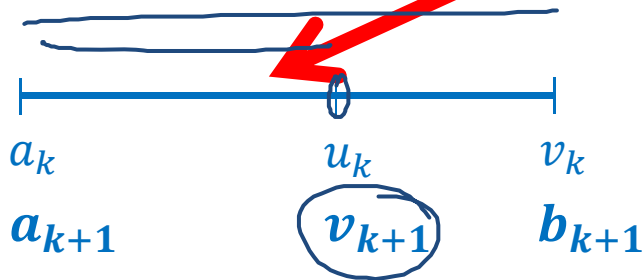
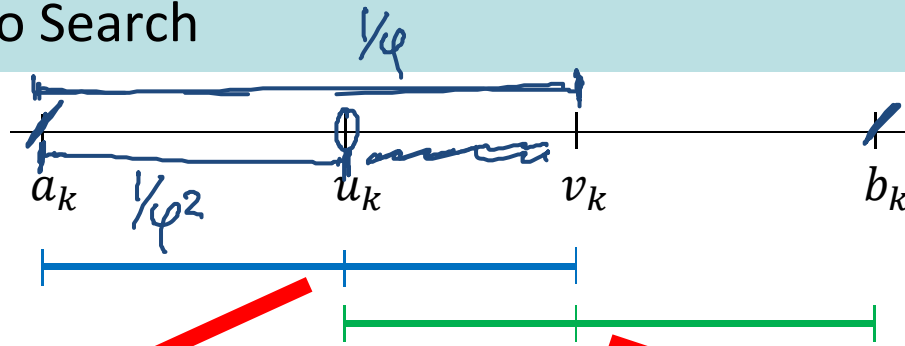


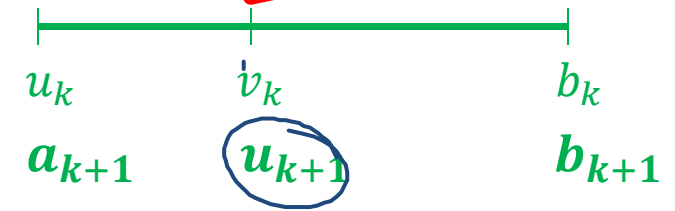
Optimierung

Vorlesung 3

Optimierung ohne Nebenbedingungen
Newton- und Quasi-Newton-Verfahren



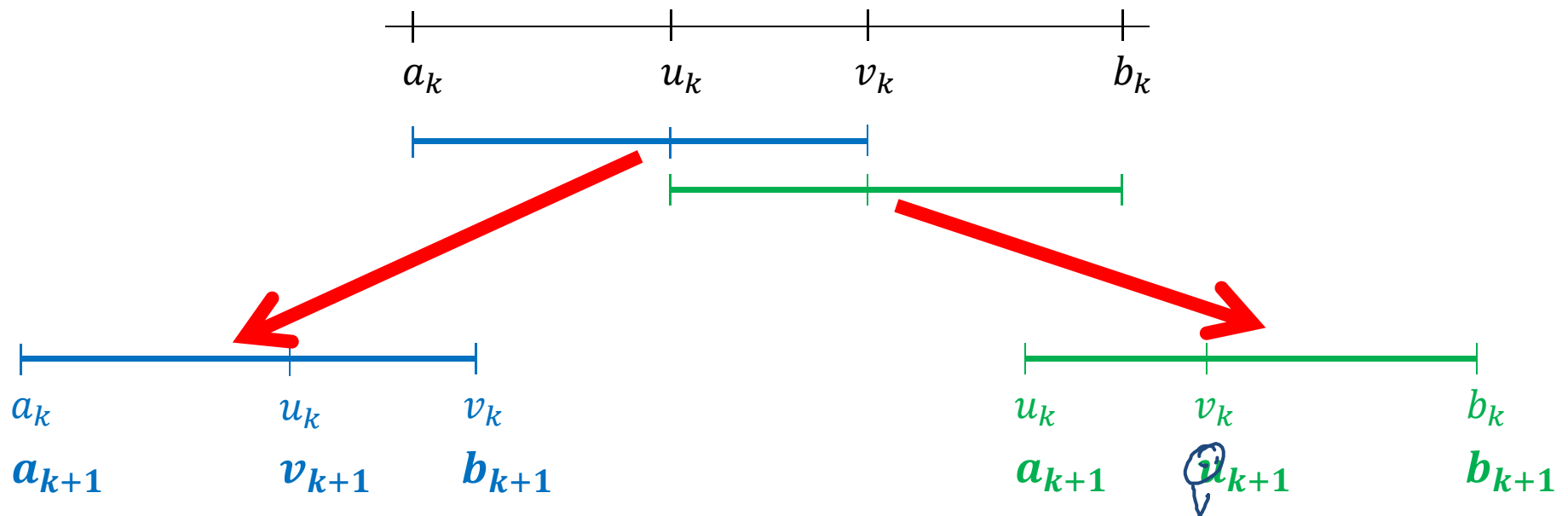
$$\varphi := \frac{b_k - a_k}{v_k - a_k}$$



$$\frac{v_k - a_k}{u_k - a_k} = \frac{b_{k+1} - a_{k+1}}{v_{k+1} - a_{k+1}} = \varphi \Rightarrow \frac{b_k - a_k}{u_k - a_k} = \varphi^2 \rightarrow \frac{b_k - u_k}{v_k - u_k} = \frac{b_{k+1} - a_{k+1}}{u_{k+1} - a_{k+1}} = \varphi^2$$

$$b_k - u_k = (b_k - a_k) \left(1 - \frac{1}{\varphi^2}\right), \quad v_k - u_k = (b_k - a_k) \left(\frac{1}{\varphi} - \frac{1}{\varphi^2}\right)$$

$$\varphi^2 = \frac{b_k - u_k}{v_k - u_k} = \frac{1 - 1/\varphi^2}{1/\varphi - 1/\varphi^2} = \frac{(1 + 1/\varphi)(1 - 1/\varphi)}{1/\varphi(1 - 1/\varphi)} = \varphi + 1$$



$$\frac{b_k - a_k}{v_k - a_k} = \varphi \quad \Rightarrow \quad \frac{b_k - a_k}{u_k - a_k} = \varphi^2 \quad \Rightarrow \quad \varphi + 1 = \varphi^2$$

$$\varphi + 1 = \varphi^2 \quad (\varphi > 1) \quad \Rightarrow \quad \varphi = \frac{1 + \sqrt{5}}{2} \approx \underline{\underline{1.618}}$$

Minimierungsproblem $\min_{x \in \mathbb{R}^n} f(x)$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

- Finde stationären Punkt x^* : $\nabla f(x^*) = 0$

Iteratives Verfahren

- Berechne Sequenz x_0, x_1, x_2, \dots ,
- Für $i \rightarrow \infty$, x_i sollte zu stationärem Punkt konvergieren

$$x_{k+1} = x_k + d_k$$

Taylor-Approximation 1. Ordnung von ∇f

$$f'(x+c) \approx f'(x) + c f''(x)$$

$$\nabla f(x_k + d_k) \approx \nabla f(x_k) + H(x_k) \cdot d_k$$

– Hesse'sche Matrix $H(x_k) = \begin{pmatrix} \frac{\partial^2}{\partial x_1^2} f(x_k) & \cdots & \frac{\partial^2}{\partial x_1 \partial x_n} f(x_k) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial x_n \partial x_1} f(x_k) & \cdots & \frac{\partial^2}{\partial x_n^2} f(x_k) \end{pmatrix}$

– Taylor (1. Ordn.): $\nabla f(x_k + d_k) = \begin{pmatrix} \frac{\partial}{\partial x_1} f(x_k + d_k) \\ \vdots \\ \frac{\partial}{\partial x_n} f(x_k + d_k) \end{pmatrix} \approx \begin{pmatrix} \frac{\partial}{\partial x_1} f(x_k) \\ \vdots \\ \frac{\partial}{\partial x_n} f(x_k) \end{pmatrix} + \left(d_k^\top \cdot \nabla \left(\frac{\partial}{\partial x_i} f(x_k) \right) \right)$

Taylor-Approximation 1. Ordnung von ∇f

$$H(x_k) \cdot d_k = -\nabla f(x_k)$$

$$\nabla f(x_k + d_k) \approx \nabla f(x_k) + H(x_k) \cdot d_k$$

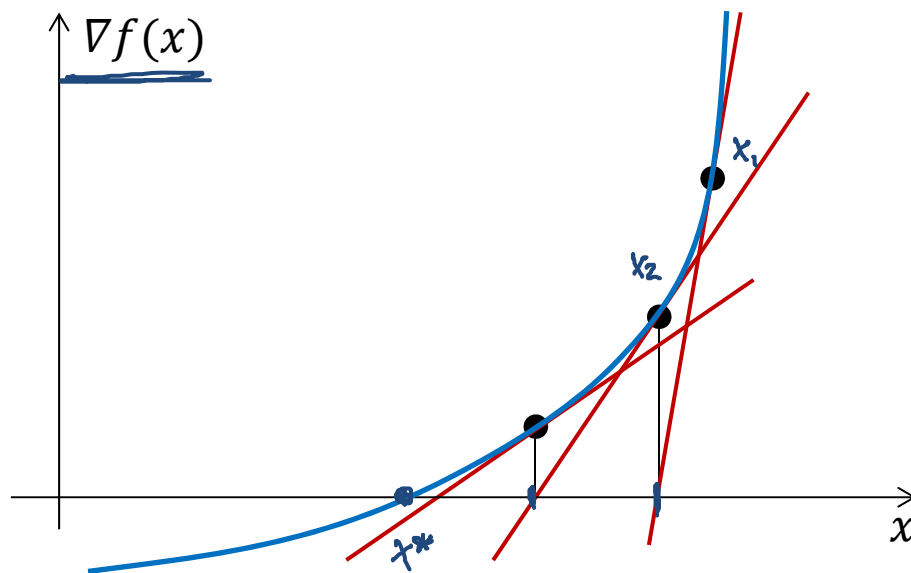
Newton-Verfahren

$$0 \stackrel{!}{=} \nabla f(x_k + d_k) \approx \nabla f(x_k) + H(x_k) \cdot d_k$$

- Setze $\nabla f(x_k) + H(x_k) d_k = 0$ und berechne $x_{k+1} = x_k + d_k$

- Löse lineares Gleichungssystem für d_k :

$$d_k = -H(x_k)^{-1} \cdot \nabla f(x_k)$$



Taylor-Approximation 2. Ordnung von f

$$\underline{f(x+c) \approx f(x) + c f'(x) + \frac{c^2}{2} f''(x)}$$

$$\underline{f(x_k + d_k) \approx f(x_k) + d_k^T \cdot \nabla f(x_k) + \frac{1}{2} d_k^T \cdot H(x_k) \cdot d_k}$$

- $H(x_k)$: Hesse'sche Matrix an der Stelle x_k

$$\frac{1}{2} \left(\frac{\partial^2}{\partial x_i^2} f(x_k) \cdot d_{k,i}^2 + \dots + 2 \frac{\partial^2}{\partial x_i \partial x_j} f(x_k) d_{k,i} d_{k,j} \right)$$

Newton-Verfahren

quadr. Approx von $f(x)$

- Neuer Wert x_{k+1} : minimiere 2. Ordnung Taylor-Approximation

$$d_k = \arg \min_{d \in \mathbb{R}^n} \left\{ \underbrace{f(x_k) + d^T \nabla f(x_k) + \frac{1}{2} d^T H(x_k) d}_{g(d)} \right\}$$

$$\frac{1}{2} c x^2$$

$$x_{k+1} = x_k + d_k$$

$$\nabla g(d) = 0$$

$$\nabla g(d) = \nabla f(x_k) + H(x_k) \cdot d = 0$$

$$H(x_k) d_k = -\nabla f(x_k)$$

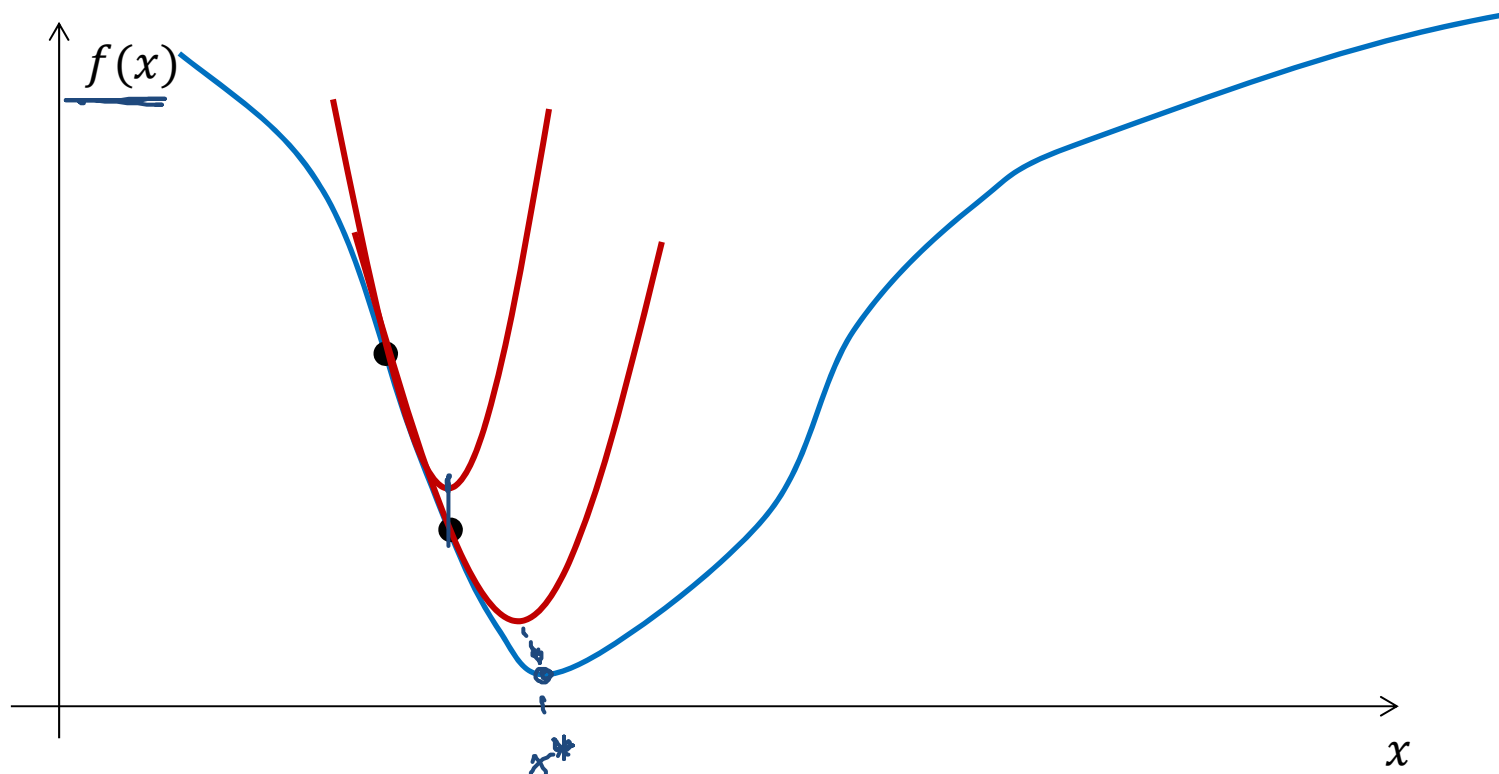
$$\underline{d_k = -H(x_k)^{-1} \nabla f(x_k)}$$

Taylor-Approximation 2. Ordnung von f

$$f(x_k + d_k) \approx f(x_k) + d_k^\top \cdot \nabla f(x_k) + \frac{1}{2} d_k^\top \cdot H(x_k) \cdot d_k$$

Newton-Verfahren

- Neuer Wert x_{k+1} : minimiere 2. Ordnung Taylor-Approximation



$$f(x) = 5x - \ln x$$

- Gradient: $\nabla f(x) = f'(x) = \underline{5 - 1/x}$, $f'(x^*) = 0 \Rightarrow \underline{x^* = 1/5}$
- Hesse'sche Matrix: $H(x) = f''(x) = 1/x^2$
- Newton-Richtung: $\underline{d_k = -H(x_k)^{-1} \nabla f(x_k) = -\frac{f'(x_k)}{f''(x_k)} = \underline{x_k - 5x_k^2}}$
 $\underline{x_{k+1} = x_k + d_k = \underline{2x_k - 5x_k^2}}$

k	x_k	x_k
1	0.05	0.3
2	0.0875000000000000	0.15
3	0.1367187500000000	0.1875
4	0.179977416992188	0.19921875
5	0.197995480848476	0.19996948242188
6	0.199979909514856	0.19999999953434
7	0.199999997981862	0.2000000000000000

$$\ln(2x) = \ln(x) + \ln 2$$

$$\underline{f(x, y)} = \ln(1 + x + y) + \ln 2x + \ln 5y$$

$$\bullet \nabla f(x, y) = \begin{pmatrix} \frac{1}{1+x+y} + \frac{1}{x} \\ \frac{1}{1+x+y} + \frac{1}{y} \end{pmatrix}, \quad \underline{\nabla f(x^*, y^*) = 0} \Rightarrow \underline{x^* = -\frac{1}{3}, y^* = -\frac{1}{3}}$$

$$\bullet H(x, y) = - \begin{pmatrix} \left(\frac{1}{1+x+y}\right)^2 + \frac{1}{x^2} & \left(\frac{1}{1+x+y}\right)^2 \\ \left(\frac{1}{1+x+y}\right)^2 & \left(\frac{1}{1+x+y}\right)^2 + \frac{1}{y^2} \end{pmatrix}$$

$$\bullet \text{Newton-Richtung: } \underline{-H(x, y)d_k = \nabla f(x, y)}$$

$$\bullet \text{Startwerte: } \underline{x_1 = -0.05}, \underline{x_2^y = -0.75}$$

$$\underline{\begin{pmatrix} 425 & 25 \\ 25 & 25 + 16/9 \end{pmatrix}} \cdot \underline{\begin{pmatrix} d_{k,x} \\ d_{k,y} \end{pmatrix}} = \underline{\begin{pmatrix} -15 \\ 5 - 4/3 \end{pmatrix}} \Rightarrow \underline{d_k = \begin{pmatrix} -0.0458677686 \\ 0.1797520661 \end{pmatrix}}$$

$$x_{k+1} = x_k + d_k$$

Abstiegsrichtung

- Eine Richtung \underline{d} ist eine Abstiegsrichtung für einen Wert $\underline{x_k}$, falls

$$\underline{d}^T \nabla f(x_k) < 0$$

- Beispiel: Neg. Gradient $\underline{d} = \underline{-\nabla f(x_k)}$

$$-\nabla f^T \nabla f < 0$$

- Taylor-Approximation: $f(\underline{x_k} + \underline{\varepsilon} \cdot \underline{d}) \approx \underline{f(x_k)} + \underline{\varepsilon \cdot d^T \nabla f(x_k)} < \underline{f(x_k)}$

Newton-Richtung $\underline{d_k} = -H(x_k)^{-1} \nabla f(x_k)$

$$\underline{\forall v: v^T H(x_k) v > 0}$$

- Falls $H(x_k)$ positiv definit ist, ist $\underline{d_k}$ eine Abstiegsrichtung

$$\underline{d_k}^T \nabla f(x_k) < 0$$

$$\begin{aligned} (-H \nabla f)^T \nabla f &= \nabla f^T H^{-1} \nabla f \\ &= \nabla f^T H^{-1} H H^{-1} \nabla f \\ &= \underbrace{(H^{-1} \nabla f)^T}_v H \underbrace{(H^{-1} \nabla f)}_v \\ &= - \underbrace{v^T H v}_{> 0} < 0 \end{aligned}$$

$$\begin{aligned} (AB)^T &= B^T A^T \\ H^{-1T} &= H^{-1} \end{aligned}$$

Konvexe Funktionen $f(x)$:

- Newton-Richtung ist eine Abstiegsrichtung
- Newton-Verfahren konvergiert zu lokalem Minimum

Konvexe, quadratische Funktionen $f(x)$:

- Newton-Schritt geht zu Minimum der quadratischen Approx. von $f(x)$
- Verfahren erreicht lokales Minimum in einem Schritt!

Allgemein:

- Konvergiert sobald man *nahe genug* bei einem lokalen Minimum ist
- Kann mit line search (wie beim Gradientenverfahren) kombiniert werden:

$$\text{Newton-Richtung: } d_k = \underline{-H(x_k)^{-1} \cdot \nabla f(x_k)}$$

Bestimme $x_{k+1} = x_k + \underline{\tau_k} \cdot \underline{d_k}$, so dass $f(\underline{x_k + \tau_k d_k})$ minimal ist

- Wie schnell kommt man zum lokalen Minimum?

Konvergenz von Folgen $s_1, s_2, \dots, \lim_{k \rightarrow \infty} s_k = \bar{s}$

Lineare Konvergenz:

$$\limsup_{k \rightarrow \infty} \frac{|s_{k+1} - \bar{s}|}{|s_k - \bar{s}|} = \underline{C} < 1$$

$$s_k = 0.9^k$$

Superlineare Konvergenz:

$$\limsup_{k \rightarrow \infty} \frac{|s_{k+1} - \bar{s}|}{|s_k - \bar{s}|} = \underline{0}$$

$$s_k = 1/k!$$

Quadratische Konvergenz:

$$\limsup_{k \rightarrow \infty} \frac{|s_{k+1} - \bar{s}|}{|s_k - \bar{s}|^2} = C$$

$$s_k = 0.9^{2^k}$$

Gradientenverfahren (steilster Abstieg)

- lineare Konvergenz

Newtonverfahren

- quadratische Konvergenz
- minimiert quadratische Funktionen in 1 Schritt

Konjugierte Gradienten

- minimiert n -dimensionale quadratische Funktionen in n Schritten
- quadratische Konvergenz alle n Schritte (superlinear)

Quasi-Newtonverfahren

- superlineare Konvergenz
- einfacher als Newton-Verfahren

pos. def.

Annahme: streng konvexe, quadratische Fkt. $f(x) = \frac{1}{2} x^T A x - b^T x$

- A symm. pos. definit ($\forall x: x^T A x > 0$, alle Eigenwerte reell und positiv)
- Eindeutiges Minimum bei x^* , so dass $\nabla f(x^*) = 0$

$$0 = \nabla f(x^*) = Ax^* - b \implies x^* = A^{-1}b$$

- Funktionswert bei x^* : $f(x^*) = \frac{1}{2} b^T A^{-1} A A^{-1} b - b^T A^{-1} b = -\frac{1}{2} b^T A^{-1} b$

1 Iterationsschritt an Stelle x

- Abstiegsrichtung d = $-\nabla f(x) = b - Ax$
- Neuer Iterationswert $x' = x + \tau \cdot d$, wobei $\tau = \arg \min_{\alpha} f(x + \alpha \cdot d)$

$$\begin{aligned} f(x+\alpha d) &= \frac{1}{2} (x+\alpha d)^T A (x+\alpha d) - b^T (x+\alpha d) \\ &= \frac{1}{2} (x^T A x + \alpha d^T A x + \alpha x^T A d + \alpha^2 d^T A d) - \underbrace{b^T x}_{\alpha d^T b} - \alpha b^T d \\ &= f(x) + \alpha d^T (Ax - b) + \frac{1}{2} d^T A d \cdot \alpha^2 = \frac{-\alpha d^T d + \frac{1}{2} \alpha^2 d^T A d + f(x)}{\alpha} \\ \frac{\partial}{\partial \alpha} f(x+\alpha d) &= d^T A d - d^T d = 0 \implies \tau = \frac{d^T d}{d^T A d} \end{aligned}$$

Quadratische konvexe Funktion

$$f(x) = \frac{1}{2} x^\top A x - b^\top x, \quad x^* = A^{-1} b, \quad f(x^*) = -\frac{1}{2} b^\top A^{-1} b$$

1 Iterationsschritt an Stelle x (neuer Wert x')

$$d = -\nabla f(x) = \underline{b - Ax}, \quad x' = \underline{x + \tau \cdot d}, \quad \tau = \underline{\frac{d^\top d}{d^\top A d}}$$

Konvergenzrate C

$$C = \left(\max_{x \text{ nahe genug bei } x^*} \right) \frac{f(x') - f(x^*)}{f(x) - f(x^*)}$$

$$\begin{aligned} \underline{f(x')} &= f(x + \tau d) = f(x) + \frac{1}{2} \tau^2 d^\top A d - \tau d^\top d \\ &= f(x) + \frac{1}{2} \frac{(d^\top d)^2}{d^\top A d} - \frac{(d^\top d)^2}{d^\top A d} \\ &= \underline{f(x) - \frac{1}{2} \frac{(d^\top d)^2}{d^\top A d}} \end{aligned}$$

Quadratische konvexe Funktion

$$f(x) = \frac{1}{2} x^T A x - b^T x, \quad x^* = A^{-1} b, \quad \underline{f(x^*)} = -\frac{1}{2} b^T A^{-1} b$$

$$\underline{d} = b - Ax, \quad x' = x + \tau d, \quad \tau = \frac{d^T d}{d^T A d}, \quad \underline{f(x')} = f(x) - \frac{1}{2} \frac{(d^T d)^2}{d^T A d}$$

Konvergenzrate C

$$C = \frac{f(x') - f(x^*)}{f(x) - f(x^*)} = \frac{f(x) - \frac{1}{2} \frac{(d^T d)^2}{d^T A d} + \frac{1}{2} b^T A^{-1} b}{f(x) + \frac{1}{2} b^T A^{-1} b}$$

$\frac{\alpha + \beta}{\alpha} = 1 + \frac{\beta}{\alpha}$

$$= 1 - \frac{\frac{1}{2} \frac{(d^T d)^2}{d^T A d}}{f(x) + \frac{1}{2} b^T A^{-1} b} = 1 - \frac{\frac{1}{2} \frac{(d^T d)^2}{d^T A d}}{\frac{1}{2} x^T A x - b^T x + \frac{1}{2} b^T A^{-1} b}$$

$$\frac{1}{2} x^T A A^{-1} A x = \frac{1}{2} (Ax)^T A^{-1} (Ax)$$

$$\frac{1}{2} (Ax)^T A^{-1} (Ax) - b^T x + \frac{1}{2} b^T A^{-1} b = \frac{1}{2} (Ax - b)^T A^{-1} (Ax - b) = \frac{1}{2} d^T A^{-1} d$$

$$\frac{1}{2} x^T A A^{-1} (-b)$$

Konvergenzrate C

$$C = \frac{f(x') - f(x^*)}{f(x) - f(x^*)} = 1 - \frac{\frac{1}{2} \frac{(d^\top d)^2}{d^\top A d}}{\frac{1}{2} d^\top A^{-1} d} = 1 - \frac{d^\top d}{d^\top A d} \cdot \frac{d^\top d}{d^\top A^{-1} d}$$

- Eigenwerte von A : $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ (positiv und reell)
- Eigenwerte von A^{-1} : $0 < 1/\lambda_n \leq 1/\lambda_{n-1} \leq \dots \leq 1/\lambda_1$
- Grösster Eigenwert μ_{\max} einer reellen, symm. $n \times n$ -Matrix M :

$$\mu_{\max} = \max_{x \in \mathbb{R}^n} \frac{x^\top M x}{x^\top x} \quad \text{Rayleigh-Koeffizient}$$

- Deshalb: $\frac{d^\top A d}{d^\top d} \leq \lambda_n$, $\frac{d^\top A^{-1} d}{d^\top d} \leq 1/\lambda_1 \Rightarrow C \leq 1 - \frac{\lambda_1}{\lambda_n}$ Konditionszahl von A

- Genauer: $C \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2$



Ziel: Die Sequenz x_1, x_2, \dots konvergiert quadratisch gegen x^*

Matrix-Norm $\|M\|$ einer Matrix M

$$\|M\| := \max_{x: \|x\|=1} \|Mx\|$$

- $\forall M, M', x: \|Mx\| \leq \|M\| \cdot \|x\|, \quad \|M \cdot M'\| \leq \|M\| \cdot \|M'\|$

Annahmen: Für x, y nahe genug bei x^* gilt $\|H(x) - H(y)\| \leq L \cdot \|x - y\|$,
sowie $\|H(x)^{-1}\| \leq C$ für Konstanten $L, C \geq 0$

Lemma:

$$\nabla f(z) - \nabla f(x) = \int_0^1 [H(x + t(z - x))](z - x) dt$$

Lemma:

$$\nabla f(z) - \nabla f(x) = \int_0^1 [H(x + t(z - x))](z - x) dt$$

Beweis:

- Definiere $\phi(t) := \nabla f(x + t(z - x))$
- Ableitung $\phi'(t) = [H(x + t(z - x))](z - x)$

Ziel: Die Sequenz x_1, x_2, \dots konvergiert quadratisch gegen x^*

Lemma:

$$\nabla f(z) - \nabla f(x) = \int_0^1 [H(x + t(z - x))](z - x) dt$$



$$x' - x^* = x - H(x)^{-1} \nabla f(x) - x^*$$

Matrix-Norm $\|M\|$ einer Matrix M

- $\forall M, M', x: \|Mx\| \leq \|M\| \cdot \|x\|, \quad \|M \cdot M'\| \leq \|M\| \cdot \|M'\|$

Annahmen: Für x, y nahe genug bei x^* gilt $\|H(x) - H(y)\| \leq L \cdot \|x - y\|$,
sowie $\|H(x)^{-1}\| \leq C$ für Konstanten $L, C \geq 0$

$$x' - x^* = H(x)^{-1} \int_0^1 [H(x + t(x^* - x)) - H(x)](x^* - x) dt$$

$$\|x' - x^*\| \leq \|H(x)^{-1}\| \int_0^1 \|H(x + t(x^* - x)) - H(x)\| \cdot \|x^* - x\| dt$$

Gradientenverfahren

- Lineare Konvergenz
 - Benötigt i.A. deutlich mehr Schritte (vor allem, wenn die Hesse'sche Matrix schlecht konditioniert ist)
- 1. Ableitung nötig
- 1 Iterationsschritt:
 - Evaluation Gradient + line search
- “Billige” Schritte, “schlechte” Konv.
- Bessere Konvergenz:
Konjugierte Gradientenverfahren

Newton-Verfahren

- Quadratische Konvergenz
 - konvergiert in wenigen Schritten (bei quadr. Problem in 1 Schritt)
- 1. & 2. Ableitung nötig
- 1 Iterationsschritt:
 - Evaluation Gradient + Hesse'sche Matrix
 - n -dim. lineares Gl.-system
- “Teure” Schritte, “gute” Konvergenz
- Einfachere Iterationen:
Quasi-Newton-Verfahren

- Newton-Verfahren + Line Search:

$$x_{k+1} = x_k + \tau_k d_k, \quad \text{wobei } d_k = -H(x_k)^{-1} \nabla f(x_k)$$

- **Idee:** Ersetze $H(x_k)^{-1}$ durch eine Matrix D_k , welche einfacher berechnet werden kann
- Taylor-Approximation: $\nabla f(x_{k+1}) - \nabla f(x_k) \approx H(x_{k+1})(x_{k+1} - x_k)$

– Idee: D_k^{-1} sollte dies auch erfüllen...

- Quasi-Newton Bedingung:

$$D_{k+1} \cdot q_k = p_k, \quad \text{wobei } p_k = x_{k+1} - x_k, \quad q_k = \nabla f(x_{k+1}) - \nabla f(x_k)$$

Quasi-Newton-Verfahren:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \tau_k \cdot \mathbf{d}_k, \quad \text{wobei } \mathbf{d}_k = -\mathbf{D}_k \nabla f(\mathbf{x}_k)$$

- τ_k wird typischerweise durch line search ermittelt

Quasi-Newton Bedingung:

$$\mathbf{D}_{k+1} \cdot \mathbf{q}_k = \mathbf{p}_k, \quad \text{wobei } \mathbf{p}_k = \mathbf{x}_{k+1} - \mathbf{x}_k, \quad \mathbf{q}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)$$

Berechnung von \mathbf{D}_{k+1} :

- Alte Matrix + Korrekturterm: $\mathbf{D}_{k+1} = \mathbf{D}_k + \mathbf{C}_k$

$$(\mathbf{D}_k + \mathbf{C}_k)\mathbf{q}_k = \mathbf{p}_k \quad \Rightarrow \quad \mathbf{C}_k\mathbf{q}_k = \mathbf{p}_k - \mathbf{D}_k\mathbf{q}_k$$

- Zum Beispiel ($\phi \in [0,1]$ ist ein Parameter)

$$\mathbf{C}(\phi) = \frac{\mathbf{p}\mathbf{p}^\top}{\mathbf{p}^\top\mathbf{q}} - \frac{\mathbf{D}\mathbf{q}\mathbf{q}^\top\mathbf{D}}{\mathbf{q}^\top\mathbf{D}\mathbf{q}} + \phi\eta\mathbf{v}\mathbf{v}^\top, \quad \text{wobei } \mathbf{v} = \frac{\mathbf{p}}{\mathbf{p}^\top\mathbf{q}} - \frac{\mathbf{D}\mathbf{q}}{\eta}, \quad \eta = \mathbf{q}^\top\mathbf{D}\mathbf{q}$$

Quasi-Newton-Verfahren:

$$x_{k+1} = x_k - \tau_k \cdot D_k \nabla f(x_k)$$

$$D_{k+1} = D_k + C_k, \quad C_k q_k = p_k - D_k q_k$$

- Parameter $\phi = 0$ (Davidson-Fletcher-Powell, erstes Quasi-Newton-Verf.)

$$C(0) = \frac{pp^T}{p^T q} - \frac{Dqq^T D}{q^T Dq}$$

- Parameter $\phi = 1$ (Broyden-Fletcher-Goldfarb-Shanno)

$$C(1) = \frac{pp^T}{p^T q} \left[1 + \frac{q^T Dq}{p^T q} \right] - \frac{Dqp^T + pq^T D}{p^T q}$$

- Quasi-Newton-Verfahren konvergieren oft superlinear
- Für quadratische Funktionen in n Schritten
 - entspricht dann dem konjugierten Gradientenverfahren