

6. Unconstrained optimization

- unconstrained minimization
- descent methods
- the gradient descent method
- Newton's method

Unconstrained minimization

$$\text{minimize } f(\mathbf{x}) \quad (6.1)$$

- $\mathbf{x} = (x_1, \dots, x_n)$ is the *variable*
- $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is the *objective function*

Solution: the point $\mathbf{x}^* = (x_1^*, \dots, x_n^*)$ is a *minimizer (minimum point)* of f or *solution* of (6.1) if

$$f(\mathbf{x}^*) \leq f(\mathbf{x})$$

for all n -vectors \mathbf{x}

Optimal value and local minimizer

Optimal value: the *optimal value* of the minimization problem is the greatest p such that $p \leq f(\mathbf{x})$, denoted by $\min f(\mathbf{x})$

- if \mathbf{x}^* is a minimizer of f , then $f(\mathbf{x}^*) = \min f(\mathbf{x})$ and we say that the optimal value is attained at \mathbf{x}^*
- if $\min f(\mathbf{x}) = -\infty$, then we say that the function is unbounded below
- the optimal value is unique even though there could be multiple solutions

Local minimizer

- the minimizer \mathbf{x}^* of f is also called a *global minimizer* of f
- a vector \mathbf{x}^o is a *local minimizer* or *local minimum point* if there exists a scalar $r > 0$ such that $f(\mathbf{x}^o) \leq f(\mathbf{x})$ for all $\|\mathbf{x} - \mathbf{x}^o\| \leq r$; it is a *strict local minimizer* if $f(\mathbf{x}^o) < f(\mathbf{x})$

First-order necessary condition

if the n -vector \mathbf{x}^o is a local minimizer of $f : \mathbb{R}^n \rightarrow \mathbb{R}$, then

$$\nabla f(\mathbf{x}^o) = \mathbf{0} \quad \left(\frac{\partial f}{\partial x_i}(\mathbf{x}^o) = 0, \quad i = 1, \dots, n \right)$$

- this condition is *necessary* but not sufficient; points that satisfy $\nabla f(\hat{\mathbf{x}}) = \mathbf{0}$ can be minimizers, maximizers, or neither (saddle points)
- points that satisfies $\nabla f(\hat{\mathbf{x}}) = \mathbf{0}$ are called *stationary points* or *critical points*
- in general, to find a global minimizer, we need to check whether the solutions of $\nabla f(\hat{\mathbf{x}}) = \mathbf{0}$ are in fact global minimizers
- it often is very difficult to solve the set of nonlinear equations and numerical algorithms are often used for finding stationary points

Example 6.1

find the stationary points of

$$f(\mathbf{x}) = x_1^3 - x_1^2x_2 + 2x_2^2$$

setting the gradient (partial derivatives) to zero, we get the FONC:

$$\frac{\partial f}{\partial x_1} = 3x_1^2 - 2x_1x_2 = 0$$

$$\frac{\partial f}{\partial x_2} = -x_1^2 + 4x_2 = 0$$

solving, we get two stationary points: $(0, 0)$ and $(6, 9)$

Second-order condition

Necessary condition: if \mathbf{x}^o is a local minimizer, then $\nabla f(\mathbf{x}^o) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}^o) \geq 0$

Sufficient condition: if $\nabla f(\mathbf{x}^o) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}^o) > 0$, then \mathbf{x}^o is a strict local minimizer

Necessary and sufficient condition: if $\nabla^2 f(\mathbf{x}) \geq 0$ for all \mathbf{x} (' f is convex'), then \mathbf{x}^* is global minimizer if and only if $\nabla f(\mathbf{x}^*) = \mathbf{0}$

- for single variable, the Hessian is just the second derivative $f''(x)$
- we can find maximizers by finding minimizers of $-f$

Example 6.2

find the stationary points of $f(\mathbf{x})$ and, if possible, determine whether they are local or global minimizers

a) $f(\mathbf{x}) = x_1^3 - x_1^2x_2 + 2x_2^2$

b) $f(\mathbf{x}) = \frac{1}{2}x_1^2 + x_1x_2 + 2x_2^2 - 4x_1 - 4x_2 - x_2^3$

c) $f(\mathbf{x}) = x_1^2 - x_1x_2 + x_2^2 - 3x_2$

a) for $f(\mathbf{x}) = x_1^3 - x_1^2x_2 + 2x_2^2$, setting the gradient, we find that the stationary points are $(0, 0)$ and $(6, 9)$ (see page 6.5); the Hessian is

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} 6x_1 - 2x_2 & -2x_1 \\ -2x_1 & 4 \end{bmatrix}$$

hence,

$$\nabla^2 f(0, 0) = \begin{bmatrix} 0 & 0 \\ 0 & 4 \end{bmatrix}, \quad \nabla^2 f(6, 9) = \begin{bmatrix} 18 & -12 \\ -12 & 4 \end{bmatrix}$$

- since $\nabla^2 f(0, 0)$ is only positive semidefinite, it is still unclear whether $(0, 0)$ is a local minimizer
- since the matrix $\nabla^2 f(6, 9)$ is indefinite, the point $(6, 9)$ is not a local minimizer/maximizer
- since $f(\epsilon, 0) > 0$ for any $\epsilon > 0$ and $f(\epsilon, 0) < 0$ for any $\epsilon < 0$, we conclude that the point $(0, 0)$ is not a local minimizer/maximizer

b) for $f(\mathbf{x}) = \frac{1}{2}x_1^2 + x_1x_2 + 2x_2^2 - 4x_1 - 4x_2 - x_2^3$, the FONC is

$$\nabla f(\mathbf{x}) = \begin{bmatrix} x_1 + x_2 - 4 \\ x_1 + 4x_2 - 4 - 3x_2^2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

solving, we get the stationary points $(4, 0)$ and $(3, 1)$; the Hessian is

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} 1 & 1 \\ 1 & 4 - 6x_2 \end{bmatrix}$$

thus,

$$\nabla^2 f(4, 0) = \begin{bmatrix} 1 & 1 \\ 1 & 4 \end{bmatrix}, \quad \nabla^2 f(3, 1) = \begin{bmatrix} 1 & 1 \\ 1 & -2 \end{bmatrix}$$

- since $\nabla^2 f(4, 0) > 0$ and $\nabla^2 f(3, 1)$ is indefinite, $\hat{\mathbf{x}} = (4, 0)$ is a local minimizer and $(3, 1)$ is not a minimizer/maximizer
- note that the point $\hat{\mathbf{x}} = (4, 0)$ is not a global minimizer since $f(0, x_2) \rightarrow -\infty$ as $x_2 \rightarrow \infty$

c) for $f(\mathbf{x}) = x_1^2 - x_1x_2 + x_2^2 - 3x_2$, the FONC is

$$\nabla f(\mathbf{x}) = \begin{bmatrix} 2x_1 - x_2 \\ -x_1 + 2x_2 - 3 \end{bmatrix} = \mathbf{0}$$

or

$$\begin{aligned} 2x_1 - x_2 &= 0 \\ -x_1 + 2x_2 &= 3 \end{aligned}$$

these two equations have a unique solution $\hat{x}_1 = 1, \hat{x}_2 = 2$, which is a candidate for a global minimizer; since the Hessian

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

is positive definite, the point $\hat{\mathbf{x}} = (1, 2)$ is a global minimizer

Example 6.3 (quadratic functions)

suppose we want to minimize $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T Q \mathbf{x} + \mathbf{r}^T \mathbf{x} + c$ where Q is an $n \times n$ symmetric matrix; the FONC is

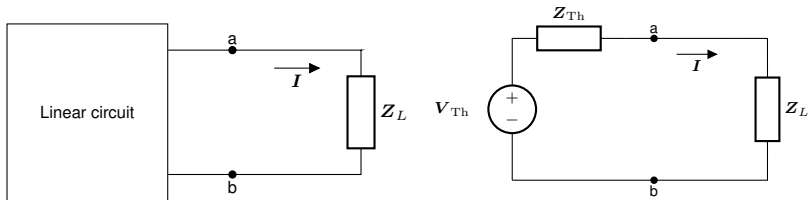
$$\nabla f(\mathbf{x}) = Q\mathbf{x} + \mathbf{r} = \mathbf{0}$$

the Hessian is $\nabla^2 f(\mathbf{x}) = Q$

- if $Q \geq 0$, then \mathbf{x}^* is a global minimizer iff $Q\mathbf{x}^* + \mathbf{r} = \mathbf{0}$
 - if $\mathbf{r} \notin \text{range}(Q)$, then there is no solution and f is unbounded below
 - if $Q > 0$, then there is a unique minimizer $\mathbf{x}^* = -Q^{-1}\mathbf{r}$
 - if Q is singular but $\mathbf{r} \in \text{range}(Q)$, then there exists multiple solutions
- if Q is indefinite, then any point satisfying the FONC is a saddle-point (not a minimizer/maximizer)
- if Q is invertible, then there is a unique stationary point:

$$\hat{\mathbf{x}} = -Q^{-1}\mathbf{r}$$

Example 6.4 (maximum power transfer)



- V_{Th} is the Thevenin voltage, $Z_{Th} = R_{Th} + jX_{Th}$ ($j = \sqrt{-1}$ is the Thevenin impedance, $Z_L = R_L + jX_L$ is the impedance of the load)
- from circuit analysis, the average power delivered to the load is

$$P = |I|^2 R_L, \quad I = \frac{V_{Th}}{R_{Th} + R_L + j(X_{Th} + X_L)}.$$

we want to find the load impedance (*i.e.*, R_L and X_L) such that average power delivered to the load P is maximized; (suppose that $V_{Th} = 1$ and $R_{Th} > 0$)

we can maximize the power by solving

$$\text{maximize } f(\mathbf{x}) = \frac{x_1}{(R_{\text{Th}} + x_1)^2 + (X_{\text{Th}} + x_2)^2}$$

setting the gradient (partial derivatives) to zero:

$$\nabla_{x_1} f(\mathbf{x}) = \frac{\partial f}{\partial x_1} = \frac{(R_{\text{Th}} + x_1)^2 + (X_{\text{Th}} + x_2)^2 - 2x_1(R_{\text{Th}} + x_1)}{[(R_{\text{Th}} + x_1)^2 + (X_{\text{Th}} + x_2)^2]^2} = 0$$

$$\nabla_{x_2} f(\mathbf{x}) = \frac{\partial f}{\partial x_2} = \frac{-2x_1(X_{\text{Th}} + x_2)}{[(R_{\text{Th}} + x_1)^2 + (X_{\text{Th}} + x_2)^2]^2} = 0$$

from the second equation, we have $x_1 = 0$ or $x_2 = -X_{\text{Th}}$; note that $x_1 = 0$ does not satisfy the first condition; using $x_2 = -X_{\text{Th}}$ into the second condition and simplifying, we get

$$(R_{\text{Th}} + x_1)^2 - 2x_1(R_{\text{Th}} + x_1) = 0 \iff x_1 = R_{\text{Th}}$$

hence, the stationary point is

$$\mathbf{x} = (R_{\text{Th}}, -X_{\text{Th}})$$

we now check the second-order conditions; to simplify derivation of the Hessian, we let

$$f(\mathbf{x}) = g(A\mathbf{x} + \mathbf{b})$$

where

$$g(y_1, y_2, y_3) = \frac{y_1}{y_2^2 + y_3^2}, \quad A = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ R_{\text{Th}} \\ X_{\text{Th}} \end{bmatrix}$$

the Hessian of f is $A^T \nabla^2 g(A\mathbf{x} + \mathbf{b}) A$; thus, we need to find the Hessian of g ; the gradient of g is

$$\nabla g(\mathbf{y}) = \begin{bmatrix} \frac{1}{y_2^2 + y_3^2} \\ \frac{-2y_1 y_2}{(y_2^2 + y_3^2)^2} \\ \frac{-2y_1 y_3}{(y_2^2 + y_3^2)^2} \end{bmatrix}$$

the Hessian of g is

$$\nabla^2 g(\mathbf{y}) = \begin{bmatrix} 0 & \frac{-2y_2}{(y_2^2+y_3^2)^2} & \frac{-2y_3}{(y_2^2+y_3^2)^2} \\ \frac{-2y_2}{(y_2^2+y_3^2)^2} & \frac{-2y_1(y_2^2+y_3^2)+8y_1y_2^2}{(y_2^2+y_3^2)^3} & \frac{8y_1y_2y_3}{(y_2^2+y_3^2)^3} \\ \frac{-2y_3}{(y_2^2+y_3^2)^2} & \frac{8y_1y_2y_3}{(y_2^2+y_3^2)^3} & \frac{-2y_1(y_2^2+y_3^2)+8y_1y_3^2}{(y_2^2+y_3^2)^3} \end{bmatrix}$$

$$= \frac{2}{(y_2^2 + y_3^2)^2} \begin{bmatrix} 0 & -y_2 & -y_3 \\ -y_2 & \frac{-y_1(y_2^2+y_3^2)+4y_1y_2^2}{(y_2^2+y_3^2)} & \frac{4y_1y_2y_3}{(y_2^2+y_3^2)} \\ -y_3 & \frac{4y_1y_2y_3}{(y_2^2+y_3^2)} & \frac{-y_1(y_2^2+y_3^2)+4y_1y_3^2}{(y_2^2+y_3^2)} \end{bmatrix}$$

note that at $\mathbf{x} = (R_{\text{Th}}, -X_{\text{Th}})$

$$A\mathbf{x} + \mathbf{b} = \begin{bmatrix} R_{\text{Th}} \\ 2R_{\text{Th}} \\ 0 \end{bmatrix}$$

hence, at $\mathbf{x} = (R_{\text{Th}}, -X_{\text{Th}})$, we have

$$\begin{aligned}\nabla^2 g(A\mathbf{x} + \mathbf{b}) &= \frac{2}{(2R_{\text{Th}})^4} \begin{bmatrix} 0 & -2R_{\text{Th}} & 0 \\ -2R_{\text{Th}} & 3R_{\text{Th}} & 0 \\ 0 & 0 & -R_{\text{Th}} \end{bmatrix} \\ &= \frac{1}{(2R_{\text{Th}})^3} \begin{bmatrix} 0 & -2 & 0 \\ -2 & 3 & 0 \\ 0 & 0 & -1 \end{bmatrix}\end{aligned}$$

the Hessian of f at $\mathbf{x} = (R_{\text{Th}}, -X_{\text{Th}})$ is

$$\begin{aligned}\nabla^2 f(\mathbf{x}) &= A^T \nabla^2 g(A\mathbf{x} + \mathbf{b}) A \\ &= \frac{1}{(2R_{\text{Th}})^3} \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & -2 & 0 \\ -2 & 3 & 0 \\ 0 & 0 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \\ &= \frac{1}{(2R_{\text{Th}})^3} \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix}\end{aligned}$$

since $R_{\text{Th}} > 0$, the above Hessian is negative definite, the point $\mathbf{x} = (R_{\text{Th}}, -X_{\text{Th}})$ is a local maximum; because it is the only point stationary point, it is a global maximum

Outline

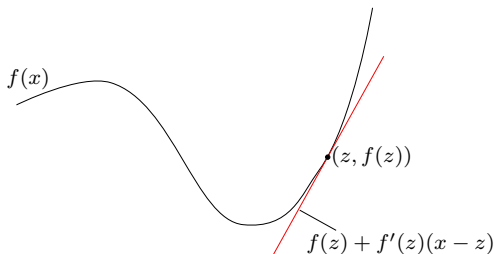
- unconstrained minimization
- **descent methods**
- the gradient descent method
- Newton's method

First-order approximation

First-order (Taylor) approximation of $f(x)$ around z :

$$\begin{aligned}\hat{f}(x) &= f(z) + f'(z)(x - z) && (f : \mathbb{R} \rightarrow \mathbb{R}) \\ \hat{f}(x) &= f(z) + \nabla f(z)^T(x - z) && (f : \mathbb{R}^n \rightarrow \mathbb{R}) \\ \hat{f}(x) &= f(z) + Df(z)(x - z) && (f : \mathbb{R}^n \rightarrow \mathbb{R}^m)\end{aligned}\tag{6.2}$$

Geometrical interpretation



approximation is good if x is close to z and bad otherwise

Descent direction

Descent direction: a vector $\mathbf{d} \in \mathbb{R}^n$ is called a *descent direction* for f if

$$f(\mathbf{x} + \alpha\mathbf{d}) < f(\mathbf{x})$$

for sufficiently small $\alpha > 0$

Directional derivative: the *directional derivative* of f at \mathbf{x} in the direction \mathbf{d} is

$$f'(\mathbf{x}; \mathbf{d}) = \lim_{\alpha \rightarrow 0} \frac{f(\mathbf{x} + \alpha\mathbf{d}) - f(\mathbf{x})}{\alpha} = \nabla f(\mathbf{x})^T \mathbf{d} \quad (6.3)$$

- directional derivative $\nabla f(\mathbf{x})^T \mathbf{d}$ gives an approximate rate of change (increase) of f in the direction \mathbf{d} at the point \mathbf{x}
- a vector $\mathbf{d} \in \mathbb{R}^n$ is a descent method if

$$f'(\mathbf{x}; \mathbf{d}) = \nabla f(\mathbf{x})^T \mathbf{d} < 0$$

Descent methods

Algorithm General descent method

choose a starting point $\mathbf{x}^{(0)}$, a solution tolerance $\epsilon > 0$, and a stopping criteria

repeat for $k \geq 1$

- (a) determine a decent direction $\mathbf{d}^{(k)}$
 - (b) choose a stepsize α_k
 - (c) update $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$ if stopping criteria is satisfied, then stop and output $\mathbf{x}^{(k+1)}$
-

- α_k is called a *learning rate* or *stepsize*
- the stepsize determines the rate that $\mathbf{x}^{(k+1)}$ changes from $\mathbf{x}^{(k)}$ in the descent direction $\mathbf{d}^{(k)}$

Outline

- unconstrained minimization
- descent methods
- **the gradient descent method**
- Newton's method

Negative gradient direction

the directional derivative of f at \mathbf{x} in the direction $\mathbf{d} = -\nabla f(\mathbf{x})$ is

$$\mathbf{d}^T \nabla f(\mathbf{x}) = -\|\nabla f(\mathbf{x})\|^2 < 0$$

for any \mathbf{x} with $\nabla f(\mathbf{x}) \neq \mathbf{0}$; thus, $-\nabla f(\mathbf{x})$ is a descent direction

- suppose $\|\mathbf{d}\| = 1$, then by Cauchy-Schwarz, we have

$$-\|\nabla f(\mathbf{x})\| \leq \nabla f(\mathbf{x})^T \mathbf{d}$$

and equality holds only if $\mathbf{d} = \nabla f(\mathbf{x}) / \|\nabla f(\mathbf{x})\|$

- hence, $-\nabla f(\mathbf{x})$ point in the *steepest descent* (maximum rate of decrease) direction at \mathbf{x}
- if we set $\mathbf{d}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$ in the general descent method, we get the *gradient descent method* or *gradient descent algorithm*

Gradient descent method

Algorithm Gradient descent algorithm

given a starting point $\mathbf{x}^{(0)}$ and a solution tolerance $\epsilon > 0$

repeat for $k \geq 1$

1. choose a stepsize α_k

2. update

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$$

if $\|\nabla f(\mathbf{x}^{(k+1)})\| \leq \epsilon$ stop and output $\mathbf{x}^{(k+1)}$

- for α_k small enough, the algorithm is a descent method
- when α_k is large enough, the algorithm may not be a descent method and often does not work

Determining the stepsize

(suppose $\mathbf{d}^{(k)}$ is any descent direction)

Constant stepsize: set $\alpha_k = \alpha$ for all k

Exact line search

$$\alpha_k = \operatorname{argmin}_{\alpha \geq 0} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$$

- it is not always possible to actually find the exact minimizer α
- called *the method of steepest descent* if $\mathbf{d}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$

Backtracking line search: choose $\beta \in (0, 1)$, and $\gamma \in (0, 1)$ and start with an initial guess α_k (e.g., $\alpha_k = 1$), set $\alpha_k := \beta \alpha_k$ until

$$f(\mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}) - f(\mathbf{x}^{(k)}) < \gamma \alpha_k \nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)}$$

this method is a compromise between the above two methods

Stopping criteria

1. $|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})| < \epsilon$
2. $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| < \epsilon$
3. $|f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)})|/|f(\mathbf{x}^{(k)})| < \epsilon$
4. $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|/|\mathbf{x}^{(k)}| < \epsilon$
5. $\|\nabla f(\mathbf{x}^{(k)})\| < \epsilon$

- the above conditions do not necessarily imply that $\mathbf{x}^{(k)}$ is a good solution since it can be a local minimizer/maximizer or a saddle-point (unless f is convex)
- it is common to run the algorithm from different starting points and choose the best solution of these multiple runs

Example 6.5

$$f(x_1, x_2, x_3) = (x_1 - 4)^4 + (x_2 - 3)^2 + 4(x_3 + 5)^4$$

the gradient of this function is

$$\nabla f(\mathbf{x}) = \begin{bmatrix} 4(x_1 - 4)^3 \\ 2(x_2 - 3) \\ 16(x_3 + 5)^3 \end{bmatrix}$$

let the initial point be $\mathbf{x}^{(0)} = (4, 2, -1)$; applying one iteration of the gradient descent with $\alpha = 0.002$ gives

$$\mathbf{x}^{(1)} = \begin{bmatrix} 4 \\ 2 \\ -1 \end{bmatrix} - 0.002 \begin{bmatrix} 4(4 - 4)^3 \\ 2(2 - 3) \\ 16(-1 + 5)^3 \end{bmatrix} = \begin{bmatrix} 4.000 \\ 2.004 \\ -3.048 \end{bmatrix}$$

notice that

$$59.06 = f(4, 2.004, -3.048) < f(4, 2, -1) = 1025$$

this shows that $\alpha = 0.002$ is a good choice

if we use exact line search, then

$$\begin{aligned}\alpha_0 &= \operatorname{argmin}_{\alpha > 0} f(\mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})) \\ &= \operatorname{argmin}_{\alpha > 0} (0 + (2 + 2\alpha - 3)^2 + 4(-1 - 1024\alpha + 5)^4) = 3.967 \times 10^{-3}\end{aligned}$$

hence,

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \alpha_0 \nabla f(\mathbf{x}^{(0)}) = (4.000, 2.008, -5.062)$$

Example 6.6

$$f(x_1, x_2) = \frac{x_1^2}{5} + x_2^2$$

- the gradient is $\nabla f(\mathbf{x}) = (\frac{2}{5}x_1, 2x_2)$
- we have

$$f(\mathbf{z} - \alpha \nabla f(\mathbf{z})) = \frac{1}{5}(z_1 - \frac{2}{5}\alpha z_1)^2 + (z_2 - 2\alpha z_2)^2$$

- finding the the stepsize at iteration k in the method of steepest descent requires solving

$$\begin{aligned}\alpha &= \operatorname{argmin}_{\alpha > 0} f(\mathbf{z} - \alpha \nabla f(\mathbf{z})) \\ &= \operatorname{argmin}_{\alpha > 0} \left(\frac{1}{5}(z_1 - \frac{2}{5}\alpha z_1)^2 + (z_2 - 2\alpha z_2)^2 \right)\end{aligned}$$

- setting the derivative with respect to α to zero, we get

$$-\frac{4}{25}z_1(z_1 - \frac{2}{5}\alpha z_1) - 4z_2(z_2 - 2\alpha z_2) = 0$$

- solving for α , gives

$$\alpha = \frac{\frac{4}{25}z_1^2 + 4z_2^2}{\frac{8}{125}z_1^2 + 8z_2^2} > 0$$

- hence, the method of steepest descent is

$$\begin{bmatrix} x_1^{(k+1)} \\ x_2^{(k+1)} \end{bmatrix} = \begin{bmatrix} x_1^{(k)} \\ x_2^{(k)} \end{bmatrix} - \frac{\frac{4}{25}(x_1^{(k)})^2 + 4(x_2^{(k)})^2}{\frac{8}{125}(x_1^{(k)})^2 + 8(x_2^{(k)})^2}} \begin{bmatrix} \frac{2}{5}x_1^{(k)} \\ 2x_2^{(k)} \end{bmatrix}$$

Exact line search for quadratic functions

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

where Q is positive definite; gradient method with exact line search requires solving:

$$\alpha_k = \underset{\alpha > 0}{\operatorname{argmin}} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$$

where $\mathbf{d}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$

Update form

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{\|\nabla f(\mathbf{x}^{(k)})\|^2}{\nabla f(\mathbf{x}^{(k)})^T Q \nabla f(\mathbf{x}^{(k)})} \nabla f(\mathbf{x}^{(k)})$$

where $\nabla f(\mathbf{x}) = Q\mathbf{x} - \mathbf{b}$

Derivation

- letting $\mathbf{d} = \mathbf{d}^{(k)}$ and using the chain rule, we have

$$\begin{aligned}g'(\alpha) &= \mathbf{d}^T \nabla f(\mathbf{x}^{(k)} + \alpha \mathbf{d}) \\ &= \mathbf{d}^T (Q(\mathbf{x}^{(k)} + \alpha \mathbf{d}) - \mathbf{b}) \\ &= \alpha \mathbf{d}^T Q \mathbf{d} + \mathbf{d}^T (Q \mathbf{x}^{(k)} - \mathbf{b}) \\ &= \alpha \mathbf{d}^T Q \mathbf{d} + \mathbf{d}^T \mathbf{d}\end{aligned}$$

- setting to zero and solving for α , we get

$$\alpha_k = \frac{\mathbf{d}^T \mathbf{d}}{\mathbf{d}^T Q \mathbf{d}}$$

where $\mathbf{d} = -\nabla f(\mathbf{x}^{(k)}) = -(Q \mathbf{x}^{(k)} - \mathbf{b})$

Convergence discussion

under certain mild assumptions, the iterates $\{\mathbf{x}^{(k)}\}$ of the gradient algorithm can be shown to converge to a stationary point, *i.e.*,

$$\lim_{k \rightarrow \infty} \nabla f(\mathbf{x}^{(k)}) = \mathbf{0}$$

- if f is convex (e.g., $\nabla^2 f(\mathbf{x}) \geq 0$ for all \mathbf{x}), then the iterates $\{\mathbf{x}^{(k)}\}$ of gradient algorithm converges to a global minimizer
- the rate of convergence is sublinear (slow) in general and linear if $\mu I \leq \nabla^2 f(\mathbf{x})$ for all \mathbf{x} and some constant $\mu > 0$

Outline

- unconstrained minimization
- descent methods
- the gradient descent method
- **Newton's method**

Newton's method

consider n nonlinear equation in n variables

$$r_1(\mathbf{x}) = 0, \quad r_2(\mathbf{x}) = 0, \quad \dots, \quad r_n(\mathbf{x}) = 0$$

where $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$; we let $r(\mathbf{x}) = (r_1(\mathbf{x}), \dots, r_n(\mathbf{x}))$

Newton's method: choose $\mathbf{x}^{(0)}$ and repeat for $k \geq 0$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - Dr(\mathbf{x}^{(k)})^{-1}r(\mathbf{x}^{(k)})$$

assumes $Dr(\mathbf{x}^{(k)})$ exists and nonsingular

Unconstrained optimization: if $r(\mathbf{x}) = \nabla f(\mathbf{x})$, we get

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

Interpretation of Newton update

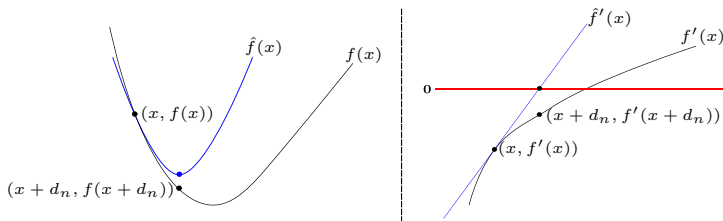
$$\mathbf{x} = \mathbf{x}^{(k)} - \nabla^2 f(\mathbf{x}^{(k)})^{-1} \nabla f(\mathbf{x}^{(k)})$$

(I) minimizing the quadratic approximation of f around $\mathbf{x}^{(k)}$:

$$\hat{f}(\mathbf{x}) = f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^T (\mathbf{x} - \mathbf{x}^{(k)}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(k)})^T \nabla^2 f(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)})$$

(II) solve approximate optimality condition around $\mathbf{x}^{(k)}$:

$$\widehat{\nabla} f(\mathbf{x}) = \nabla f(\mathbf{x}^{(k)}) + \nabla^2 f(\mathbf{x}^{(k)}) (\mathbf{x} - \mathbf{x}^{(k)}) = \mathbf{0}$$



Damped Newton's method

Algorithm Damped Newton method

given a starting point $\mathbf{x}^{(0)}$, a solution tolerance $\epsilon > 0$

repeat for $k \geq 1$

1. choose a stepsize α_k
2. update

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

if $\|\nabla f(\mathbf{x}^{(k+1)})\| \leq \epsilon$ stop and output $\mathbf{x}^{(k+1)}$

- assumes $\nabla^2 f(\mathbf{x})$ exists and is invertible
- $\mathbf{d}_n = -(\nabla^2 f(\mathbf{x}^{(k)}))^{-1} \nabla f(\mathbf{x}^{(k)})$ is called *Newton step* at $\mathbf{x}^{(k)}$
- similar stepsize selection and stopping criteria as before can be used

single-variable update

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

Convergence discussion

- under certain assumptions one can prove local quadratic rate of convergence (fast), *i.e.*, near the optimal solution the error $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|$ (where \mathbf{x}^* is an optimal solution) satisfies

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

for some positive $c > 0$

- if $\nabla^2 f(\mathbf{x}) > 0$ then $\mathbf{d}_n = -(\nabla^2 f(\mathbf{x}^{(k)}))^{-1} \nabla f(\mathbf{x}^{(k)})$ is a descent direction and quadratic convergence to the global minimizer is guaranteed under certain conditions
- does not work well if $\nabla^2 f(\mathbf{x})$ is not positive-definite since it is not a descent method in this case
- can use hybrid gradient-Newton method by setting $\mathbf{d}^{(k)} = -\nabla^2 f(\mathbf{x}^{(k)})^{-1} \nabla f(\mathbf{x}^{(k)})$ if $\nabla^2 f(\mathbf{x}^{(k)})$ is positive-definite and $\mathbf{d}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$ otherwise

Numerical example I

$$f(\mathbf{x}) = e^{x_1+x_2-1} + e^{x_1-x_2-1} + e^{-x_1-1}$$

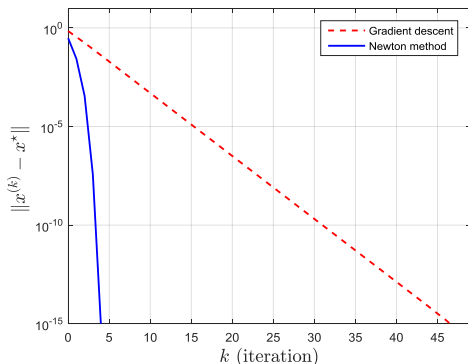
the gradient and Hessian are

$$\nabla f(\mathbf{x}) = \begin{bmatrix} e^{x_1+x_2-1} + e^{x_1-x_2-1} - e^{-x_1-1} \\ e^{x_1+x_2-1} - e^{x_1-x_2-1} \end{bmatrix}$$

and

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} e^{x_1+x_2-1} + e^{x_1-x_2-1} + e^{-x_1-1} & e^{x_1+x_2-1} - e^{x_1-x_2-1} \\ e^{x_1+x_2-1} - e^{x_1-x_2-1} & e^{x_1+x_2-1} + e^{x_1-x_2-1} \end{bmatrix}$$

we apply gradient descent and Newton method with initial starting point $x^{(0)} = (-1, 1)$ and step-size $\alpha = 1$



- both algorithms converge to $x^* = (-0.34657, 0)$
- Newton method is much faster since it uses second-order information

Matlab implementation

```
g=@(x) [exp(x(1)+x(2)-1)+exp(x(1)-x(2)-1)-exp(-x(1)-1);...
exp(x(1)+x(2)-1)-exp(x(1)-x(2)-1)]; % gradient
hess=@(x) [exp(x(1)+x(2)-1)+exp(x(1)-x(2)-1)+exp(-x(1)-1) ...
exp(x(1)+x(2)-1)-exp(x(1)-x(2)-1);...
exp(x(1)+x(2)-1)-exp(x(1)-x(2)-1) ...
exp(x(1)+x(2)-1)+exp(x(1)-x(2)-1)] % hessain
%% Newton and GD iterations
x = [-1; 1];%GD initialization
xn = [-1; 1];%Newton initialization
alpha=1; %step-size
for k=1:50
%%%Gradient descent update%%%
grad=g(x);
if (norm(grad) < 1e-16), break; end;
x = x - alpha*grad;
%%%Newton update%%%
dn=-hess(xn)\g(xn);
xn = xn + alpha*dn;
end
```


Alternative way to construct gradient and Hessian

$$f(\mathbf{x}) = e^{x_1+x_2-1} + e^{x_1-x_2-1} + e^{-x_1-1}$$

we can write f as $f(\mathbf{x}) = g(A\mathbf{x} + \mathbf{b})$, where $g(\mathbf{y}) = e^{y_1} + e^{y_2} + e^{y_3}$, and

$$A = \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ -1 & 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix}$$

the gradient and Hessian of g are

$$\nabla g(\mathbf{y}) = \begin{bmatrix} e^{y_1} \\ e^{y_2} \\ e^{y_3} \end{bmatrix}, \quad \nabla^2 g(\mathbf{y}) = \begin{bmatrix} e^{y_1} & 0 & 0 \\ 0 & e^{y_2} & 0 \\ 0 & 0 & e^{y_3} \end{bmatrix}$$

it follows that

$$\begin{aligned} \nabla f(\mathbf{x}) &= A^T \nabla g(A\mathbf{x} + \mathbf{b}) \\ \nabla^2 f(\mathbf{x}) &= A^T \nabla^2 g(A\mathbf{x} + \mathbf{b}) A \end{aligned}$$

Matlab implementation

```
A=[1 1;1 -1;-1 0];  
b=[1;1;1];  
  
for k=1:50  
%% Gradient descent update %%  
v=exp(A*x-b);  
grad=A'*v;  
if (norm(grad) < 1e-16), break; end;  
x = x - alpha*grad;  
%% Newton's update %%  
vn=exp(A*xn-b);  
gradn=A'*vn;  
D = diag(vn);  
H=A'*D*A;  
dn=-H\gradn;  
xn = xn + alpha*dn;  
end;
```

Numerical example II

$$\text{minimize } f(\mathbf{x}) = \sum_{i=1}^m \log(e^{\mathbf{a}_i^T \mathbf{x} - b_i} + e^{-\mathbf{a}_i^T \mathbf{x} + b_i})$$

- $\mathbf{a}_i \in \mathbb{R}^n$ and $b_i \in \mathbb{R}$ are the problem data
- m and n can be very large
- suppose that we want to solve this problem using Newton's method with initialization $\mathbf{x}^{(0)} = \mathbf{1}$, stopping criteria $\|\nabla f(\mathbf{x}^{(k)})\| < 10^{-5}$, and line search parameters: $\alpha_0 = 1$, $\beta = 1/2$, and $\gamma = 0.01$
- to implement the algorithm, we first need to find the gradient and Hessian of the function f

the function f can be written as

$$f(\mathbf{x}) = g(A\mathbf{x} - \mathbf{b}) \quad \text{where} \quad g(\mathbf{y}) = \sum_{i=1}^m \log(e^{y_i} + e^{-y_i})$$

and

$$A = \begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_m^T \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}$$

the gradient and Hessian of h are:

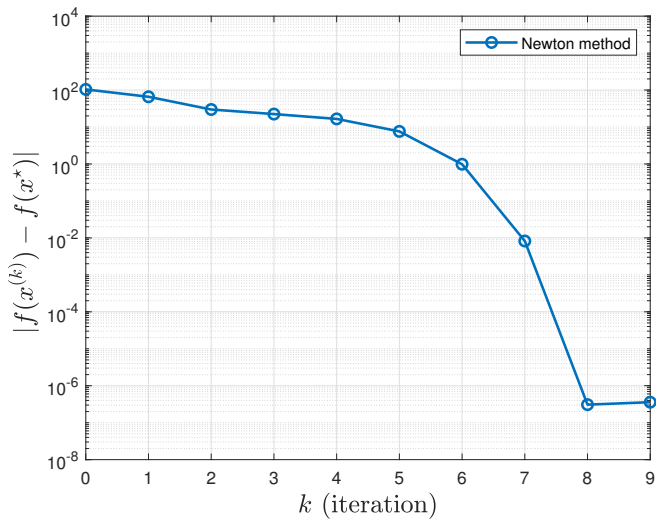
$$\nabla g(\mathbf{y}) = \begin{bmatrix} (e^{y_1} - e^{-y_1}) / (e^{y_1} + e^{-y_1}) \\ \vdots \\ (e^{y_m} - e^{-y_m}) / (e^{y_m} + e^{-y_m}) \end{bmatrix}$$
$$\nabla^2 g(\mathbf{y}) = \text{diag}(4 / (e^{y_1} + e^{-y_1})^2, \dots, 4 / (e^{y_m} + e^{-y_m})^2)$$

using the composition with affine function property, we have

$$\nabla f(\mathbf{x}) = A^T \nabla g(A\mathbf{x} - \mathbf{b}), \quad \nabla^2 f(\mathbf{x}) = A^T \nabla^2 g(A\mathbf{x} - \mathbf{b}) A$$

MATLAB code

```
alpha_0=1;
beta=0.5;
gamma=0.01;
x = ones(n,1); %initialization
k=1;
y = A*x-b;
grad = A'*((exp(y)-exp(-y))./(exp(y)+exp(-y)));
while (norm(grad) >= 1e-5)
k=k+1; %iteration counter
hess = 4*A'*diag(1./(exp(y)+exp(-y)).^2)*A;
d = -hess\grad;
alpha = alpha_0;
f = sum(log(exp(y)+exp(-y)));
while (sum(log(exp(A*(x+alpha*d)-b)+exp(-A*(x+alpha*d)+b))) ...
> f + gamma*alpha*grad'*d)
alpha = beta*alpha;
end
x = x+alpha*d;
y = A*x-b;
f = sum(log(exp(y)+exp(-y)));
grad = A'*((exp(y)-exp(-y))./(exp(y)+exp(-y)));
end
```



References and further readings

- A. Beck. *Introduction to Nonlinear Optimization: Theory, Algorithms, and Applications with MATLAB*, SIAM, 2014.
- E. KP. Chong and S. H. Zak. *An Introduction to Optimization*, John Wiley & Sons, 2013.
- L. Vandenberghe. *EE133A lecture notes*, Univ. of California, Los Angeles.
(<http://www.seas.ucla.edu/~vandenbe/ee133a.html>)
- Uri M. Ascher. *A First Course on Numerical Methods*. Society for Industrial and Applied Mathematics, 2011.