x^*, x^k) \text{ and} \\ |g'(x^k)| &> \alpha_2 \quad \text{for } x^k \text{ sufficiently close to } x^*, \end{aligned}$$

then

$$|x^{k+1} - x^*| \leq \frac{\alpha_1}{2\alpha_2} |x^k - x^*|^2 \quad (\text{order two convergence if } x^k \rightarrow x^*)$$

Note that

$$|x^{k+1} - x^*| \leq \underbrace{\frac{\alpha_1}{2\alpha_2}}_{\text{required to be } < 1} |x^k - x^*|$$

If $\frac{\alpha_1}{2\alpha_2} |x^k - x^*| < 1 \forall k$, then

$$|x^{k+1} - x^*| < |x^k - x^*| \forall k$$

How to choose α_1 and α_2 ?

At x^* , $g(x^*) = 0$, and $g'(x^*) > 0$

Since $g' \in C^0$, $\exists \eta > 0 \ni g'(x) > 0 \forall x \in (x^* - \eta, x^* + \eta)$

Let

$$\alpha_1 = \max_{x \in (x^* - \eta, x^* + \eta)} |g''(x)|$$

$$\alpha_2 = \min_{x \in (x^* - \eta, x^* + \eta)} g'(x)$$

Therefore,

$$\left| \frac{1}{2} \frac{g''(\bar{x}^k)}{g'(\bar{x}^k)} \right| \leq \frac{\alpha_1}{2\alpha_2} = \beta, \text{ say.}$$

Preferable to choose $x^0 \in (x^* - \eta, x^* + \eta)$

Also, we want $\beta|x^k - x^*| < 1 \forall k$. That is,

$$\begin{aligned} |x^k - x^*| &< 1/\beta \forall k \\ \Rightarrow x^k &\in (x^* - 1/\beta, x^* + 1/\beta) \end{aligned}$$

Therefore, choose $x^0 \in (x^* - \eta, x^* + \eta) \cap (x^* - 1/\beta, x^* + 1/\beta)$

Does $\{x^k\}$ converge to x^* if x^0 is chosen using this approach?

We have

$$\begin{aligned} |x^k - x^*| &\leq \beta|x^{k-1} - x^*|^2 \\ \therefore \beta|x^k - x^*| &\leq (\beta|x^0 - x^*|)^{2^k} \\ \therefore |x^k - x^*| &\leq \frac{1}{\beta} \underbrace{(\beta|x^0 - x^*|)^{2^k}}_{<1} \end{aligned}$$

Therefore,

$$\lim_{k \rightarrow \infty} |x^k - x^*| = 0$$

Not a practical approach to choose x^0

Theorem

Let $f : \mathbb{R} \rightarrow \mathbb{R}$, $f \in \mathcal{C}^3$. Let $x^* \in \mathbb{R}$ be such that $g(x^*) = 0$ and $g'(x^*) > 0$. Then, provided x^0 is *sufficiently close to x^** , the sequence $\{x^k\}$ generated by classical Newton algorithm converges to x^* with an *order of convergence two*.

Initialization of x^0 requires knowledge of x^* !

Modified Newton Method

Modifications:

- Given \mathbf{x}^k and $\mathbf{d}_N^k = -(\mathbf{H}^k)^{-1}\mathbf{g}^k$,
Fix some constant $\delta > 0$.
Find the smallest $\zeta_k \geq 0$ such that the smallest eigenvalue of the matrix $(\mathbf{H}^k + \zeta_k \mathbf{I})$ is greater than δ .
Therefore, $\mathbf{d}^k = -(\mathbf{H}^k + \zeta_k \mathbf{I})^{-1}\mathbf{g}^k$ is a descent direction.
- Given \mathbf{x}^k and $\mathbf{d}^k = -(\mathbf{H}^k + \zeta_k \mathbf{I})^{-1}\mathbf{g}^k$, use line search techniques to determine α^k and \mathbf{x}^{k+1}

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha^k \mathbf{d}^k$$

Modified Newton Algorithm

- (1) Initialize \mathbf{x}^0 , ϵ and δ , set $k := 0$.
 - (2) **while** $\|\mathbf{g}^k\| > \epsilon$
 - (a) Find the smallest $\zeta_k \geq 0$ such that the smallest eigenvalue of $\mathbf{H}^k + \zeta_k \mathbf{I}$ is greater than δ
 - (b) Set $\mathbf{d}^k = -(\mathbf{H}^k + \zeta_k \mathbf{I})^{-1} \mathbf{g}^k$
 - (c) Find $\alpha^k (> 0)$ along \mathbf{d}^k such that
 - (i) $f(\mathbf{x}^k + \alpha^k \mathbf{d}^k) < f(\mathbf{x}^k)$
 - (ii) α^k satisfies Armijo-Wolfe (or Armijo-Goldstein) conditions
 - (d) $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha^k \mathbf{d}^k$
 - (e) $k := k + 1$
- endwhile**

Output : $\mathbf{x}^* = \mathbf{x}^k$, a stationary point of $f(\mathbf{x})$.

- Modified Newton algorithm has global convergence properties and has order of convergence equal to two