

CHEM-E7225
2023

Preliminaries

The Newton
method

Newton-type
methods

Convergence



Aalto University

Root-finding with Newton-type methods

CHEM-E7225 (was E7195), 2023

Francesco Corona (☹_☹)

Chemical and Metallurgical Engineering
School of Chemical Engineering

CHEM-E7225
2023

Preliminaries

The Newton
method

Newton-type
methods

Convergence

Overview

Some notions in mathematical and numerical analysis that are used in optimisation

- Only instrumental concepts, to solve optimal control problems

Optimisation refers to the problem of finding the value of the inputs (independent variables) to some function such that the corresponding outputs (dependent variables) take an optimal value, where optimality is defined in some sense by the function itself

- The task can be formulated as a root-finding problem
- As the problem of finding the zeros of a function

We will focus to a specific class of solution approaches known as Newton-type methods

Preliminaries

The Newton
method

Newton-type
methods

Convergence

Preliminaries

Root-finding with Newton-type methods

Preliminaries

The Newton
method

Newton-type
methods

Convergence

Preliminaries

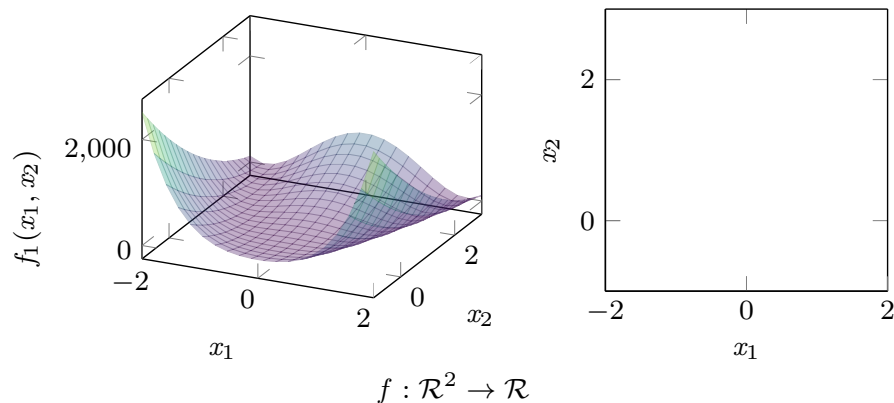
Let function f be a twice-differentiable function (first and second derivatives) on \mathcal{R}^N

$$f : \mathcal{R}^N \rightarrow \mathcal{R}, \quad f \in \mathcal{C}^2(\mathcal{R}^N)$$

We shall use function f to refresh some basic notions from multivariate calculus

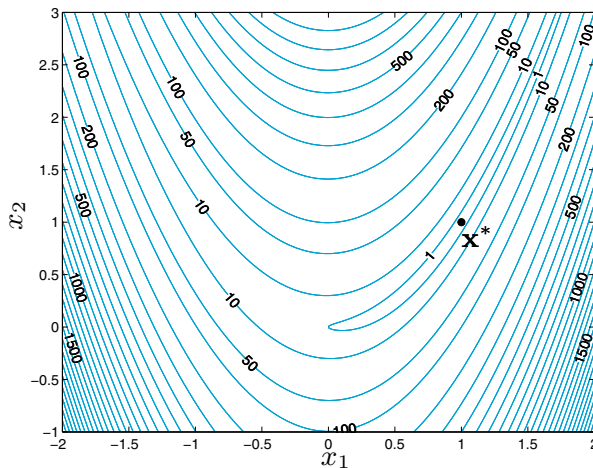
- We are mainly interested in its gradient vector and Hessian matrix

We consider a **Rosenbrock's function**, classic benchmark for optimisation methods



Preliminaries | A scalar function

Example

The **Rosenbrock's function**

$$f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

Function $f(x)$ has a global minimum

$$x^* = (1, 1)$$

□

Preliminaries | Gradient

Let the symbol $\nabla f(x)$ denote the **gradient vector** of function f at some point $x \in \mathcal{R}^N$

$$\nabla f(x_1, x_2, \dots, x_N) = \underbrace{\begin{bmatrix} \frac{\partial f(x_1, x_2, \dots, x_N)}{\partial x_1} \\ \frac{\partial f(x_1, x_2, \dots, x_N)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x_1, x_2, \dots, x_N)}{\partial x_N} \end{bmatrix}}_{N \times 1}$$

At any point $x \in \mathcal{R}^N$, we can define a vector of first derivatives, the gradient of f at x

$$\nabla f(x) = \left[\frac{\partial f(x)}{\partial x_1} \quad \frac{\partial f(x)}{\partial x_2} \quad \cdots \quad \frac{\partial f(x)}{\partial x_N} \right]_x^T$$

The symbol $\nabla = \left[\frac{\partial}{\partial x_1} \quad \frac{\partial}{\partial x_2} \quad \cdots \quad \frac{\partial}{\partial x_N} \right]^T$ denotes the **gradient operator**

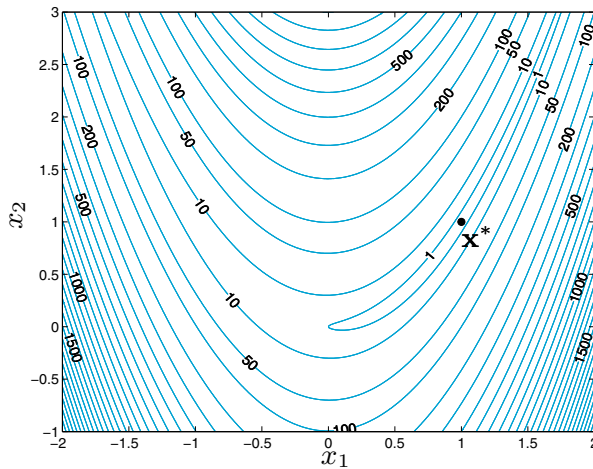
Preliminaries

The Newton
methodNewton-type
methods

Convergence

Preliminaries | Gradient (cont.)

Example

The **Rosenbrock's function**

$$f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

$$\frac{\partial f(x)}{\partial x_1} = -400x_1(x_2 - x_1^2) - 2(1 - x_1)$$

$$\frac{\partial f(x)}{\partial x_2} = 200(x_2 - x_1^2)$$

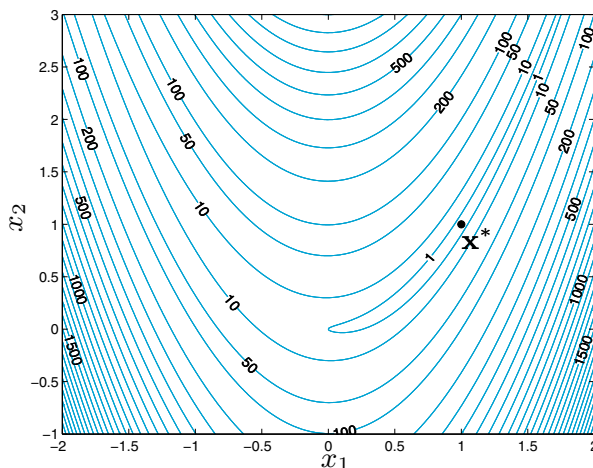
$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x_1, x_2)}{\partial x_1} \\ \frac{\partial f(x_1, x_2)}{\partial x_2} \end{bmatrix}$$

Preliminaries

The Newton
methodNewton-type
methods

Convergence

Preliminaries | Gradient (cont.)



$$\begin{aligned} \nabla f(x) &= \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \end{bmatrix} \\ &= \begin{bmatrix} -400x_1(x_2 - x_1^2) - 2(1 - x_1) \\ 200(x_2 - x_1^2) \end{bmatrix} \end{aligned}$$

Consider some point x' , say $x' = (0, 0)$, we can evaluate the gradient vector of f at x'

$$\begin{aligned} \nabla f(x') &= \begin{bmatrix} -400x'_1(x'_2 - x'^2_1) - 2(1 - x'_1) \\ 200(x'_2 - x'^2_1) \end{bmatrix} \\ &= \begin{bmatrix} -2 \\ 0 \end{bmatrix} \end{aligned}$$

□

Preliminaries

The Newton
methodNewton-type
methods

Convergence

Preliminaries | Hessian

Let the symbol $\nabla^2 f(x)$ denote the **Hessian matrix** of function f at point $x \in \mathcal{R}^N$

$$\nabla^2 f(x) = \underbrace{\begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_N} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_N} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_N \partial x_1} & \frac{\partial^2 f(x)}{\partial x_N \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_N \partial x_N} \end{bmatrix}}_{N \times N}$$

At point $x \in \mathcal{R}^N$, we can define a matrix of second derivatives, the Hessian of f at x

$$\nabla^2 f(x) = [h_{ij}]_{i,j=1}^N,$$

with $h_{ij} = \frac{\partial^2 f(x_1, x_2, \dots, x_N)}{\partial x_j \partial x_i}$, a symmetric ($N \times N$) matrix

Preliminaries

The Newton
methodNewton-type
methods

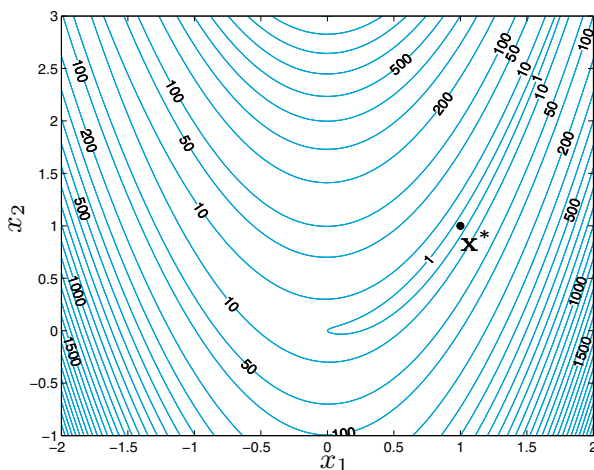
Convergence

Preliminaries | Hessian (cont.)

Example

The **Rosenbrock's function**

$$f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$



$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_2} \end{bmatrix}$$

$$\frac{\partial^2 f(x)}{\partial x_1 \partial x_1} = 1200x_1^2 - 400x_2 + 2$$

$$\frac{\partial^2 f(x)}{\partial x_1 \partial x_2} = -400x_1$$

$$\frac{\partial^2 f(x)}{\partial x_2 \partial x_1} = -400x_1$$

$$\frac{\partial^2 f(x)}{\partial x_2 \partial x_2} = 200$$

Preliminaries | Hessian (cont.)

Function $f(x)$,

$$f(x_1, x_2) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

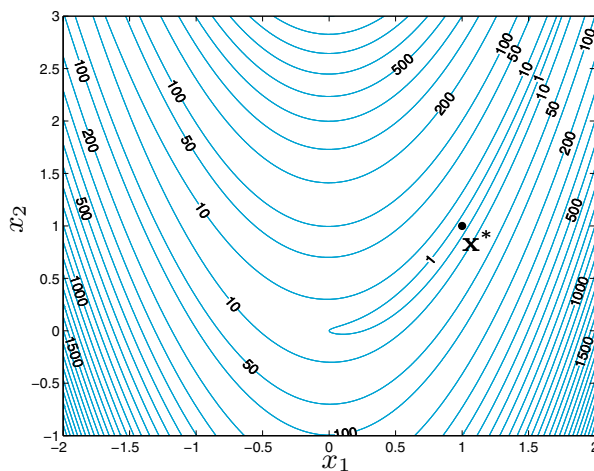
Gradient vector $\nabla f(x)$,

$$\begin{aligned}\nabla f(x) &= \left[\frac{\partial f(x)}{\partial x_1} \quad \frac{\partial f(x)}{\partial x_2} \right]^T \\ &= \begin{bmatrix} -400x_1(x_2 - x_1^2) - 2(1 - x_1) \\ 200(x_2 - x_1^2) \end{bmatrix}\end{aligned}$$

Hessian matrix $\nabla^2 f(x)$,

$$\begin{aligned}\nabla^2 f(x) &= \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_2} \end{bmatrix} \\ &= \begin{bmatrix} 1200x_1^2 - 400x_2 + 2 & -400x_1 \\ -400x_1 & 200 \end{bmatrix}\end{aligned}$$

Preliminaries | Hessian (cont.)

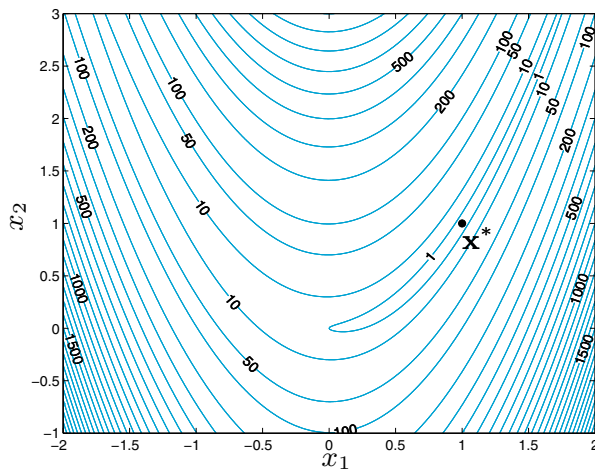


$$\begin{aligned}\nabla^2 f(x) &= \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_2} \end{bmatrix} \\ &= \begin{bmatrix} 1200x_1^2 - 400x_2 + 2 & -400x_1 \\ -400x_1 & 200 \end{bmatrix}\end{aligned}$$

Consider some point x' , say $x' = (0, 0)$, we can evaluate the Hessian matrix of f at x'

$$\begin{aligned}\nabla^2 f(x') &= \begin{bmatrix} 1200x_1'^2 - 400x_2' + 2 & -400x_1' \\ -400x_1' & 200 \end{bmatrix} \\ &= \begin{bmatrix} 2 & 0 \\ 0 & 200 \end{bmatrix}\end{aligned}$$

Preliminaries | Hessian (cont.)

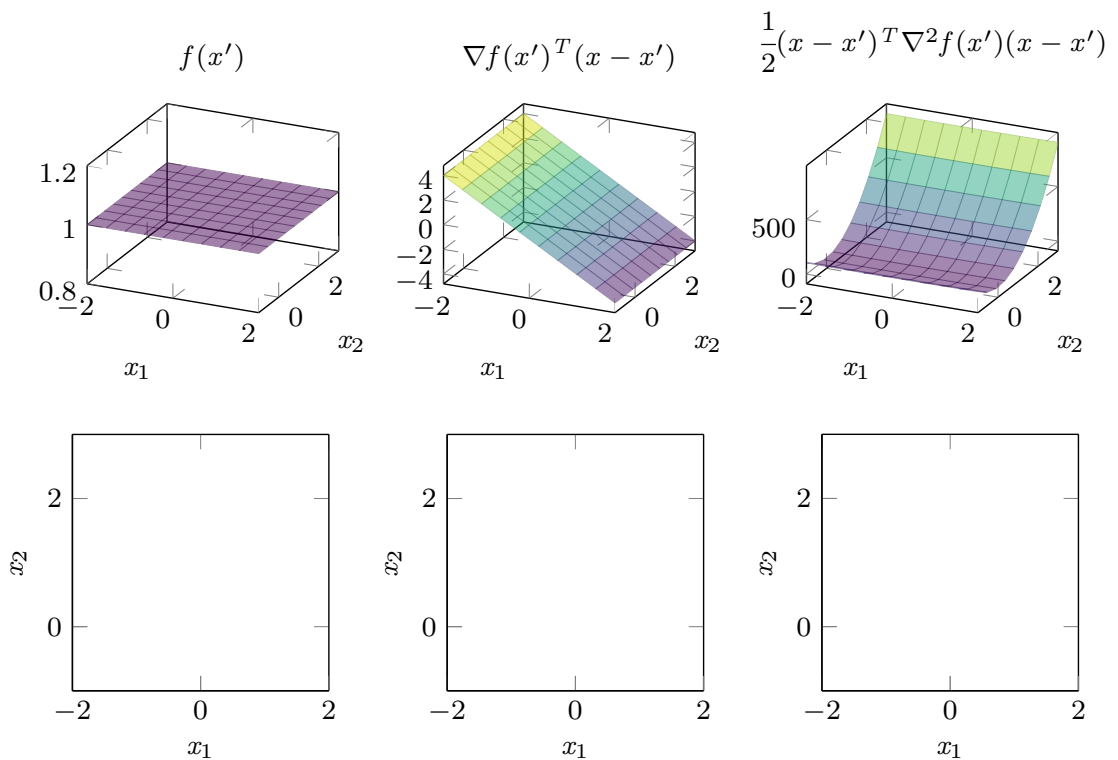


$$\begin{aligned} \nabla^2 f(x) &= \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_2} \end{bmatrix} \\ &= \begin{bmatrix} 1200x_1^2 - 400x_2 + 2 & -400x_1 \\ -400x_1 & 200 \end{bmatrix} \end{aligned}$$

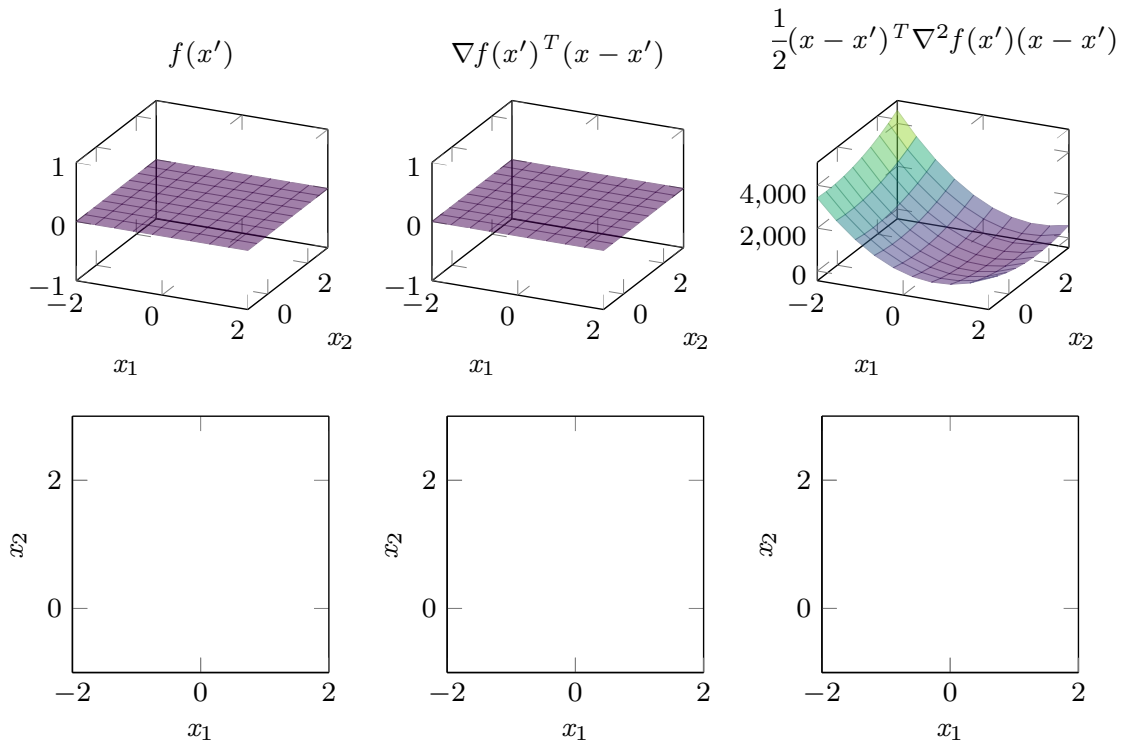
Consider some point x' , say $x' = (1, 1)$, we can evaluate the Hessian matrix of f at x'

$$\begin{aligned} \nabla^2 f(x') &= \begin{bmatrix} 1200x_1'^2 - 400x_2' + 2 & -400x_1' \\ -400x_1' & 200 \end{bmatrix} \\ &= \begin{bmatrix} 802 & -400 \\ -400 & 200 \end{bmatrix} \end{aligned}$$

Preliminaries | Taylor's expansion



Preliminaries | Taylor's expansion(cont.)



The Newton method

Root-finding with Newton-type methods

The Newton method

Let function $g : \mathcal{R}^N \rightarrow \mathcal{R}$ with $N \geq 1$ be some function belonging to class $\mathcal{C}^2(\mathcal{R}^N)$

- We are interested in solving the system of nonlinear equations $\nabla g(x) = 0$

That is, we are interested in the point(s) $x^* = (x_1^*, x_2^*, \dots, x_N^*)$ such that

$$\nabla g(x^*) = \begin{bmatrix} \frac{\partial g(x_1, x_2, \dots, x_N)}{\partial x_1} \\ \frac{\partial g(x_1, x_2, \dots, x_N)}{\partial x_2} \\ \vdots \\ \frac{\partial g(x_1, x_2, \dots, x_N)}{\partial x_N} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Points in which all the partial derivatives of g are zero are **stationary points**

- We want to know where these fixed points or extrema are located

The Newton method (cont.)

$$\begin{aligned} \nabla g(x^*) &= \begin{bmatrix} \frac{\partial g(x_1, x_2, \dots, x_N)}{\partial x_1} \\ \frac{\partial g(x_1, x_2, \dots, x_N)}{\partial x_2} \\ \vdots \\ \frac{\partial g(x_1, x_2, \dots, x_N)}{\partial x_N} \end{bmatrix} \\ &= \begin{bmatrix} f_1(x_1, x_2, \dots, x_N) \\ f_2(x_1, x_2, \dots, x_N) \\ \vdots \\ f_N(x_1, x_2, \dots, x_N) \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \end{aligned}$$

Points x^* are points in which all the functions $\{f_1, f_2, \dots, f_N\}$ are jointly equal to zero

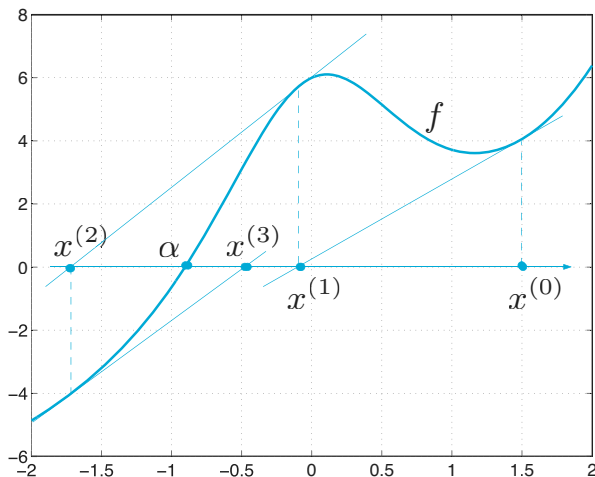
The Newton method | Baby Newton

Tangent's method (Baby Newton)

Consider the problem of finding the zero of a differentiable function $f : [a, b] \subset \mathcal{R} \rightarrow \mathcal{R}$

\rightsquigarrow We are interested in point(s) $\alpha \in [a, b]$ such that $f(\alpha) = 0$

We know the function corresponding to the tangent to $f(x)$ at any point $x^{(k)} \in [a, b]$



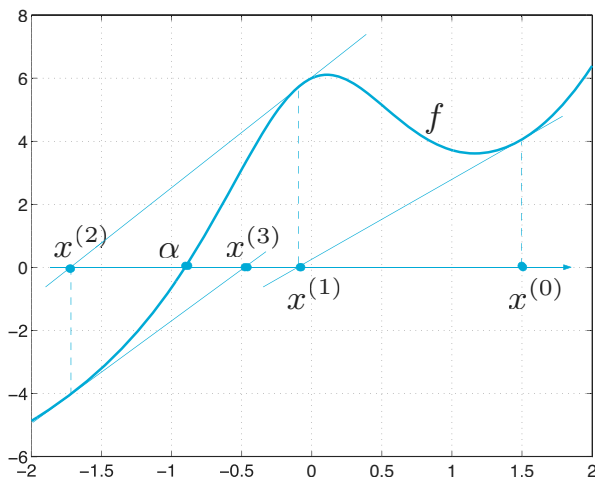
$$y(x) = f(x^{(k)}) + \underbrace{f'(x^{(k)})}_{\left. \frac{df(x)}{dx} \right|_{x^{(k)}}} (x - x^{(k)})$$

$\rightsquigarrow k = 0$, the tangent to f at $x^{(0)}$

$$y(x) = f(x^{(0)}) + f'(x^{(0)}) (x - x^{(0)})$$

The Newton method | Baby Newton (cont.)

We can use the tangent to find the point $x = x^{(k+1)}$, the point such that $y(x^{(k+1)}) = 0$



$$0 = f(x^{(k)}) + f'(x^{(k)}) (x^{(k+1)} - x^{(k)})$$

$$\rightsquigarrow x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}$$

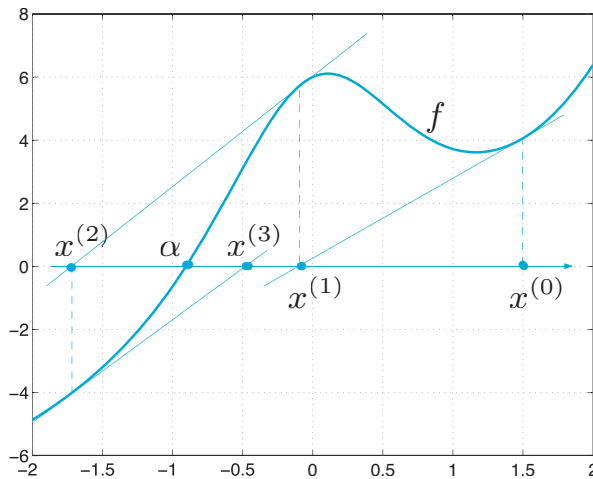
\rightsquigarrow From $x^{(0)}$, we solve for point $x^{(1)}$

$$x^{(1)} = x^{(0)} - \frac{f(x^{(0)})}{f'(x^{(0)})}$$

The Newton method | Baby Newton (cont.)

The operation can be repeated for $k = 0, 1, 2, \dots$, to find points $x^{(k+1)}$ approaching α

- As $f(x^{(k)}) \rightarrow 0$ also differences $|x^{(k+1)} - x^{(k)}| \rightarrow 0$ (asymptotically, $x^{(k)} \rightarrow \alpha$)



$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}$$

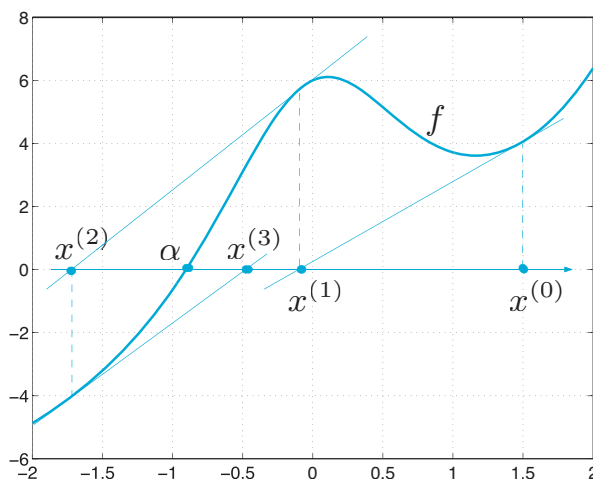
↪ Importantly, the derivative $f'(x^{(k)})$ must exist and must be non-zero

The Newton method | Baby Newton (cont.)

$$y(x) = f(x^{(k)}) + f'(x^{(k)})(x - x^{(k)})$$

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}$$

The recursion defines the sequence $\{x^{(k)}\}$ that is generated by the **Newton's method**



↪ The method reduces to locally substituting f with its tangent

↪ The substitution is repeated until **convergence**

The Newton method | Convergence tests

Newton's method returns the exact value α after an infinite number of iterations

- When it converges

In general, we would be satisfied to obtain an approximation which is ε -accurate

- We would end the recursion when the user-specified tolerance ε is reached

$$\underbrace{\|\alpha - x^{(k)}\|}_{e^{(k)}} < \varepsilon$$

To perform the test, we would need to know the exact value of α (the unknown)

In practice, an estimator of the error based on measurable quantities is used

- For the Newton's method, this could be the difference between iterates

$$\underbrace{\|x^{(k)} - x^{(k-1)}\|}_{\hat{e}^{(k)}} < \varepsilon$$

- Alternatively, there is also the possibility to use the residuals

$$\|f(x^{(k)})\| < \varepsilon$$

The Newton method (cont.)

Consider the set of nonlinear equations from the vector-valued function $f : \mathcal{R}^N \rightarrow \mathcal{R}^N$

$$\begin{aligned} f_1(x_1, x_2, \dots, x_N) &= 0 \\ f_2(x_1, x_2, \dots, x_N) &= 0 \\ &\vdots \\ f_N(x_1, x_2, \dots, x_N) &= 0 \end{aligned}$$

As we are letting $f = (f_1, \dots, f_N)^T$ and $x = (x_1, \dots, x_N)^T$, we can re-write the system

$$f(x) = 0$$

We are interested in solving this system of equations, by extending Newton's method

$$\begin{aligned} x^{(k+1)} &= x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})} \\ &= x^{(k)} - \left(f'(x^{(k)})\right)^{-1} f(x^{(k)}) \end{aligned}$$

For $N > 1$, we firstly need to replace the first derivative with a set of first derivatives

The Newton method | The Jacobian

$$\begin{aligned} f_1(x_1, x_2, \dots, x_N) &= 0 \\ f_2(x_1, x_2, \dots, x_N) &= 0 \\ &\vdots \\ f_N(x_1, x_2, \dots, x_N) &= 0 \end{aligned}$$

We have N functions, each in N variables, we collect all the derivatives of function f

Let the symbol $J_f(x)$ denote the **Jacobian matrix** of function f

$$J_f(x) = \underbrace{\begin{bmatrix} \frac{\partial f_1(x)}{\partial x_1} & \frac{\partial f_1(x)}{\partial x_2} & \dots & \frac{\partial f_1(x)}{\partial x_N} \\ \frac{\partial f_2(x)}{\partial x_1} & \frac{\partial f_2(x)}{\partial x_2} & \dots & \frac{\partial f_2(x)}{\partial x_N} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_N(x)}{\partial x_1} & \frac{\partial f_N(x)}{\partial x_2} & \dots & \frac{\partial f_N(x)}{\partial x_N} \end{bmatrix}}_{N \times N}$$

The Jacobian J_f of function f is the multivariate equivalent of the first derivative f'

The Newton method | The Jacobian (cont.)

Note how the N rows of the Jacobian matrix of the vector-valued function f correspond to the transpose of the N gradient vectors of the (scalar) components of function f

$$\begin{aligned} J_f(x) &= \underbrace{\begin{bmatrix} \frac{\partial f_1(x)}{\partial x_1} & \frac{\partial f_1(x)}{\partial x_2} & \dots & \frac{\partial f_1(x)}{\partial x_N} \\ \frac{\partial f_2(x)}{\partial x_1} & \frac{\partial f_2(x)}{\partial x_2} & \dots & \frac{\partial f_2(x)}{\partial x_N} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_N(x)}{\partial x_1} & \frac{\partial f_N(x)}{\partial x_2} & \dots & \frac{\partial f_N(x)}{\partial x_N} \end{bmatrix}}_{N \times N} \\ &= \begin{bmatrix} \nabla f_1(x)^T \\ \nabla f_2(x)^T \\ \vdots \\ \nabla f_N(x)^T \end{bmatrix} \end{aligned}$$

The Newton method (cont.)

Consider the recursion of the univariate Newton's method, as we derived it earlier

$$x^{(k+1)} = x^{(k)} - \left(f' \left(x^{(k)} \right) \right)^{-1} f \left(x^{(k)} \right)$$

We can define $\delta x = x^{(k+1)} - x^{(k)}$ and understand the recursion as follows

$$\text{Solve for } \delta x^{(k)} \quad f' \left(x^{(k)} \right) \delta x^{(k)} = -f \left(x^{(k)} \right)$$

$$\text{Compute } x^{(k+1)} = x^{(k)} + \delta x^{(k)}$$

For the general case, let $x^{(0)} \in \mathcal{R}^N$ be an initial solution then for $k = 0, 1, \dots$

$$\text{Solve for } \delta x^{(k)} \quad J_f \left(x^{(k)} \right) \delta x^{(k)} = -f \left(x^{(k)} \right)$$

$$\text{Compute } x^{(k+1)} = x^{(k)} + \delta x^{(k)}$$

The operation is repeated until convergence

- That is, until δx is small enough

The Newton method (cont.)

$$\text{Solve for } \delta x \quad \underbrace{J_f \left(x^{(k)} \right)}_A \underbrace{\delta x^{(k)}}_x = \underbrace{-f \left(x^{(k)} \right)}_b$$

$$\text{Compute } x^{(k+1)} = x^{(k)} + \delta x^{(k)}$$

A system of linear equations with coefficient matrix $J_f \left(x^{(k)} \right)$ is solved at each iteration

$$\underbrace{\delta x^{(k)}}_x = -J_f^{-1} \left(x^{(k)} \right) f \left(x^{(k)} \right)$$

It is possible to re-write the Newton method as we derived it, as an iteration scheme

$$\text{Solve for } \delta x \quad J_f \left(x^{(k)} \right) \underbrace{\delta x^{(k)}}_{(x^{(k+1)} - x^{(k)})} = -f \left(x^{(k)} \right)$$

$$\text{Compute } x^{(k+1)} = x^{(k)} + \underbrace{\delta x^{(k)}}_{-J_f^{-1} \left(x^{(k)} \right) f \left(x^{(k)} \right)}$$

The Newton method (cont.)

Starting at $k = 0$ from some initial solution $x^{(0)}$, we have the Newton's step

$$x^{(k+1)} = x^{(k)} - J_f^{-1}(x^{(k)}) f(x^{(k)}), \quad k = 0, 1, \dots$$

Each Newton step moves $x^{(k)}$ in the (opposite) direction of vector $f(x^{(k)})$

- The direction is also rotated according to matrix $J_f^{-1}(x^{(k)})$

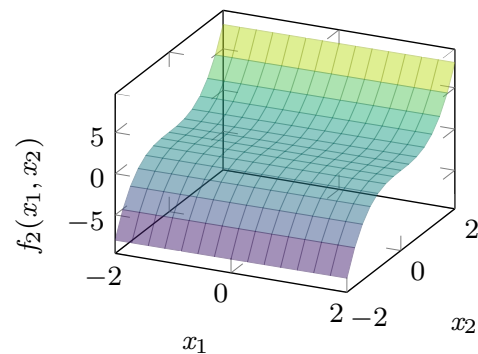
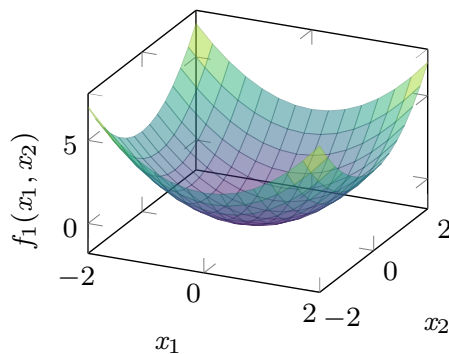
The Newton method (cont.)

Example

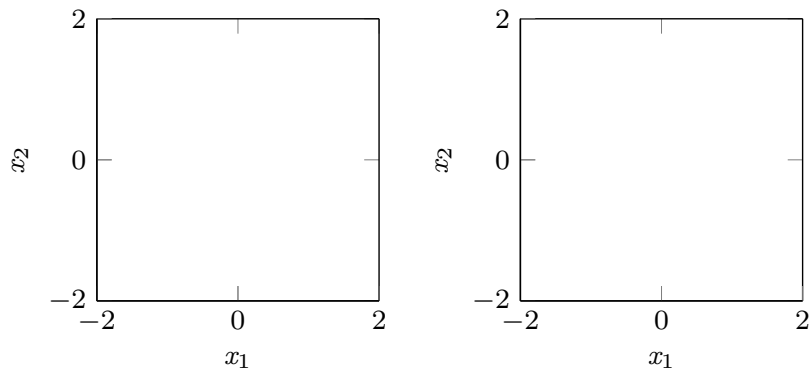
Consider the function $f(x) = (f_1(x_1, x_2), f_2(x_1, x_2))^T$

$$\begin{cases} f_1(x_1, x_2) = x_1^2 + x_2^2 - 1 \\ f_2(x_1, x_2) = \sin\left(\frac{\pi}{2}x_1\right) + x_2^3 \end{cases}$$

We are interested in point(s) $x^* = (x_1^*, x_2^*)$ where $f(x^*) = 0$



The Newton method (cont.)



The problem has two solutions, two points where both function f_1 and f_2 equal to zero

$$\approx (0.47, -0.88)^T$$

$$\approx (-0.47, 0.88)^T$$

□

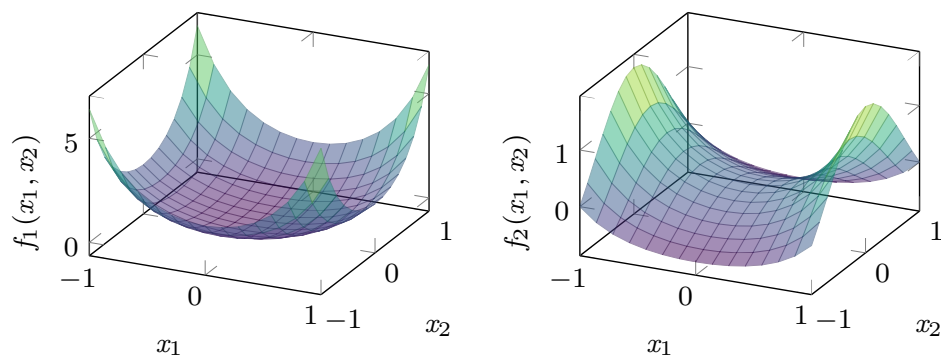
The Newton method (cont.)

Example

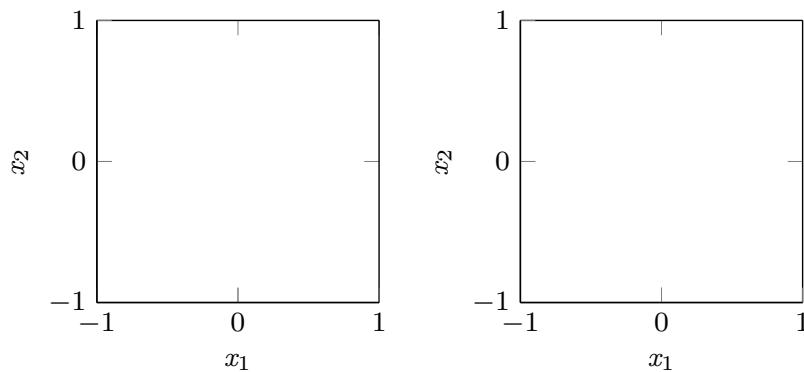
Consider the function $f(x) = (f_1(x_1, x_2), f_2(x_1, x_2))^T$

$$\begin{cases} f_1(x_1, x_2) = e^{(x_1^2 + x_2^2)} - 1 \\ f_2(x_1, x_2) = e^{(x_1^2 - x_2^2)} - 1 \end{cases}$$

We are interested in point(s) $x^* = (x_1^*, x_2^*)$ where $f(x^*) = 0$



The Newton method (cont.)



The problem has a unique solution, the point where both function f_1 and f_2 are zero
 $\rightsquigarrow (0, 0)^T$

□

The Newton method | Towards optimisation

Our main interest in root-finding is on certain vector-valued functions $f : \mathcal{R}^N \rightarrow \mathcal{R}^N$

- Functions f that are the gradient of twice-differentiable functions $g : \mathcal{R}^N \rightarrow \mathcal{R}$

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

We are interested in point(s) $x^* \in \mathcal{R}^N$ such that $\nabla g(x) = 0$

- Extrema: Minima, maxima, and saddle points of $g(x)$

$$\text{Solve for } \delta x \quad J_f(x^{(k)}) \delta x^{(k)} = -f(x^{(k)})$$

$$\text{Compute } x^{(k+1)} = x^{(k)} + \delta x^{(k)}$$

\rightsquigarrow Function f is gradient $\nabla g(x)$ of g , the Jacobian $J_f(x^{(k)})$ is its Hessian $\nabla^2 g(x^{(k)})$

$$\text{Solve for } \delta x \quad \nabla^2 g(x^{(k)}) \delta x^{(k)} = -\nabla g(x^{(k)})$$

$$\text{Compute } x^{(k+1)} = x^{(k)} + \delta x^{(k)}$$

The Newton method | Towards optimisation (cont.)

$$\text{Solve for } \delta x \quad \nabla^2 g(x^{(k)}) \delta x^{(k)} = -\nabla g(x^{(k)})$$

$$\text{Compute } x^{(k+1)} = x^{(k)} + \delta x^{(k)}$$

We can re-write the Newton method as we derived it, as an explicit iteration scheme

$$\text{Solve for } \delta x \quad \nabla^2 g(x^{(k)}) (x^{(k+1)} - x^{(k)}) = -\nabla g(x^{(k)})$$

$$\text{Compute } x^{(k+1)} = x^{(k)} - \left(\nabla^2 g(x^{(k)})\right)^{-1} \nabla g(x^{(k)})$$

That is, starting from some initial solution (guess) $x^{(0)}$

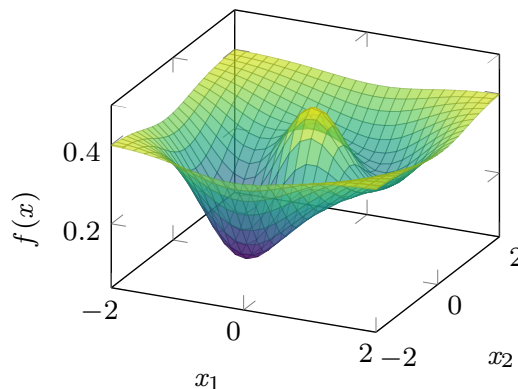
$$x^{(k+1)} = x^{(k)} - \left(\nabla^2 g(x^{(k)})\right)^{-1} \nabla g(x^{(k)}), \quad k = 0, 1, \dots$$

The Newton method | Towards optimisation (cont.)

Example

Consider function $f : \mathcal{R}^2 \rightarrow \mathcal{R}$

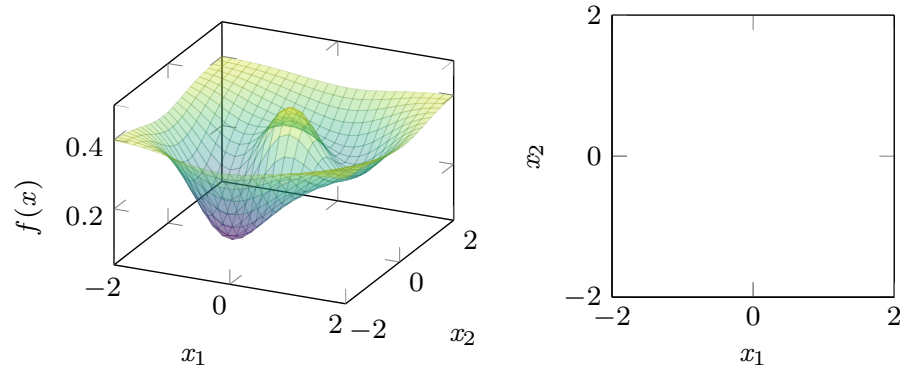
$$f(x) = \frac{2}{5} - \frac{1}{10} (5x_1^2 + 5x_2^2 + 3x_1x_2 - x_1 - 2x_2) e^{-(x_1^2 + x_2^2)}$$



We are interested in those points (x_1, x_2) where $\nabla f(x) = 0$

We can easily identify the points where the gradient vanishes

- Two minima
- A maximum
- A saddle



We can analyse solutions from the Newton's method, from different initial points $x^{(0)}$

Suppose that we let $x^{(0)} = (-0.9, -0.9)$

↪ After 5 iterations the method converges to $x = (-0.63058, -0.70074)$

Suppose that we let $x^{(0)} = (-1.0, -1.0)$

↪ After 400 iterations the stopping criterion is still not fulfilled

Suppose that we let $x^{(0)} = (+0.5, -0.5)$

↪ After 5 iterations the method converges to the saddle point

- $x = (0.80659, -0.54010)$

□

The Newton method | Towards optimisation (cont.)

$$x^{(k+1)} = x^{(k)} - \left(J_f \left(x^{(k)} \right) \right)^{-1} f \left(x^{(k)} \right)$$

$$x^{(k+1)} = x^{(k)} - \left(\nabla^2 g \left(x^{(k)} \right) \right)^{-1} \nabla g \left(x^{(k)} \right)$$

In spite of a simple implementation, the Newton method is demanding for a large N

- ↪ The method requires analytic expressions of the derivatives
- ↪ The method also requires inverting the Jacobian (Hessian)
- ↪ Naive inversion of a $N \times N$ matrix is $\mathcal{O}(N^3)$

For the method to converge, it is also important that $x^{(0)}$ is chosen near enough x^*

The Newton method | Towards optimisation (cont.)

$$x^{(k+1)} = x^{(k)} - \left(J_f(x^{(k)}) \right)^{-1} f(x^{(k)})$$

$$x^{(k+1)} = x^{(k)} - \left(\nabla^2 g(x^{(k)}) \right)^{-1} \nabla g(x^{(k)})$$

Flexibility is achieved by replacing Jacobians (Hessians) by invertible approximations

$$\underbrace{M^{(k)}}_{N \times N}$$

The use of invertible approximations leads the family of **Newton-type methods**

$$x^{(k+1)} = x^{(k)} - \left(\begin{array}{c} \underbrace{M^{(k)}} \\ (J_f(x^{(k)}))^{-1} \end{array} \right)^{-1} f(x^{(k)})$$

$$x^{(k+1)} = x^{(k)} - \left(\begin{array}{c} \underbrace{M^{(k)}} \\ (\nabla^2 g(x^{(k)}))^{-1} \end{array} \right)^{-1} \nabla g(x^{(k)})$$

Newton-type methods

Root-finding with Newton-type methods

Newton-type methods

Consider some function $f \in \mathcal{C}^2(\mathcal{R}^N)$ bounded below, we are interested in its minima

- That is, we are interested in point $x^* \in \mathcal{R}^N$ such that $f(x^*)$ is the smallest
- The minima of f occur at points x where the gradient $\nabla f(x)$ is zero

We can use Newton and Newton-type recursions to find the zeros of function $\nabla f(x)$

- From some initial approximate solution $x^{(0)}$, we have

$$x^{(k+1)} = x^{(k)} + \underbrace{(-1) \left(M^{(k)} \right)^{-1} \nabla f(x^{(k)})}_{d^{(k)}}, \quad k = 0, 1, \dots$$

At iteration steps $k \geq 0$, let $x^{(k+1)}$ be the next point of the sequence

- Point $x^{(k+1)}$ depends on point $x^{(k)}$ and some vector $d^{(k)}$

The vector (direction) $d^{(k)}$ depends on two terms

- ↪ The gradient vector $\nabla f(x^{(k)})$ of f
- ↪ The Hesse matrix $\nabla^2 f(x^{(k)})$ of f
- ↪ (Or, an approximation $M^{(k)}$)

Newton-type methods (cont.)

$$x^{(k+1)} = x^{(k)} + \underbrace{\alpha^{(k)}}_{\alpha^{(k)}} \underbrace{\left(M^{(k)} \right)^{-1} (-\nabla f(x^{(k)}))}_{d^{(k)}}, \quad k = 0, 1, \dots$$

We can also introduce a dependence on some parameter $\alpha_k \in \mathcal{R}_{\geq 0}$, the **step-length**

- In the basic implementation $\alpha^{(k)}$ is constant and equal to positive one
- $\nabla f(x^{(k)})$ gives the direction of max positive growth of f from $x^{(k)}$
- $\nabla^2 f(x^{(k)})$ applies a transformation to the gradient direction f

The negative sign sets the iterates to move downwards

Newton-type methods (cont.)

The algorithmic formulation of a general Newton-type (line-search/descent) method

Let $x^{(0)} \in \mathcal{R}^N$ be an initial approximation of the minimiser

Determine descent direction $d^{(k)} \in \mathcal{R}^N$

Compute step-length $\alpha^{(k)} \in \mathcal{R}_{\geq 0}$

Compute new approximation $x^{(k+1)} = x^{(k)} + \alpha^{(k)} d^{(k)}$

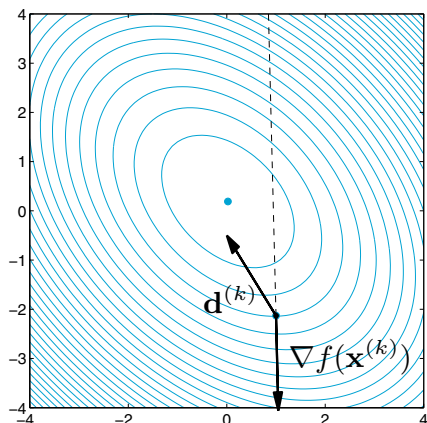
Newton-type methods (cont.)

Because vector $d^{(k)}$ needs be a **descent direction**, it must satisfy certain conditions

$$\begin{cases} d^{(k)} \nabla f(x^{(k)}) < 0 & (\nabla f(x^{(k)}) \neq 0) \\ d^{(k)} = 0 & (\nabla f(x^{(k)}) = 0) \end{cases}$$

$\rightsquigarrow d^{(k)} \nabla f(x^{(k)})$ is the directional derivative of f along $d^{(k)}$

For example, consider a function f (a quadratic form) and its gradient vector at $x^{(k)}$



$$\underbrace{(1)}_{\alpha^{(k)}} \underbrace{\left(M^{(k)} \right)^{-1} (-\nabla g(x^{(k)}))}_{d^{(k)}}$$

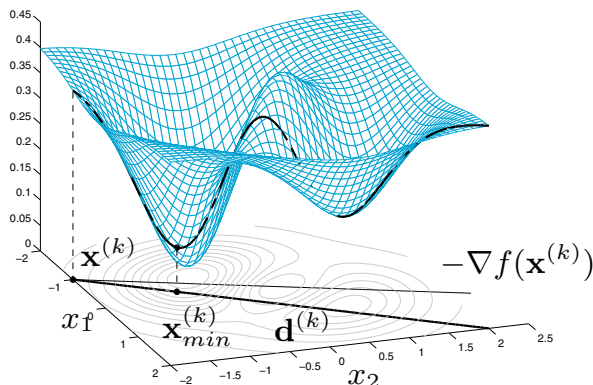
Direction $d^{(k)}$ must be a suitable descent direction

- Parameter $\alpha^{(k)}$ defines the step-size

Newton-type methods (cont.)

The step-size $\alpha^{(k)}$ can be computed by solving a one-dimensional minimisation problem

- The minimisation of the restriction of function $f(x)$ along direction $d^{(k)}$
- The idea is to set $\alpha^{(k)}$ to reach $x_{\min}^{(k)}$, the minimiser along $d^{(k)}$

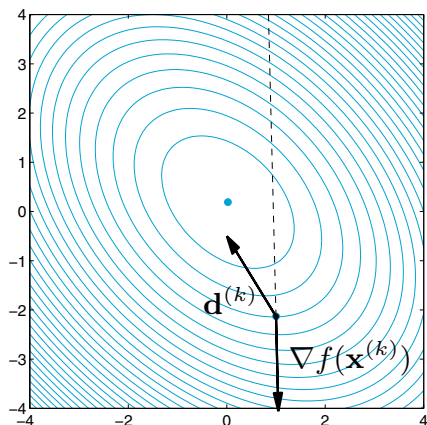


When f or its restriction is not quadratic, the computation of $\alpha^{(k)}$ can be involved

- Certain (Wolfe's) conditions on $\alpha^{(k)}$ must be satisfied before it is accepted

Newton-type methods | Step-lengths

Given a descent direction $d^{(k)}$, we would want the step-length $\alpha^{(k)}$ to be optimal



- Such that function f is smallest along $d^{(k)}$
- (Along the restriction $f(x^{(k)} + \alpha^{(k)}d^{(k)})$)

$$f(x^{(k)} + \alpha^{(k)}d^{(k)}) \approx f(x^{(k)}) + \alpha^{(k)}\nabla^T f(x^{(k)})d^{(k)} + \frac{1}{2}(\alpha^{(k)})^2 d^{(k)}\nabla^2 f(x^{(k)})d^{(k)}$$

By setting to zero the derivative with respect to $\alpha^{(k)}$ of a second-order approximation around $x^{(k)}$ of the restriction of function f along the descent direction $d^{(k)}$, we get

$$\alpha^{(k)} = \frac{-\nabla^T f(x^{(k)})d^{(k)}}{d^{(k)}\nabla^2 f(x^{(k)})d^{(k)}}$$

Newton-type methods | Newton directions

$$x^{(k+1)} = x^{(k)} + \alpha^{(k)} \underbrace{\left(M^{(k)} \right)^{-1} \left(-\nabla f(x^{(k)}) \right)}_{d^{(k)}}, \quad k = 0, 1, \dots$$

Newton's directions, $M^{(k)} = \nabla^2 f(x^{(k)})$

$$d^{(k)} = - \left(\underbrace{\nabla^2 f(x^{(k)})}_{M^{(k)} \succ 0} \right)^{-1} \nabla f(x^{(k)})$$

Newton-type methods | Newton directions (cont.)

$$d^{(k)} = - \left(\underbrace{\nabla^2 f(x^{(k)})}_{M^{(k)}} \right)^{-1} \nabla f(x^{(k)})$$

Consider functions f such that the Hessian matrices $\{\nabla^2 f(x^{(k)})\}$ are positive definite

- Also, suppose that for some $\kappa > 0$ the condition number of $\nabla^2 f(x^{(k)})$

$$\begin{aligned} \mathcal{K}(\nabla^2 f(x^{(k)})) &= \frac{\lambda_{\max}(\nabla^2 f(x^{(k)}))}{\lambda_{\min}(\nabla^2 f(x^{(k)}))} \\ &\leq \kappa \quad (\text{for all } k) \end{aligned}$$

Under these conditions, the sequence $\{x^{(k)}\}$ converges to a minimum x^* of function f

Newton-type methods | Newton directions (cont.)

The positive definiteness of $M^{(k)}$ must be over a large enough neighbourhood of $x^{(k)}$

If $M^{(k)} \succ 0$, then $d^{(k)} = - (M^{(k)})^{-1} \nabla f(x^{(k)})$ is a descent direction

$$d^{(k)} \nabla f(x^{(k)}) < 0 \quad (\nabla f(x^{(k)}) \neq 0)$$

We have,

$$\begin{aligned} d^{(k)} \nabla f(x^{(k)}) &= - \underbrace{\nabla^T f(x^{(k)})}_{>0} \underbrace{M^{(k)}^{-1}}_{>0} \nabla f(x^{(k)}) \\ &< 0 \end{aligned}$$

Note that a descent direction do not necessarily imply a reduction in function value

- The step-length $\alpha^{(k)}$ may lead to $f(x^{(k+1)}) > f(x^{(k)})$
- The step-length can be reduced, to avoid this risk

Newton-type methods | Newton directions (cont.)

$$d^{(k)} = - \left(\underbrace{\nabla^2 f(x^{(k)})}_{M^{(k)}} \right)^{-1} \nabla f(x^{(k)})$$

For Hessians that are not positive definite, direction $d^{(k)}$ may not be a descent direction

$$d^{(k)} \nabla f(x^{(k)}) \geq 0 \quad (\nabla f(x^{(k)}) \neq 0)$$

- (Also Wolfe's conditions on the step-length may lose validity/meaning)

Against this, it is possible to add a diagonal or full matrix $E^{(k)}$ to the Hessian

$$\underbrace{\nabla^2 f(x^{(k)}) + E^{(k)}}_{M^{(k)}} \succ 0$$

Newton-type methods | Quasi-Newton

$$x^{(k+1)} = x^{(k)} + \alpha^{(k)} \underbrace{\left(M^{(k)} \right)^{-1} \left(-\nabla f(x^{(k)}) \right)}_{d^{(k)}}, \quad k = 0, 1, \dots$$

Quasi-Newton directions, $M^{(k)} = \tilde{\nabla}^2 f(x^{(k)})$

$$d^{(k)} = - \left(\underbrace{\tilde{\nabla}^2 f(x^{(k)})}_{M^{(k)} \succ 0} \right)^{-1} \nabla f(x^{(k)})$$

A common approach for constructing approximations of the Hessian matrix is **BFGS**

- The **Broyden**, **Fletcher**, **Galfarb**, and **Shanno**'s method

Newton-type methods | Quasi-Newton (cont.)

Given an initial symmetric matrix $M^{(0)} \succ 0$, the BFGS method recursively computes

$$\begin{aligned} M^{(k+1)} = & M^{(k)} \\ & + \frac{(\nabla f(x^{(k+1)}) - \nabla f(x^{(k)})) (\nabla f(x^{(k+1)}) - \nabla f(x^{(k)}))^T}{(\nabla f(x^{(k+1)}) - \nabla f(x^{(k)}))^T (x^{(k+1)} - x^{(k)})^T} \\ & - \frac{M^{(k)} (x^{(k+1)} - x^{(k)}) (x^{(k+1)} - x^{(k)})^T M^{(k)}}{(x^{(k+1)} - x^{(k)})^T M^{(k)} (x^{(k+1)} - x^{(k)})} \end{aligned}$$

The matrices from rank-one updates, as BFGS, need be symmetric and positive definite

This is guaranteed by the following condition,

$$(x^{(k+1)} - x^{(k)})^T (\nabla f(x^{(k)}) - \nabla f(x^{(k+1)})) > 0$$

From a quadratic approximation of f about $x^{(k)}$

Newton-type methods | Quasi-Newton (cont.)

Given an approximate solution $x^{(0)}$ and a positive definite approximate Hessian $M^{(0)}$

We have the general formulation of the quasi-Newton's method

$$\begin{aligned} \text{Solve for } d^{(k)} \quad & M^{(k)} d^{(k)} = -\nabla f(x^{(k)}) \\ \text{(Compute } \alpha^{(k)} \quad & \text{Verify Wolfe's conditions)} \\ \text{Set } x^{(k+1)} \quad & x^{(k+1)} = x^{(k)} + \alpha^{(k)} d^{(k)} \\ \text{Compute} \quad & x^{(k+1)} - x^{(k)} \\ \text{Compute} \quad & \nabla f(x^{(k+1)}) - \nabla f(x^{(k)}) \\ \text{Compute} \quad & M^{(k+1)} \end{aligned}$$

Newton-type methods | Gradient directions (cont.)

$$x^{(k+1)} = x^{(k)} - \underbrace{\left(M^{(k)}\right)^{-1} \nabla f(x^{(k)})}_{d^{(k)}}, \quad k = 0, 1, \dots$$

Gradient directions (gradient descent, steepest descent, ...), $M^{(k)} = I$

$$d^{(k)} = - \left(\underbrace{I}_{M^{(k)}} \right)^{-1} \nabla f(x^{(k)})$$

This approach is successfully utilised for large-scale optimisation problems

- Where large Hessian matrices are expensive to invert

Newton-type methods | Gradient directions (cont.)

Conjugate-gradient directions

$$d^{(0)} = -\nabla f(x^{(0)})$$

$$d^{(k+1)} = -\nabla f(x^{(k+1)}) + \beta^{(k)} d^{(k)}, \quad k = 0, 1, \dots$$

There exist alternatives for computing parameter $\beta^{(k)}$, some commonly used ones

↪ **Fletcher-Reeves**

$$\beta_{\text{FR}}^{(k)} = \frac{\|\nabla f(x^{(k)})\|^2}{\|\nabla f(x^{(k-1)})\|^2}$$

↪ **Hestenes-Stiefel**

$$\beta_{\text{HS}}^{(k)} = \frac{\nabla f(x^{(k)})^T (\nabla f(x^{(k)}) - \nabla f(x^{(k-1)}))}{d^{(k-1)T} (\nabla f(x^{(k)}) - \nabla f(x^{(k-1)}))}$$

↪ **Polak-Ribière**

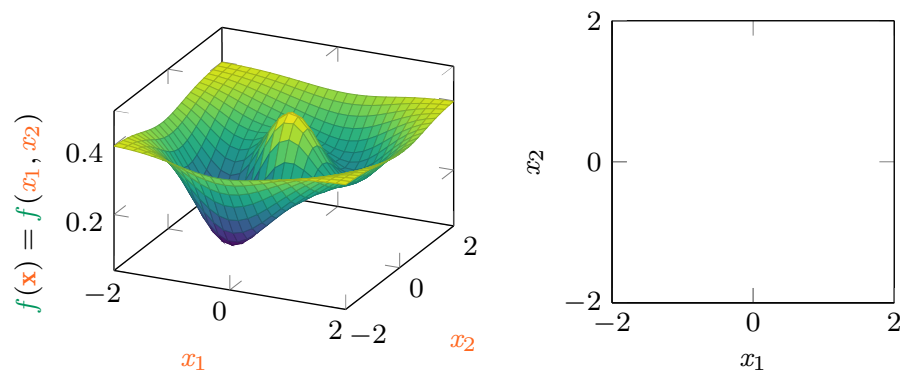
$$\beta_{\text{PR}}^{(k)} = \frac{\nabla f(x^{(k)})^T (\nabla f(x^{(k)}) - \nabla f(x^{(k-1)}))}{\|\nabla f(x^{(k-1)})\|^2}$$

Newton-type methods (cont.)

Example

Consider function $f : \mathcal{R}^2 \rightarrow \mathcal{R}$

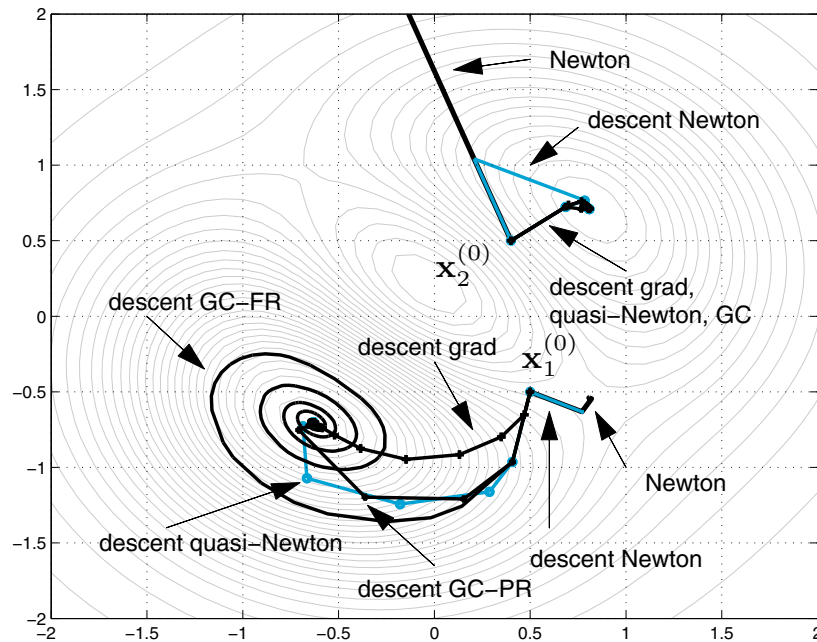
$$f(x) = \frac{2}{5} - \frac{1}{10} (5x_1^2 + 5x_2^2 + 3x_1x_2 - x_1 - 2x_2) e^{-(x_1^2 + x_2^2)}$$



Compare sequences $\{x^{(k)}\}$ with Newton and quasi-Newton direction, from various $x^{(0)}$

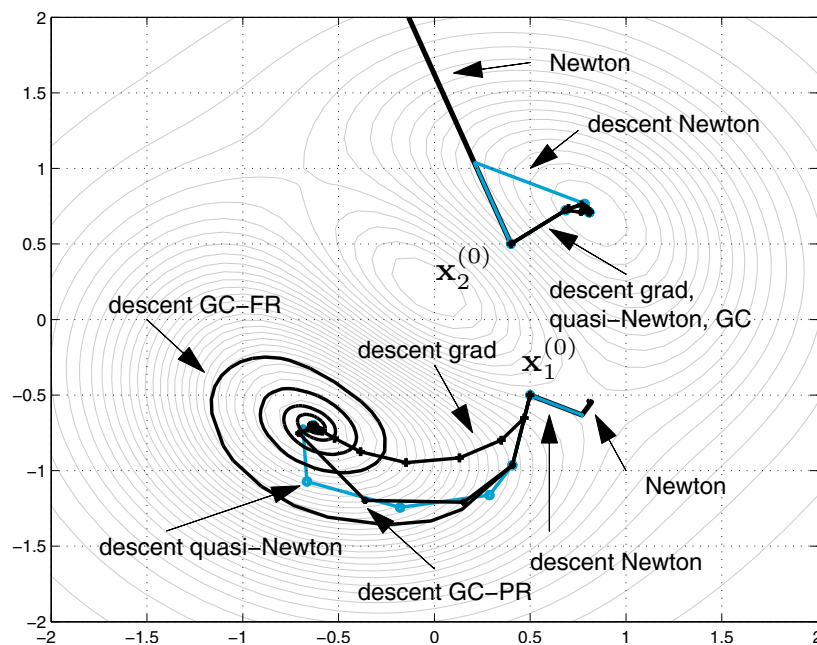
Newton-type methods (cont.)

$x_1^{(0)} = (0.5, -0.5)$, Newton converges (to a saddle) and some inexact methods collapse



Newton-type methods (cont.)

$x_2^{(0)} = (0.4, 0.5)$, Newton diverges but starts well together with some inexact methods



Convergence rates

Root-finding with Newton-type methods

Convergence

Consider the set of nonlinear equations from the vector-valued function $f : \mathcal{R}^N \rightarrow \mathcal{R}^N$

$$f(x) = 0$$

We are interested in solving this system of equations, a root-finding problem

- That is, find x^* such that $f(x^*) = 0$

The exact Newton method, solve

$$f(x^{(k)}) + J_f(x^{(k)})(x^{(k+1)} - x^{(k)}) = 0$$

We get the exact iterates,

$$x^{(k+1)} = x^{(k)} - \left(J_f(x^{(k)}) \right)^{-1} f(x^{(k)})$$

The Newton-type iterates,

$$x^{(k+1)} = x^{(k)} - \left(M^{(k)} \right)^{-1} f(x^{(k)})$$

$M^{(k)}$ must be an invertible and positive definite approximation of the Jacobian $J_f(x^{(k)})$

Convergence (cont.)

Let $x^{(k)}$ be the approximated solution at iteration k and let x^* denote the solution

Consider converging sequences of iterates $\{x^{(k)}\}$,

$$\lim_{k \rightarrow \infty} x^{(k)} = x^*$$

Or, equivalently,

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

We are interested in characterising the rate at which iterates the $x^{(k)}$ converge to x^*

Consider the convergence condition,

$$\frac{\|x^{(k+1)} - x^*\|}{\|x^{(k)} - x^*\|^p} \leq C \quad (\text{for some } C \in (0, 1) \text{ and } k \geq k_0)$$

Order of convergence is denoted by p and C is known as the **convergence factor**

- The condition is defined for an error, which is based on unknown x^*

The necessary condition for convergence is that $x^{(k_0)}$ is chosen sufficiently close to x^*

- Because of this, only **local convergence** properties can be established

Convergence (cont.)

We define the following **(local) convergence rates** for the sequence $\{x^{(k)}\}$ of iterates

- **Q-linear**, for some $C^{(k)} \in (0, 1)$ and for all $k = 0, 1, \dots$

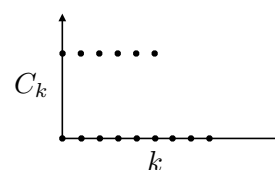
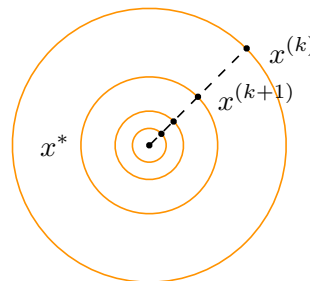
$$\frac{\|x^{(k+1)} - x^*\|}{\|x^{(k)} - x^*\|} \leq C^{(k)}$$

$C^{(k)}$ is the rate of convergence

- It remains constant with k

An equivalent form

$$\limsup_{k \rightarrow \infty} \frac{\|x^{(k+1)} - x^*\|}{\|x^{(k)} - x^*\|} < 1$$



Linear contraction rates characterise an exponential decay of the approximation error

- An exponential decay (or growth) is not necessarily rapid, it depends on C

Convergence (cont.)

- **Q-superlinear**, for some stable sequence $C^{(k)} \rightarrow 0$

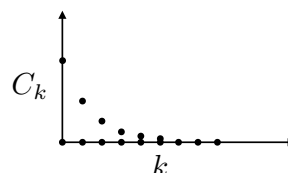
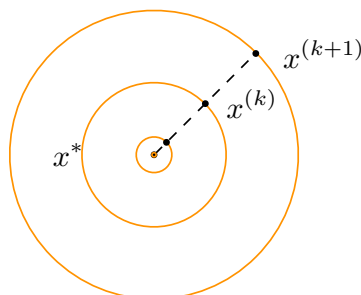
$$\frac{\|x^{(k+1)} - x^*\|}{\|x^{(k)} - x^*\|} \leq C^{(k)}$$

$C^{(k)}$ is the rate of convergence

- It shrinks with k

In the limit form,

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



The rate of the exponential decay is not constant, but decays with the iteration count

- It is equivalent to an always increasing linear contraction rate

Convergence (cont.)

- **Q-quadratic**, for some $C^{(k)} < \infty$ and for all $k = 0, 1, \dots$

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

Rearranging terms, we have

$$\frac{\|x^{(k+1)} - x^*\|}{\|x^{(k)} - x^*\|} \leq \underbrace{C^{(k)} \|x^{(k)} - x^*\|}_{C^{(k)}(x^{(k)})}$$

$C^{(k)}(x^{(k)})$ is a local rate of convergence

- It shrinks with k and $x^{(k)}$

