

### 3. Iterative Methoden für lineare Gleichungssysteme

(1) Problemstellung:  $A \underline{x} = \underline{b}$

$$A \in \mathbb{R}^{n \times n}, \text{ regulär}$$

$$\underline{b} \in \mathbb{R}^n$$

Speziell:  $n$  sehr gross, z.B.  $n = 10'000$   
 $A$  sparse (schwach besetzt)

Gesucht: Lösung  $\underline{x} = A^{-1} \underline{b}$ , ev. approximativ  
 unter Ausnutzung der Sparsity:  
 Matrix  $A$  (ein Unding!) nicht  
 aufschreiben, nicht speichern.  
 Nur die Abbildung  $\underline{x} \mapsto A \underline{x}$  brauchen!

#### 3.1. Die Jacobi-Iteration

Sei  $A = C - B$  eine additive Zerlegung  
 von  $A$

$C$  so, dass Gl'systeme mit  $C$  einfach lösbar

(2) Aus (1) folgt  $C \underline{x} = \underline{b} + B \underline{x}$

eine Vektorgleichung  
 in Fixpunktform; lösen  
 durch Iteration:

Startvektor  $\underline{x}_0$ , z.B.  $\underline{x}_0 = 0$

(3)  $\underline{x}_{k+1}$  aus  $C \underline{x}_{k+1} = \underline{b} + B \underline{x}_k$ ,  $k = 0, 1, \dots$   
 einfach lösbar!

Fehlertheorie

Subtrahiere (2) von (3) :

$$C(\underline{x}_{k+1} - \underline{x}) = B(\underline{x}_k - \underline{x})$$

Def : Fehlervektor  $\underline{e}_k := \underline{x}_k - \underline{x}$ Fehlergesetz :  $C \underline{e}_{k+1} = B \underline{e}_k$ 

oder  $\underline{e}_{k+1} = (C^{-1}B) \underline{e}_k$

Normen :

$$\|e_{k+1}\| \leq \|C^{-1}B\| \cdot \|e_k\|$$

SATZ : Hinreichend für die Konvergenz des Iterationsverfahrens
$$C \underline{x}_{k+1} = \underline{b} + B \underline{x}_k \text{ zur Lösung}$$

von  $(C-B)\underline{x} = \underline{b}$  ist  $\|C^{-1}B\| < 1$ .

Beispiele :(i) Jacobi - Iteration .  $A = (a_{jk})$ 

Sei A "diagonal dominant", d.h.

$$(4) \quad |a_{jj}| > \sum_{k \neq j} |a_{jk}|, \quad j=1, 2, \dots, n.$$

Wähle  $C := \text{diag}(a_{jj})$ Aus (4) folgt sofort  $\|C^{-1}B\|_{\infty} < 1$ 

(siehe S. 21)

(ii) Gauss-Seidel-Iteration

$$A = \begin{pmatrix} & & -B \\ & & \\ C & & \end{pmatrix}$$

häufig leicht besser als Jacobi

### 3.2. Das Verfahren der konjugierten Gradienten

cg ; M. Hestenes, E. Stiefel 1951/52

Sei  $A \in \mathbb{R}^{n \times n}$  **symmetrisch**,  $A^T = A$ , und **positiv definit**, d.h. die zu A gehörige quadratische Form erfüllt

$$Q(\underline{x}) := \underline{x}^T A \underline{x} > 0 \quad \forall \underline{x} \neq 0$$

#### Bemerkungen.

(i) A hat lauter positive Eigenwerte

Beweis. Sei  $\lambda$  Ew. von A,  $\underline{v}$  zugehöriger Ev.

$$\begin{aligned} A \underline{v} &= \underline{v} \cdot \lambda \Rightarrow 0 < \underline{v}^T A \underline{v} = \underbrace{\underline{v}^T \underline{v}}_{\|\underline{v}\|^2} \cdot \lambda \\ &\Rightarrow \lambda > 0, \text{ g.e.d.} \end{aligned}$$

(ii) Graphische Darstellungen von  $Q(\underline{x}) = \underline{x}^T A \underline{x}$  für  $\underline{x} = (x_1, x_2)^T \in \mathbb{R}^2$  siehe S. 50.

Fälle: A **positiv definit**, **pos. semidefinit**, **indefinit**

#### Problemstellung

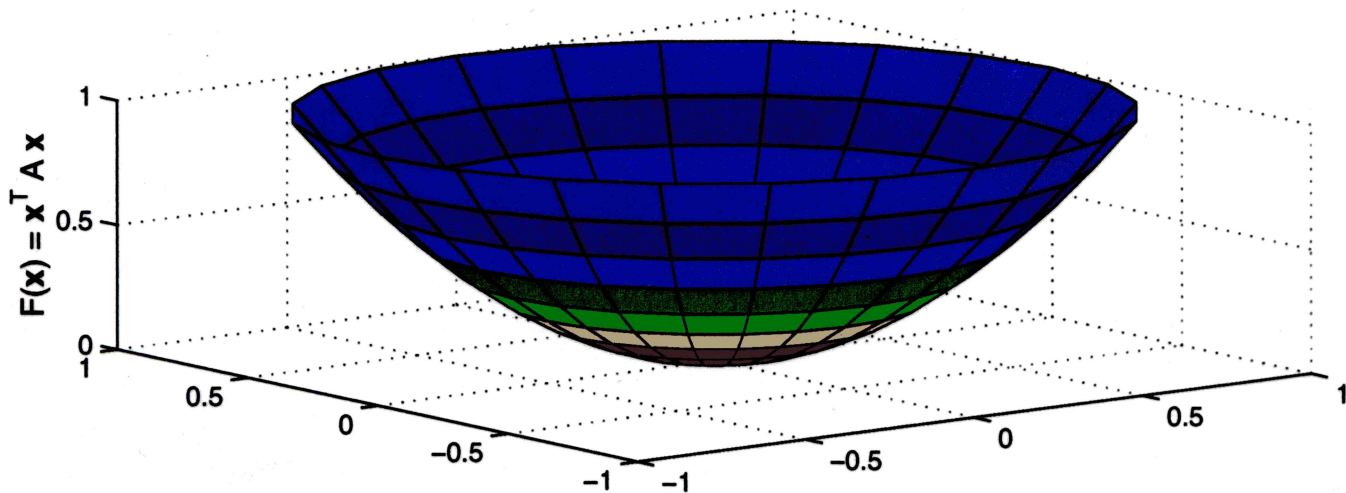
Löse lineares Gleichungssystem

$$A \underline{x} + \underline{b} = 0,$$

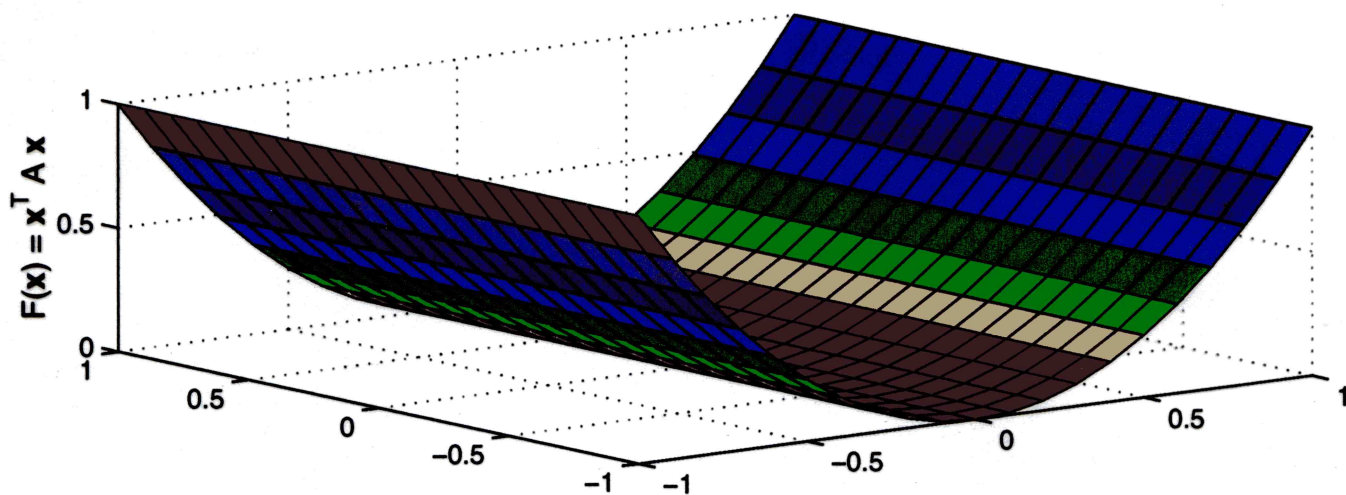
$$A^T = A, \text{ pos. def.}, \text{ **sparse**, } \underline{x}, \underline{b} \in \mathbb{R}^n$$

Bem: Differenzenverfahren für elliptische part. DGl. führen häufig auf solche Gl'systeme.

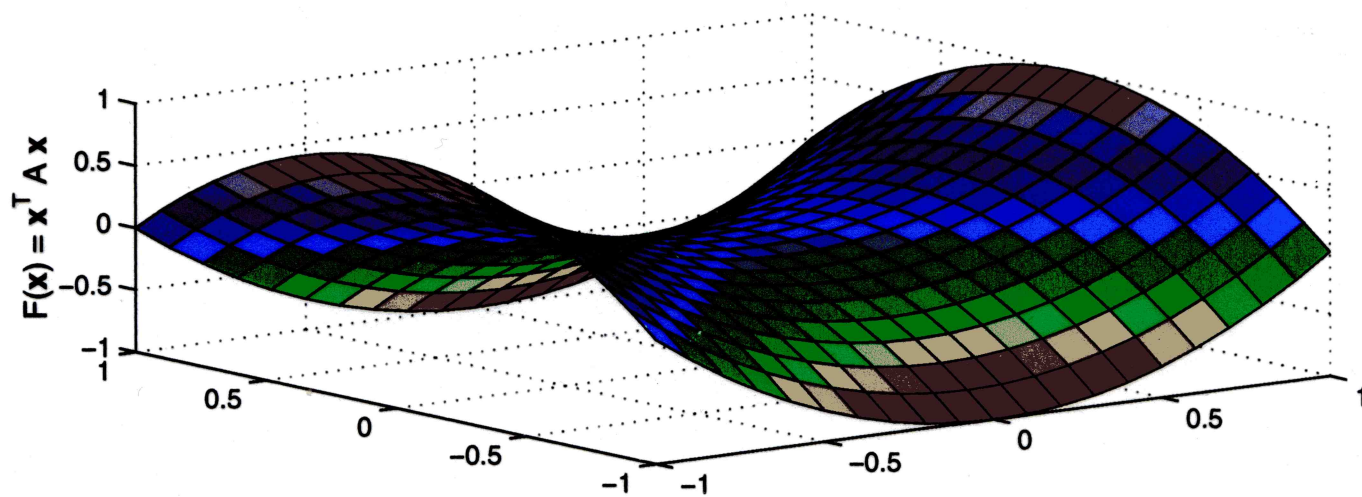
# Quadratische Formen in $x = (x_1, x_2)^T$



$A=A^T$  positiv definit



$A=A^T$  positiv semidefinit



$A=A^T$  indefinit

SATZ :  $A \underline{x} + \underline{b} = 0$ ,  $A$  sym. pos. def. (51)  
ist äquivalent mit

$$F(\underline{u}) := \frac{1}{2} \underline{u}^T A \underline{u} + \underline{u}^T \cdot \underline{b}$$

hat in  $\underline{x}$  ein Minimum.

Beweis.

(i) Differentiation ergibt

$$\text{grad } F = \nabla F = A \underline{u} + \underline{b} =: \underline{r} = \text{Residuum}$$

Damit folgt " $\Leftarrow$ " sowie " $\underline{x}$  ist stationärer Punkt".

(ii) Minimaleigenschaft:

Betrachte  $\underline{v} \in \mathbb{R}^n$ ,  $\underline{v} \neq 0$ :

$$\begin{aligned} F(\underline{x} + \underline{v}) &= \frac{1}{2} (\underline{x} + \underline{v})^T A (\underline{x} + \underline{v}) + (\underline{x} + \underline{v})^T \underline{b} \\ &= F(\underline{x}) + \underbrace{\frac{1}{2} \underline{x}^T A \underline{v} + \frac{1}{2} \underline{v}^T A \underline{x}}_{= \frac{1}{2} \underline{v}^T A \underline{x} \text{ wegen } A^T = A} + \frac{1}{2} \underline{v}^T A \underline{v} + \underline{v}^T \underline{b} \\ &= F(\underline{x}) + \underbrace{\underline{v}^T (A \underline{x} + \underline{b})}_0 + \frac{1}{2} \underline{v}^T A \underline{v} > F(\underline{x}) \quad \text{q.e.d.} \end{aligned}$$

Grundidee: Minimiere  $F$ !  $\rightarrow$

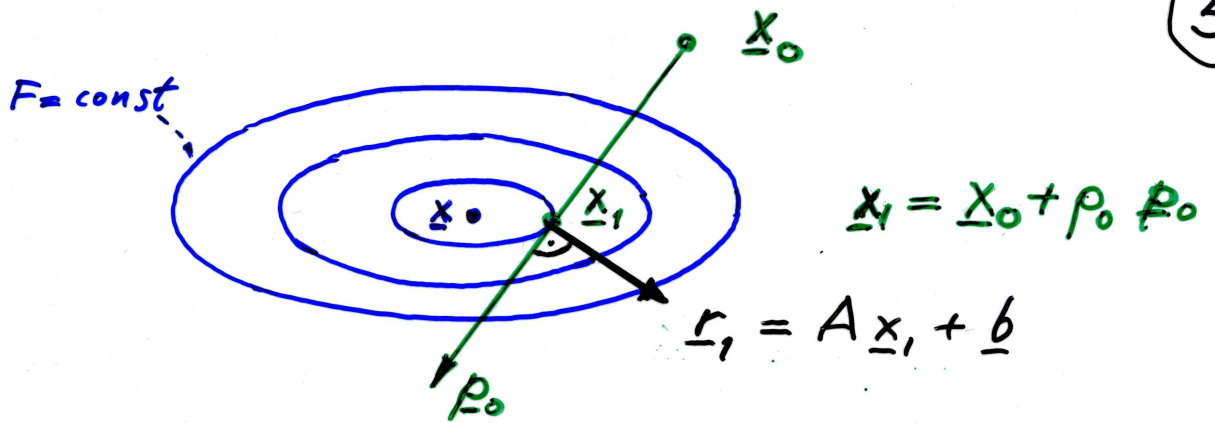
Dies löst automatisch  $A \underline{x} + \underline{b} = 0$

Vorbereitung: 1-dimensionale Minimierung

Gegeben: Startpunkt  $\underline{x}_0$ , Suchrichtung  $\underline{p}_0$

Gesucht:  $\rho = \rho_0 \in \mathbb{R}$  so, dass

$$F(\underline{x}_0 + \rho \underline{p}_0) \stackrel{!}{=} \min$$



Notwendige Bedingung:

$$\frac{d}{d\rho} \left[ \frac{1}{2} (\underline{x}_0^T + \rho \underline{p}_0^T) A (\underline{x}_0 + \rho \underline{p}_0) + (\underline{x}_0^T + \rho \underline{p}_0^T) \underline{b} \right] = 0$$

$$(*) \Rightarrow \rho (\underline{p}_0^T A \underline{p}_0) + \underbrace{\underline{p}_0^T (A \underline{x}_0 + \underline{b})}_{\underline{r}_0} = 0$$

$$\Rightarrow \underline{p}_0 = - \frac{\underline{p}_0^T \underline{r}_0}{\underline{p}_0^T A \underline{p}_0}, \quad \underline{x}_1 = \underline{x}_0 + \rho_0 \underline{p}_0$$

## Strategien

### (a) Steilster Abstieg

Wähle  $\underline{p}_k := - \underline{r}_k = - \nabla F(\underline{x}_k), \quad k \geq 0$

Algorithmus: Gegeben  $\underline{x}_0 \in \mathbb{R}^n$

for  $k = 0, 1, \dots$

$$\underline{r}_k = A \underline{x}_k + \underline{b}$$

$$\rho_k = \frac{\underline{r}_k^T \underline{r}_k}{\underline{r}_k^T A \underline{r}_k}$$

$$\underline{x}_{k+1} = \underline{x}_k - \rho_k \underline{r}_k$$

Leider: In  $n > 2$  Dimensionen viel zu langsame Konvergenz!

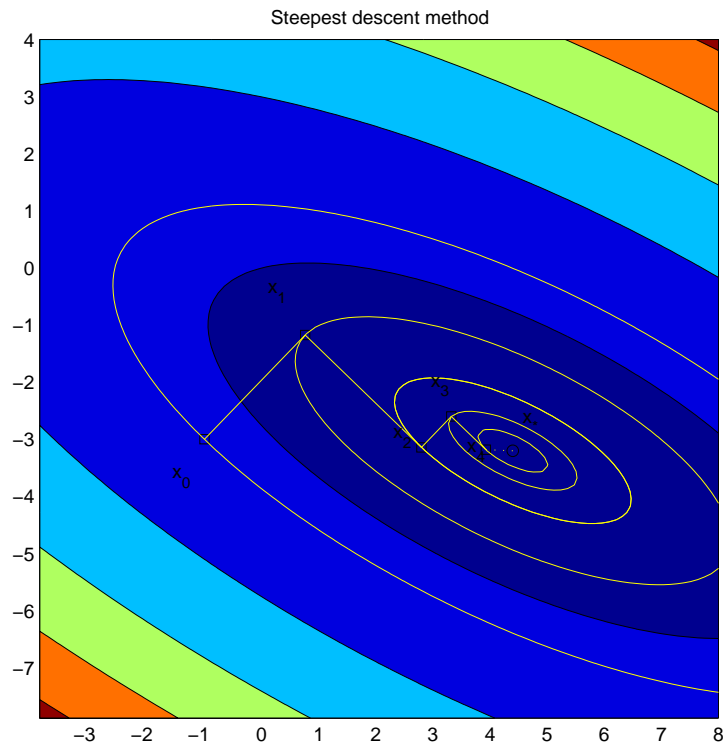


Figure 1: The steepest descent method for  $N = 2$ .

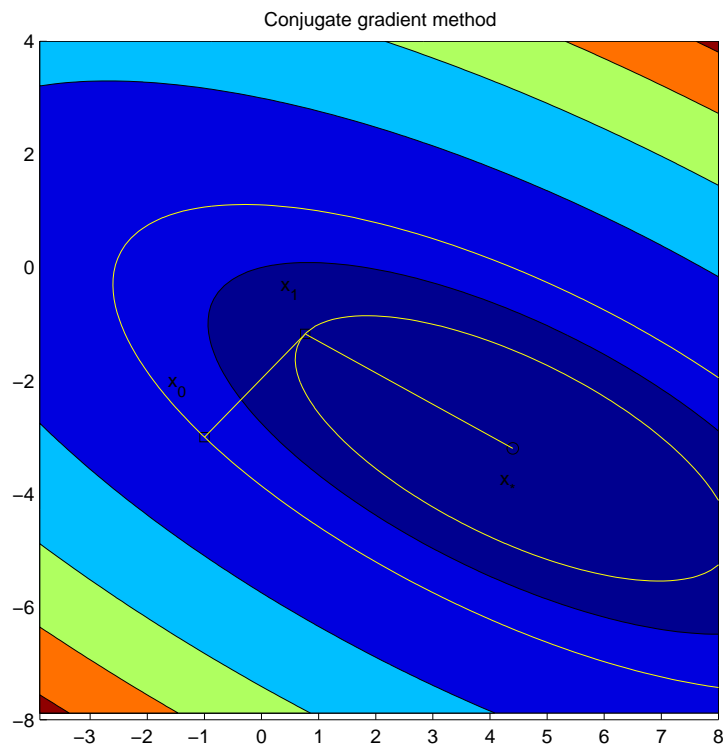


Figure 2: Conjugate directions — the CG method for  $N = 2$ .

(b) Konjugierte Gradienten

Die (zündende!) Idee :

Simultane Minimierung von  $F$  mit zwei (lin. unabh.) Suchrichtungen,  $p_0, p_1$

Gesucht :  $p_0, p_1$  so, dass

$$F(\underline{x}_0 + \rho_0 p_0 + \rho_1 p_1) \stackrel{!}{=} \min$$

Analog zu (\*) ergibt sich :

$$\begin{aligned}
 (**) \quad & \rho_0 (p_0^T A p_0) + \rho_1 (p_1^T A p_0) + p_0^T r_0 = 0 \\
 & \rho_0 (p_0^T A p_1) + \rho_1 (p_1^T A p_1) + p_1^T r_0 = 0
 \end{aligned}$$

Dieses System von 2 linearen Gleichungen ist diagonal (entkoppelt), falls

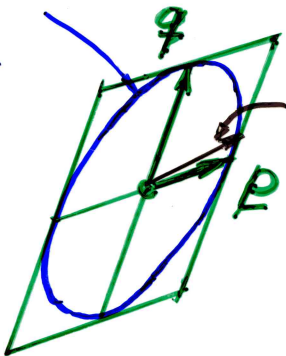
$$p_0^T A p_1 = p_1^T A p_0 = 0$$

Definition

Die beiden Richtungen  $p, q \in \mathbb{R}^n$  heissen konjugiert bezüglich  $A=A^T$ , falls

$$p^T A q = 0$$

Niveaufläche  
 $p^T A p = \text{const}$



geometrisch

$p + \epsilon q$   
mit  $\epsilon \rightarrow 0$

Algebraisch :

Geg  $p$  ; finde eine Tangentenrichtung  $q$  an Niveaufläche im Endpunkt von  $p$ .



Bedingung an  $q$ , im  $\lim_{\epsilon \rightarrow 0}$ :

$$(p + \epsilon q)^T A (p + \epsilon q) = 1$$

subtrahiere  $p^T A p = 1$

---

$$2 \epsilon \cdot \underbrace{p^T A q}_{=0} + O(\epsilon^2) = 0$$

Das cg-Verfahren:

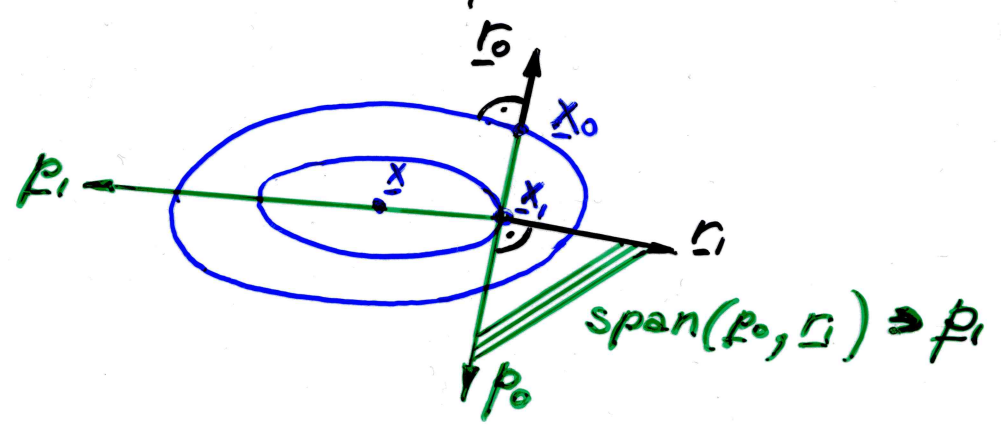
Wähle Suchrichtungen  $p_0, p_1, p_2, \dots$

so, dass  $p_1$  konj. zu  $p_0$ ,  $p_1^T A p_0 = 0$

$p_2$  konj. zu  $p_1$ ,  $p_2^T A p_1 = 0, \dots$

Konsequenzen:

- Algebraisch: Gleichungen für  $p_0, p_1$  entkoppelt
- Geometrisch, z.B. für Dimension  $n=2$ :



$p_1$  zeigt direkt auf das Minimum  $\underline{x}$

Berechnung (Wahl!) von  $p_1$ :

Wähle  $p_1 \in \text{span}(p_0, r_1)$

Ansatz:  $p_1 = \sigma_1 p_0 - r_1$

$$p_0, p_1 \text{ konj} \Rightarrow 0 = p_0^T A p_1 = \sigma_1 p_0^T A p_0 - p_0^T A r_1$$

$$\Rightarrow \sigma_1 = \frac{p_0^T A r_1}{p_0^T A p_0}$$

Da  $p_0, p_1$  konjugiert, folgt aus (\*\*, S. 54):

$$p_1 = - \frac{p_1^T r_1}{p_1^T A p_1}$$

und weiter

$$\underline{x}_2 = \underline{x}_1 + p_1 p_1$$

Zusammengefasst:

cg-Algorithmus zur Lösung von

$$A \underline{x} + \underline{b} = 0, \quad A \text{ sym. pos. def.}$$

Initialisierung:  $\underline{x}_0 = 0, \quad \underline{r}_0 = \underline{b}, \quad p_{-1} = 0$

for  $k=0, 1, \dots$

$$(1) \quad \sigma_k = \begin{cases} 0, & (k=0) \\ \frac{p_{k-1}^T A r_k}{p_{k-1}^T A p_{k-1}} & = \frac{r_k^T r_k}{r_{k-1}^T r_{k-1}} \quad (k>0) \end{cases}$$

$$(2) \quad p_k = \sigma_k p_{k-1} - r_k$$

$$(3) \quad p_k = - \frac{p_k^T r_k}{p_k^T A p_k} = \frac{r_k^T r_k}{p_k^T A p_k}$$

$$(4) \quad \underline{x}_{k+1} = \underline{x}_k + p_k p_k$$

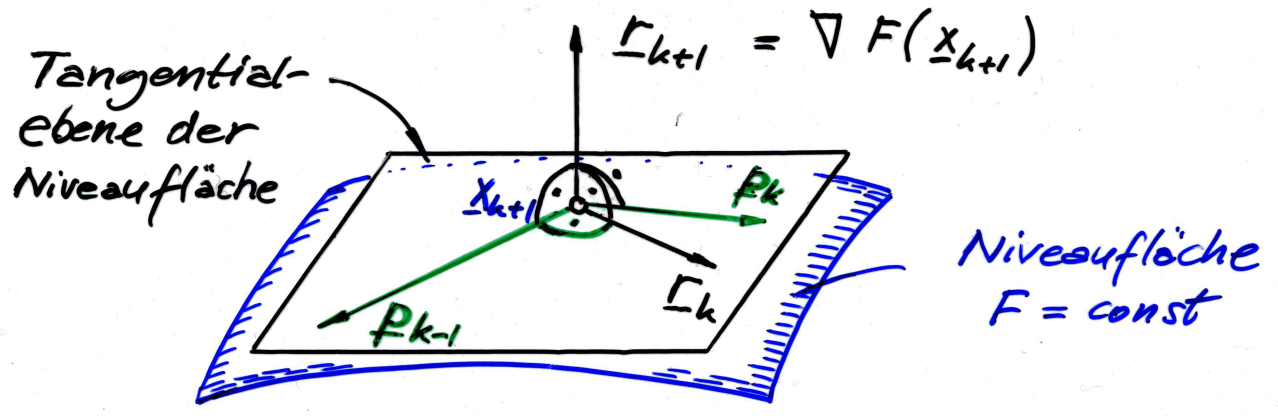
$$(5) \quad \underline{r}_{k+1} = A \underline{x}_{k+1} + \underline{b} = \underline{r}_k + p_k A p_k$$

äquivalent, aber "handlicher", s. S. 57

Für den oben formulierten cg-Algorithmus gilt der

SATZ:  $r_{k+1} \perp r_k, r_{k+1} \perp p_{k-1}, r_{k+1} \perp p_k$   
für  $k = 0, 1, \dots$

Beweis: formal durch vollständige Induktion;  
besser geometrisch:



Damit Umrechnung der Beziehungen (3), (5), (1):

$$(3): \text{Zähler} = p_k^T r_k \stackrel{(2)}{=} (\sigma_k p_{k-1} - r_k)^T r_k$$

$$\stackrel{\text{Satz}}{=} -r_k^T r_k$$

$$(5): r_{k+1} \stackrel{(4)}{=} A(x_k + \rho_k p_k) + b$$

$$= r_k + \rho_k A p_k, \text{ da } r_k := A x_k + b$$

$$(1): \sigma_k = \frac{(p_{k-1} p_{k-1}^T A) r_k}{(p_{k-1} p_{k-1}^T A) p_{k-1}}$$

$$\stackrel{(5)}{=} \frac{(r_k - r_{k-1})^T r_k}{(r_k - r_{k-1})^T p_{k-1}} \stackrel{\text{Satz (3)}}{=} \frac{r_k^T r_k}{r_{k-1}^T r_{k-1}}$$

Folgerungen aus dem SATZ (S.57):

- Die Suchrichtungen  $p_0, p_1, \dots$  sind paarweise konjugiert
- Die Residuen (Gradienten)  $r_0, r_1, \dots$  sind paarweise orthogonal.

$\Rightarrow r_0 \in \mathbb{R}^n, r_1 \in \mathbb{R}^n, \dots, r_{n-1} \in \mathbb{R}^n$   
 bilden vollständige orthog. Basis in  $\mathbb{R}^n$

$\Rightarrow r_n = 0$  (da  $\in \mathbb{R}^n$  und  $\perp \mathbb{R}^n$ )

$\Rightarrow \underline{x}_n$  ist exakte Lösung,  $A \underline{x}_n + \underline{b} = 0$

Praktische Aspekte

(i) Wegen Rundungsfehlern bei numerischer Rechnung wird die volle Genauigkeit erst nach ein paar zusätzlichen Schritten erreicht. Einfach weiter rechnen!

(ii) Zielprobleme:  $n > 10'000$

Durchführung von  $\approx n$  Iterationsschritten kommt nicht in Frage!

Das Wunder:

Bei guter Kondition von A genügen schon  $\sqrt{n}$  Iterationsschritte!

(sonst wäre der cg-Algorithmus längst vergessen).

(c) Vorkonditionierung

Def. Sei  $\underline{x} \in \mathbb{R}^n$ ,  $A \in \mathbb{R}^{n \times n}$ , symm. pos. def  
A-Norm von  $\underline{x}$ :  $\|\underline{x}\|_A := \sqrt{\underline{x}^T A \underline{x}}$

SATZ (Konvergenz des cg-Algorithmus S.56):

$$\|\underline{x}_k - \underline{x}\|_A \leq 2 \mu^k \|\underline{x}_0 - \underline{x}\|_A$$

exakte Lösung,  
 $A\underline{x} + \underline{b} = 0$

mit  $\mu = \frac{\sqrt{\text{cond}_2(A)} - 1}{\sqrt{\text{cond}_2(A)} + 1} \quad (< 1)$

Idee der Vorkonditionierung (preconditioning)

Betrachte neues Gleichungssystem

(1)  $\tilde{A} \tilde{\underline{x}} + \tilde{\underline{b}} = 0$

mit

später:  $\tilde{A} = C^{-T} A C^{-1}$

$$\left. \begin{aligned} \tilde{\underline{p}} = C \underline{p}, \quad \tilde{\underline{z}} = C \underline{z}, \quad \tilde{\underline{x}} = C \underline{x} \\ \tilde{\underline{r}} = C^{-T} \underline{r}, \quad \tilde{\underline{b}} = C^{-T} \underline{b} \end{aligned} \right\} (2)$$

Notation:  
 $C^{-T} = (C^{-1})^T$   
 $= (C^T)^{-1}$

Das Gl'system (1) mit den Trsf. (2) ist äquivalent zu  $A \underline{x} + \underline{b} = 0$ .

Ziel: Wahl der regulären Matrix C so, dass  $\text{cond}_2(\tilde{A}) \approx 1$ .

Vorbereitung: cg-Algorithmus S.56 neu schreiben mit zusätzlicher Vektorfolge  $\underline{z}_k := \underline{r}_k$

# cg-Algorithmus mit Vorkonditionierung

zur Lösung von  $A\underline{x} + \underline{b} = 0$

Initialisierung:  $\underline{x}_0 = 0, \underline{r}_0 = \underline{b}, \underline{p}_1 = 0$

for  $k = 0, 1, \dots$

(3)

~~$\underline{z}_k = \underline{r}_k$~~   $\xrightarrow{\text{ersetzen durch}}$

Berechne  $\underline{z}_k$  aus  
 $M \underline{z}_k = \underline{r}_k$

$$\sigma_k = \begin{cases} 0, & k = 0 \\ \underline{z}_k^T \underline{r}_k / \underline{z}_{k-1}^T \underline{r}_{k-1}, & k > 0 \end{cases}$$

$$\underline{p}_k = \sigma_k \underline{p}_{k-1} - \underline{z}_k$$

$$\rho_k = \underline{z}_k^T \underline{r}_k / \underline{p}_k^T A \underline{p}_k$$

$$\underline{x}_{k+1} = \underline{x}_k + \rho_k \underline{p}_k$$

$$\underline{r}_{k+1} = \underline{r}_k + \rho_k A \underline{p}_k$$

Unser Vorgehen: (i) Überall  $\sim$  aufsetzen, d.h. wir wollen System (1) lösen  
(ii) Trsf. (2) anwenden

Resultat: Alles bleibt invariant (bitte nachrechnen!)  
ausser Gleichung (3):

$$(3) \xrightarrow{(i)} \underline{\tilde{z}}_k = \underline{\tilde{r}}_k \xrightarrow{(ii)} \underbrace{(C^T C)}_M \underline{z}_k = \underline{r}_k$$

M, "preconditioner"

Anforderungen an M:

- Gleichungssystem  $M \underline{z}_k = \underline{r}_k$  einfach lösbar
- $M \approx A$   $M = A$  würde Lösung für  $k=0$  bringen!
- $M = I$  ist cg ohne Vorkond