

# 03 CG-Verfahren

Nichtlineare Optimierung

WS 2020/21

# A-konjugierte Vektoren

A-orthogonal

$$d_i^T A d_j = 0 \text{ für } i \neq j$$

Es sei  $A$  spd. A-konjugierte Vektoren  $d_0, d_1, \dots, d_k$  sind linear unabhängig.

$d_j \neq 0$

$$\sum_{j=0}^k \alpha_j d_j = 0$$

$$\Rightarrow \underbrace{d_i^T A}_{\substack{=0 \\ \text{für } i \neq j}} \sum_{j=0}^k \alpha_j d_j = \sum_{j=0}^k \alpha_j \underbrace{d_i^T A d_j}_{=0 \text{ für } i \neq j} = 0$$

$$\Rightarrow \alpha_i \underbrace{d_i^T A d_i}_{= \|d_i\|_A^2 > 0} = 0$$

$$\Rightarrow \alpha_i = 0 \text{ für } 0 \leq i \leq k$$

# Quizfrage

Grad. verfahren



Welche besondere Eigenschaft haben Verfahren zur Minimierung von  $\frac{1}{2}x^T A x - b^T x + c$ , die mit  $A$ -konjugierten Richtungen arbeiten und in jeder Iteration die exakte eindimensionale Minimierung ausführen?

→ Umfrage

A Sie konvergieren in einem Schritt.

B Sie konvergieren in  $\dim A$  Schritten.

C  $x_k$  minimiert die Zielfunktion über  $x_0 +$

D Sie konvergieren in zwei Schritten.

*lin. Hülle*   $\text{span}\{d_0, d_1, \dots, d_{k-1}\}$

# Das CG-Verfahren

$$r = Ax - b$$

$$= \nabla \phi(x)$$

$$M^{-1}r = \nabla_n \phi(x)$$

Verwende die Richtungen

$$\blacktriangleright d_0 := -M^{-1}r_0$$

wie im Gradientenver.

$$\blacktriangleright d_k := -M^{-1}r_k + \beta_k d_{k-1} \quad \text{für } k \geq 1$$

Korrektur für  $\lambda$ -Konjugiertheit

$$d_{k-1}^T A d_k = -d_{k-1}^T A M^{-1} r_k + \beta_k d_{k-1}^T A d_{k-1} \stackrel{!}{=} 0$$

$$\Rightarrow \beta_k = \frac{d_{k-1}^T A M^{-1} r_k}{d_{k-1}^T A d_{k-1}} = \dots = \frac{\|r_k\|_{M^{-1}}^2}{\|r_{k-1}\|_{M^{-1}}^2}$$

$\|r_k\|_{M^{-1}}^2$  wird sukzessive im Verfahren beseitigt

# Das CG-Verfahren

- 1: Setze  $k := 0$
- 2: Setze  $r_0 := Ax_0 - b$
- 3: Setze  $d_0 := -M^{-1}r_0$
- 4: Setze  $\delta_0 := -r_0^T d_0$  //  $\delta_0 = \|\nabla_M \phi(x_0)\|_M^2$   
*=  $\|r_0\|_{M^{-1}}^2$*
- 5: **while** Abbruchkriterium nicht erfüllt **do**
- 6:     Setze  $q_k := Ad_k$
- 7:     Setze  $\alpha_k := \delta_k / (d_k^T q_k)$
- 8:     Setze  $x_{k+1} := x_k + \alpha_k d_k$
- 9:     Setze  $r_{k+1} := r_k + \alpha_k q_k$
- 10:    Setze  $d_{k+1} := -M^{-1}r_{k+1}$
- 11:    Setze  $\delta_{k+1} := -r_{k+1}^T d_{k+1}$  //  $\delta_{k+1} = \|\nabla_M \phi(x_{k+1})\|_M^2$
- 12:    Setze  $\beta_{k+1} := \frac{\delta_{k+1}}{\delta_k}$
- 13:    Setze  $d_{k+1} := d_{k+1} + \beta_{k+1} d_k$
- 14:    Setze  $k := k + 1$
- 15: **end while**
- 16: **return**  $x_k$

# Mitrechnen von Funktionswerten

$$r_0 = Ax_0 - b$$

vermeide zusätzliche Matrix-  
vektor - Multipl. mit  $x_k$

$$\phi(x_0) = \frac{1}{2} x_0^T \underbrace{A x_0}_{\downarrow} - b^T x_0 + c$$

$$= \frac{1}{2} x_0^T r_0 - \frac{1}{2} x_0^T b + c$$

$$= \frac{1}{2} x_0^T (r_0 - b) + c \quad \leftarrow \text{allgemeiner: für } x_k, r_k$$

$$\phi(x_{k+1}) - \phi(x_k) \stackrel{(4.8)}{=} -\frac{1}{2} \frac{(d_k^T r_k)^2}{d_k^T A d_k}$$

$$= -\frac{1}{2} (d_k^T r_k) \frac{d_k^T r_k}{d_k^T A d_k}$$

$$= -\frac{1}{2} \delta_k \alpha_k$$



# Quizfrage

Welche Bausteine muss der Anwender zur Verfügung stellen, um das CG-Verfahren anwenden zu können?

→ Umfrage

- A die Matrizen  $A$  und  $M$        B die Matrizen  $A$  und  $M^{-1}$   
 C Matrix-Vektor-Produkte mit  $A$  und  $M$        D Matrix-Vektor-Produkte mit  $A$  und  $M^{-1}$

*$N$  oder Matrix-Vektor-Produkte mit  $N$  werden nicht benötigt*

# Quizfrage

Welche Eigenschaft hat das CG-Verfahren?

*solange nicht  $x_k = x^*$*

→ Umfrage

0 (A)  $\phi(x_k)$  fällt streng  
monoton

5 (C) benötigt ein  
Matrix-Vektor-Produkt  
mit  $A$  und eines mit  
 $M^{-1}$  pro Iteration

3 (B) die erste Iterierte  $x_1$  ist  
identisch zu der des  
Gradientenverfahrens

2 (D)  ~~$\|x_k - x_0\|$~~   $M$  ist streng  
monoton wachsend

*Lemma 4.20*

*Berechnung über (4.35)*

~~$\|x_k - x_0\|$~~

# Quizfrage

Welche Rolle spielt der Vorkonditionierer  $M$ ?

→ Umfrage

hoffentlich !

1 **A** legt das  $M$ -Innenprodukt im Raum  $\mathbb{R}^n$  der Optimierungsvariablen fest

~~4~~ **B** verbessert die Konditionszahl  $\frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$  zu

3 **C** legt das  $M^{-1}$ -Innenprodukt im Raum  $\mathbb{R}^n$  der Residuen fest

$$\frac{\lambda_{\max}(A; M)}{\lambda_{\min}(A; M)}$$

2 **D** bestimmt die affinen Unterräume, in denen  $\phi$  sukzessive minimiert wird

$$x_{k+1} = x_k + \alpha_k d_k = \dots = x_0 + \sum_{j=0}^k \alpha_j d_j$$

$$\hookrightarrow \in \text{span} \{d_0, \dots, d_k\} = M^{-1} \text{span} \{r_0, (AM^{-1})r_0, \dots, (AM^{-1})^k r_0\}$$

# Demonstration CG-Verfahren

`test_cg_quadratic.m`

# Zeit für Ihre Fragen

Was sind Ihre Fragen zu den Themen der Woche?

→ Benutzen Sie den **Chat**.

# Fragen und Antworten 1

Warum wird beim Training von neuronalen Netzen oft ein (stochastisches) Gradientenverfahren verwendet, wenn doch das (nichtlineare) CG-Verfahren bessere Eigenschaften hat?

Die Eigenschaften, die wir bisher verglichen haben, beziehen sich auf das Gradienten- und CG-Verfahren für quadratische Aufgaben. Bei NN treten aber nicht-lineare Aufgaben auf. Hier ist der Unterschied zwischen Gradienten- und nichtlinearem CG-Verfahren nicht immer so gravierend wie bei quadratischen Funktionen. Außerdem zeigt sich, dass bei NN eine schnelle Konvergenz zu einem lokalen Min. sich oft negativ auf die Fähigkeit auswirkt, vom Trainingsdatensatz zu generalisieren.

# Fragen und Antworten 2