

# Conjugate Direction minimization

Lectures for PHD course on  
Unconstrained Numerical Optimization

Enrico Bertolazzi

DIMS – Università di Trento

May 2008

# Outline

- 1 The Steepest Descent iterative scheme
- 2 Conjugate direction method
- 3 Conjugate Gradient method
- 4 Conjugate Gradient convergence rate
- 5 Preconditioning the Conjugate Gradient method
- 6 Nonlinear Conjugate Gradient extension



# Generic minimization algorithm

In the following we study the convergence rate of the Generic minimization algorithm applied to a quadratic function  $q(\mathbf{x})$  with **exact** line search. The function

$$q(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x} + c$$

can be viewed as a  $n$ -dimensional generalization of the 1-dimensional parabolic model.

## Generic minimization algorithm

Given an initial guess  $\mathbf{x}_0$ , let  $k = 0$ ;

**while not converged do**

    Find a descent direction  $\mathbf{p}_k$  at  $\mathbf{x}_k$ ;

    Compute a step size  $\alpha_k$  using a line-search along  $\mathbf{p}_k$ .

    Set  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$  and increase  $k$  by 1.

**end while**

## Assumption (Symmetry)

The matrix  $\mathbf{A}$  is assumed to be symmetric, in fact,

$$\mathbf{A} = \mathbf{A}^{Symm} + \mathbf{A}^{Skew}$$

where

$$\mathbf{A}^{Symm} = \frac{1}{2}[\mathbf{A} + \mathbf{A}^T], \quad \mathbf{A}^{Symm} = (\mathbf{A}^{Symm})^T$$

$$\mathbf{A}^{Skew} = \frac{1}{2}[\mathbf{A} - \mathbf{A}^T], \quad \mathbf{A}^{Skew} = -(\mathbf{A}^{Skew})^T$$

moreover

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{A}^{Symm} \mathbf{x} + \mathbf{x}^T \mathbf{A}^{Skew} \mathbf{x} = \mathbf{x}^T \mathbf{A}^{Symm} \mathbf{x}$$

so that only the symmetric part of  $\mathbf{A}$  contribute to  $q(\mathbf{x})$ .

## Assumption (SPD)

The matrix  $\mathbf{A}$  is assumed to be symmetric and positive definite, in fact,

$$\nabla q(\mathbf{x})^T = \frac{1}{2}(\mathbf{A} + \mathbf{A}^T)\mathbf{x} - \mathbf{b} = \mathbf{A}\mathbf{x} - \mathbf{b}$$

and

$$\nabla^2 q(\mathbf{x}) = \frac{1}{2}(\mathbf{A} + \mathbf{A}^T) = \mathbf{A}$$

From the **sufficient** condition for a minimum we have that  $\nabla q(\mathbf{x}_\star)^T = \mathbf{0}$ , i.e.

$$\mathbf{A}\mathbf{x}_\star = \mathbf{b}$$

and  $\nabla^2 q(\mathbf{x}_\star) = \mathbf{A}$  is SPD.

- In the following we study the convergence rate of the Steepest Descent and Conjugate Gradient methods applied to

$$q(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} - \mathbf{b}^T \mathbf{x} + c$$

where  $\mathbf{A}$  is an SPD matrix.

- This assumption simplify the analysis but it is also useful in the non linear case. In fact, by expanding a generic function  $f(\mathbf{x})$  near its minimum  $\mathbf{x}_*$  we have

$$\begin{aligned} f(\mathbf{x}) &= f(\mathbf{x}_*) + \nabla f(\mathbf{x}_*)(\mathbf{x} - \mathbf{x}_*) \\ &\quad + \frac{1}{2}(\mathbf{x} - \mathbf{x}_*)^T \nabla^2 f(\mathbf{x}_*)(\mathbf{x} - \mathbf{x}_*) + \mathcal{O}(\|\mathbf{x} - \mathbf{x}_*\|^3) \end{aligned}$$



- By setting

$$\mathbf{A} = \nabla^2 f(\mathbf{x}_*),$$

$$\mathbf{b} = \nabla^2 f(\mathbf{x}_*)\mathbf{x}_* - \nabla f(\mathbf{x}_*)$$

$$c = f(\mathbf{x}_*) - \nabla f(\mathbf{x}_*)\mathbf{x}_* + \frac{1}{2}\mathbf{x}_*^T \nabla^2 f(\mathbf{x}_*)\mathbf{x}_*$$

we have

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} - \mathbf{b}^T \mathbf{x} + c + \mathcal{O}(\|\mathbf{x} - \mathbf{x}_*\|^3)$$

- So that we expect that when an iterate  $\mathbf{x}_k$  is near  $\mathbf{x}_*$  then we can neglect  $\mathcal{O}(\|\mathbf{x} - \mathbf{x}_*\|^3)$  and the asymptotic behavior is the same of the quadratic problem.



- we can rewrite the quadratic problem in many different way as follows

$$\begin{aligned}q(\mathbf{x}) &= \frac{1}{2}(\mathbf{x} - \mathbf{x}_*)^T \mathbf{A}(\mathbf{x} - \mathbf{x}_*) + c' \\ &= \frac{1}{2}(\mathbf{A}\mathbf{x} - \mathbf{b})^T \mathbf{A}^{-1}(\mathbf{A}\mathbf{x} - \mathbf{b}) + c'\end{aligned}$$

where

$$c' = c + \frac{1}{2}\mathbf{x}_*^T \mathbf{A}\mathbf{x}_*$$

- This last forms are useful in the study of the steepest descent method.





# Outline

- 1 The Steepest Descent iterative scheme
- 2 Conjugate direction method
- 3 Conjugate Gradient method
- 4 Conjugate Gradient convergence rate
- 5 Preconditioning the Conjugate Gradient method
- 6 Nonlinear Conjugate Gradient extension

# The steepest descent for quadratic functions

(1/3)

## The steepest descent minimization algorithm

Given an initial guess  $\mathbf{x}_0$ , let  $k = 0$ ;

**while not converged do**

    Choose as descent direction  $\mathbf{p}_k = -\nabla q(\mathbf{x}_k)^T = \mathbf{b} - \mathbf{A}\mathbf{x}_k$ ;

    Compute a step size  $\alpha_k$  using a line-search along  $\mathbf{p}_k$ .

    Set  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$  and increase  $k$  by 1.

**end while**

## Definition (Residual)

*The expressions*

$$\mathbf{r}(\mathbf{x}) = \mathbf{b} - \mathbf{A}\mathbf{x}, \quad \mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$$

*are called the residual. We obviously have  $\mathbf{r}(\mathbf{x}) = -\nabla q(\mathbf{x})^T$  and  $\mathbf{r}(\mathbf{x}_\star) = \mathbf{0}$ .*



## The steepest descent for quadratic functions

(2/3)

## Lemma

The solution of the minimization problem:

$$\alpha_k = \arg \min_{\alpha \geq 0} q(\mathbf{x}_k - \alpha \mathbf{r}_k) \quad \text{is} \quad \alpha_k = -\frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k}.$$

## Proof.

Because  $p(\alpha) = q(\mathbf{x}_k - \alpha \mathbf{r}_k)$  the minimum is a stationary point:

$$\begin{aligned} \frac{dp(\alpha)}{d\alpha} &= \frac{dq(\mathbf{x}_k - \alpha \mathbf{r}_k)}{d\alpha} = -\nabla q(\mathbf{x}_k - \alpha \mathbf{r}_k) \mathbf{r}_k \\ &= \mathbf{r}(\mathbf{x}_k - \alpha \mathbf{r}_k)^T \mathbf{r}_k = (\mathbf{b} - \mathbf{A}(\mathbf{x}_k - \alpha \mathbf{r}_k))^T \mathbf{r}_k \\ &= (\mathbf{r}_k + \alpha \mathbf{A} \mathbf{r}_k)^T \mathbf{r}_k = 0 \end{aligned}$$

and solving for  $\alpha$  the result follows. □



## The steepest descent for quadratic functions

(3/3)

## The steepest descent minimization algorithm

Given an initial guess  $\mathbf{x}_0$ , let  $k = 0$ ;

**while not converged do**

    Compute  $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$ ;

    Compute the step size  $\alpha_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k}$ ;

    Set  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{r}_k$  and increase  $k$  by 1.

**end while**

Or more compactly

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k} \mathbf{r}_k$$



# The steepest descent reduction step

(1/4)

The next lemma bound the reduction of  $q(\mathbf{x}_{k+1})$  by the value of  $q(\mathbf{x}_k)$ :

## Lemma

*Consider the steepest descent for quadratic function, than we have the following estimate*

$$\|\mathbf{x}_\star - \mathbf{x}_{k+1}\|_{\mathbf{A}}^2 = \|\mathbf{x}_\star - \mathbf{x}_k\|_{\mathbf{A}}^2 \left( 1 - \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{(\mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k)(\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k)} \right)$$

where

$$\|\mathbf{x}\|_{\mathbf{A}} = \sqrt{\mathbf{x}^T \mathbf{A} \mathbf{x}}$$

is the *energy norm* induced by the SPD matrix  $\mathbf{A}$ .

## The steepest descent reduction step

(2/4)

Proof.

(1/3).

We want bound  $q(\mathbf{x}_{k+1})$  by  $q(\mathbf{x}_k)$ :

$$\begin{aligned}q(\mathbf{x}_{k+1}) &= q(\mathbf{x}_k + \alpha_k \mathbf{r}_k) \\&= \frac{1}{2} (\mathbf{A}\mathbf{x}_k + \alpha_k \mathbf{A}\mathbf{r}_k - \mathbf{b})^T \mathbf{A}^{-1} (\mathbf{A}\mathbf{x}_k + \alpha_k \mathbf{A}\mathbf{r}_k - \mathbf{b}) + c' \\&= \frac{1}{2} (\alpha_k \mathbf{A}\mathbf{r}_k - \mathbf{r}_k)^T \mathbf{A}^{-1} (\alpha_k \mathbf{A}\mathbf{r}_k - \mathbf{r}_k) + c' \\&= \frac{1}{2} \mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k + \frac{1}{2} \alpha_k^2 \mathbf{r}_k^T \mathbf{A} \mathbf{r}_k - \alpha_k \mathbf{r}_k^T \mathbf{r}_k + c' \\&= q(\mathbf{x}_k) + \frac{1}{2} \alpha_k (\alpha_k \mathbf{r}_k^T \mathbf{A} \mathbf{r}_k - 2\mathbf{r}_k^T \mathbf{r}_k)\end{aligned}$$



## The steepest descent reduction step

(3/4)

Proof.

(2/3).

Substituting  $\alpha_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k}$  we obtain

$$q(\mathbf{x}_{k+1}) = q(\mathbf{x}_k) - \frac{1}{2} \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k}$$

this shows that the steepest descent method reduce at each step the objective function  $q(\mathbf{x})$ .

Using the expression  $q(\mathbf{x}) = \frac{1}{2} \mathbf{r}(\mathbf{x})^T \mathbf{A}^{-1} \mathbf{r}(\mathbf{x}) + c'$  we can write:

$$\frac{1}{2} \mathbf{r}_{k+1}^T \mathbf{A}^{-1} \mathbf{r}_{k+1} = \frac{1}{2} \mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k - \frac{1}{2} \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k}$$



## The steepest descent reduction step

(4/4)

Proof.

(3/3).

or better

$$\mathbf{r}_{k+1}^T \mathbf{A}^{-1} \mathbf{r}_{k+1} = \mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k \left( 1 - \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{(\mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k)(\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k)} \right)$$

noticing that  $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k = \mathbf{A}\mathbf{x}_\star - \mathbf{A}\mathbf{x}_k = \mathbf{A}(\mathbf{x}_\star - \mathbf{x}_k)$  we have

$$\|\mathbf{x}_\star - \mathbf{x}_{k+1}\|_{\mathbf{A}}^2 = \|\mathbf{x}_\star - \mathbf{x}_k\|_{\mathbf{A}}^2 \left( 1 - \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{(\mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k)(\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k)} \right)$$

where

$$\|\mathbf{x}\|_{\mathbf{A}} = \sqrt{\mathbf{x}^T \mathbf{A} \mathbf{x}}$$

is the **energy norm** induced by the SPD matrix  $\mathbf{A}$ .





The estimate of the convergence rate for the **steepest descent** method is linked to the estimate of the term

$$\frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{(\mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k)(\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k)}$$

in particular we can prove

### Lemma (Kantorovic)

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  an SPD matrix then the following inequality is valid

$$1 \leq \frac{(\mathbf{x}^T \mathbf{A} \mathbf{x})(\mathbf{x}^T \mathbf{A}^{-1} \mathbf{x})}{(\mathbf{x}^T \mathbf{x})^2} \leq \frac{(M + m)^2}{4 M m}$$

for all  $\mathbf{x} \neq \mathbf{0}$ . Where  $m = \lambda_1$  is the smallest eigenvalue of  $\mathbf{A}$  and  $M = \lambda_n$  is the biggest eigenvalue of  $\mathbf{A}$ .



Proof.

(1/5).

**STEP 1: problem reformulation.** First of all notice that

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

for all  $\mathbf{y} = \alpha \mathbf{x}$  with  $\alpha \neq 0$ . Choosing  $\alpha = \|\mathbf{x}\|^{-1}$  have:

$$\min_{\|\mathbf{z}\|=1} (\mathbf{z}^T \mathbf{A} \mathbf{z})(\mathbf{z}^T \mathbf{A}^{-1} \mathbf{z}) \leq$$

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

$$\leq \max_{\|\mathbf{z}\|=1} (\mathbf{z}^T \mathbf{A} \mathbf{z})(\mathbf{z}^T \mathbf{A}^{-1} \mathbf{z})$$



Proof.

(2/5).

**STEP 2: eigenvector expansions.** Matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is an SPD matrix so that there exists  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$  a complete orthonormal eigenvectors set with  $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  corresponding eigenvalues. Let be  $\mathbf{x} \in \mathbb{R}^n$  then

$$\mathbf{x} = \sum_{k=1}^n \alpha_k \mathbf{u}_k, \quad \mathbf{x}^T \mathbf{x} = \sum_{k=1}^n \alpha_k^2$$

so that  $(\mathbf{x}^T \mathbf{A} \mathbf{x})(\mathbf{x}^T \mathbf{A}^{-1} \mathbf{x}) = h(\alpha_1, \dots, \alpha_n)$  where

$$h(\alpha_1, \dots, \alpha_n) = \left( \sum_{k=1}^n \alpha_k^2 \lambda_k \right) \left( \sum_{k=1}^n \alpha_k^2 \lambda_k^{-1} \right)$$

then the lemma can be reformulated:

- Find maxima and minima of  $h(\alpha_1, \dots, \alpha_n)$
- subject to  $\sum_{k=1}^n \alpha_k^2 = 1$ .



Proof.

(3/5).

**STEP 3: problem reduction.** By using Lagrange multiplier maxima and minima are the stationary points of:

$$g(\alpha_1, \dots, \alpha_n, \mu) = h(\alpha_1, \dots, \alpha_n) + \mu \left( \sum_{k=1}^n \alpha_k^2 - 1 \right)$$

setting  $A = \sum_{k=1}^n \alpha_k^2 \lambda_k$  and  $B = \sum_{k=1}^n \alpha_k^2 \lambda_k^{-1}$  we have

$$\frac{\partial g(\alpha_1, \dots, \alpha_n, \mu)}{\partial \alpha_k} = 2\alpha_k (\lambda_k B + \lambda_k^{-1} A + \mu) = 0$$

so that

- ① Or  $\alpha_k = 0$ ;
- ② Or  $\lambda_k$  is a root of the quadratic polynomial  $\lambda^2 B + \lambda \mu + A$ .

in any case there are at most 2 coefficients  $\alpha$ 's not zero. <sup>a</sup>

---

<sup>a</sup>the argument should be improved in the case of multiple eigenvalues

Proof.

(4/5).

**STEP 4: problem reformulation.** say  $\alpha_i$  and  $\alpha_j$  are the only non zero coefficients, then  $\alpha_i^2 + \alpha_j^2 = 1$  and we can write

$$\begin{aligned}h(\alpha_1, \dots, \alpha_n) &= (\alpha_i^2 \lambda_i + \alpha_j^2 \lambda_j)(\alpha_i^2 \lambda_i^{-1} + \alpha_j^2 \lambda_j^{-1}) \\&= \alpha_i^4 + \alpha_j^4 + \alpha_i^2 \alpha_j^2 \left( \frac{\lambda_i}{\lambda_j} + \frac{\lambda_j}{\lambda_i} \right) \\&= \alpha_i^2(1 - \alpha_j^2) + \alpha_j^2(1 - \alpha_i^2) + \alpha_i^2 \alpha_j^2 \left( \frac{\lambda_i}{\lambda_j} + \frac{\lambda_j}{\lambda_i} \right) \\&= 1 + \alpha_i^2 \alpha_j^2 \left( \frac{\lambda_i}{\lambda_j} + \frac{\lambda_j}{\lambda_i} - 2 \right) \\&= 1 + \alpha_i^2(1 - \alpha_i^2) \frac{(\lambda_i - \lambda_j)^2}{\lambda_i \lambda_j}\end{aligned}$$



Proof.

(5/5).

**STEP 5: bounding maxima and minima.** notice that

$$0 \leq \beta(1 - \beta) \leq \frac{1}{4}, \quad \forall \beta \in [0, 1]$$

$$1 \leq 1 + \alpha_i^2(1 - \alpha_i^2) \frac{(\lambda_i - \lambda_j)^2}{\lambda_i \lambda_j} \leq 1 + \frac{(\lambda_i - \lambda_j)^2}{4\lambda_i \lambda_j} = \frac{(\lambda_i + \lambda_j)^2}{4\lambda_i \lambda_j}$$

to bound  $(\lambda_i + \lambda_j)^2 / (4\lambda_i \lambda_j)$  consider the function

$f(x) = (1 + x)^2 / x$  which is increasing for  $x \geq 1$  so that we have

$$\frac{(\lambda_i + \lambda_j)^2}{4\lambda_i \lambda_j} \leq \frac{(M + m)^2}{4 M m}$$

and finally

$$1 \leq h(\alpha_1, \dots, \alpha_n) \leq \frac{(M + m)^2}{4 M m}$$



# Convergence rate of Steepest Descent

The Kantorovich inequality permits to prove:

## Theorem (Convergence rate of Steepest Descent)

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  an SPD matrix then the *steepest descent* method:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k} \mathbf{r}_k$$

converge to the solution  $\mathbf{x}_* = \mathbf{A}^{-1} \mathbf{b}$  with at least linear  $q$ -rate in the norm  $\|\cdot\|_{\mathbf{A}}$ . Moreover we have the error estimate

$$\|\mathbf{x}_{k+1} - \mathbf{x}_*\|_{\mathbf{A}} \leq \frac{\kappa - 1}{\kappa + 1} \|\mathbf{x}_k - \mathbf{x}_*\|_{\mathbf{A}}$$

$\kappa = M/m$  is the *condition number* where  $m = \lambda_1$  is the smallest eigenvalue of  $\mathbf{A}$  and  $M = \lambda_n$  is the biggest eigenvalue of  $\mathbf{A}$ .

## Proof.

Remember from slide  $N^{\circ}16$

$$\|\mathbf{x}_{\star} - \mathbf{x}_{k+1}\|_{\mathbf{A}}^2 = \|\mathbf{x}_{\star} - \mathbf{x}_k\|_{\mathbf{A}}^2 \left( 1 - \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{(\mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k)(\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k)} \right)$$

from Kantorovich inequality

$$1 - \frac{(\mathbf{r}_k^T \mathbf{r}_k)^2}{(\mathbf{r}_k^T \mathbf{A}^{-1} \mathbf{r}_k)(\mathbf{r}_k^T \mathbf{A} \mathbf{r}_k)} \leq 1 - \frac{4 M m}{(M + m)^2} = \frac{(M - m)^2}{(M + m)^2}$$

so that

$$\|\mathbf{x}_{\star} - \mathbf{x}_{k+1}\|_{\mathbf{A}} \leq \frac{M - m}{M + m} \|\mathbf{x}_{\star} - \mathbf{x}_k\|_{\mathbf{A}}$$





## Remark (One step convergence)

The steepest descent method can converge in one iteration if  $\kappa = 1$  or when  $\mathbf{r}_0 = \mathbf{u}_k$  where  $\mathbf{u}_k$  is an eigenvector of  $\mathbf{A}$ .

- 1 In the first case ( $\kappa = 1$ ) we have  $\mathbf{A} = \beta \mathbf{I}$  for some  $\beta > 0$  so it is not interesting.
- 2 In the second case we have

$$\frac{(\mathbf{u}_k^T \mathbf{u}_k)^2}{(\mathbf{u}_k^T \mathbf{A}^{-1} \mathbf{u}_k)(\mathbf{u}_k^T \mathbf{A} \mathbf{u}_k)} = \frac{(\mathbf{u}_k^T \mathbf{u}_k)^2}{\lambda_k^{-1}(\mathbf{u}_k^T \mathbf{u}_k) \lambda_k(\mathbf{u}_k^T \mathbf{u}_k)} = 1$$

in both cases we have  $\mathbf{r}_1 = \mathbf{0}$  i.e. we have found the solution.

# Outline

- 1 The Steepest Descent iterative scheme
- 2 Conjugate direction method**
- 3 Conjugate Gradient method
- 4 Conjugate Gradient convergence rate
- 5 Preconditioning the Conjugate Gradient method
- 6 Nonlinear Conjugate Gradient extension

# Conjugate direction method

## Definition (Conjugate vector)

Given two vectors  $\mathbf{p}$  and  $\mathbf{q}$  in  $\mathbb{R}^n$  are *conjugate respect to  $\mathbf{A}$*  if they are orthogonal respect the scalar product induced by  $\mathbf{A}$ ; i.e.,

$$\mathbf{p}^T \mathbf{A} \mathbf{q} = \sum_{i,j=1}^n A_{ij} p_i q_j = 0.$$

Clearly,  $n$  vectors  $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n \in \mathbb{R}^n$  that are pair wise conjugated respect to  $\mathbf{A}$  form a base of  $\mathbb{R}^n$ .



## Problem (Linear system)

Find the minimum of  $q(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} - \mathbf{b}^T \mathbf{x} + c$  is equivalent to solve the first order necessary condition, i.e.

$$\text{Find } \mathbf{x}_\star \in \mathbb{R}^n \text{ such that: } \mathbf{A}\mathbf{x}_\star = \mathbf{b}.$$

## Observation

Consider  $\mathbf{x}_0 \in \mathbb{R}^n$  and decompose the error  $\mathbf{e}_0 = \mathbf{x}_\star - \mathbf{x}_0$  by the conjugate vectors  $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n \in \mathbb{R}^n$ :

$$\mathbf{e}_0 = \mathbf{x}_\star - \mathbf{x}_0 = \sigma_1 \mathbf{p}_1 + \sigma_2 \mathbf{p}_2 + \dots + \sigma_n \mathbf{p}_n.$$

Evaluating the coefficients  $\sigma_1, \sigma_2, \dots, \sigma_n \in \mathbb{R}$  is equivalent to solve the problem  $\mathbf{A}\mathbf{x}_\star = \mathbf{b}$ , because knowing  $\mathbf{e}_0$  we have

$$\mathbf{x}_\star = \mathbf{x}_0 + \mathbf{e}_0.$$

## Observation

Using conjugacy the coefficients  $\sigma_1, \sigma_2, \dots, \sigma_n \in \mathbb{R}$  can be computed as

$$\sigma_i = \frac{\mathbf{p}_i^T \mathbf{A} \mathbf{e}_0}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i}, \quad \text{for } i = 1, 2, \dots, n.$$

In fact, for all  $1 \leq i \leq n$ , we have

$$\begin{aligned} \mathbf{p}_i^T \mathbf{A} \mathbf{e}_0 &= \mathbf{p}_i^T \mathbf{A} (\sigma_1 \mathbf{p}_1 + \sigma_2 \mathbf{p}_2 + \dots + \sigma_n \mathbf{p}_n), \\ &= \sigma_1 \mathbf{p}_i^T \mathbf{A} \mathbf{p}_1 + \sigma_2 \mathbf{p}_i^T \mathbf{A} \mathbf{p}_2 + \dots + \sigma_n \mathbf{p}_i^T \mathbf{A} \mathbf{p}_n, \\ &= \sigma_i \mathbf{p}_i^T \mathbf{A} \mathbf{p}_i, \end{aligned}$$

because  $\mathbf{p}_i^T \mathbf{A} \mathbf{p}_j = 0$  for  $i \neq j$ .

The conjugate direction method evaluate the coefficients  $\sigma_1, \sigma_2, \dots, \sigma_n \in \mathbb{R}$  recursively in  $n$  steps, solving for  $k \geq 0$  the minimization problem:

## Conjugate direction method

Given  $\mathbf{x}_0$ ;  $k \leftarrow 0$ ;

**repeat**

$k \leftarrow k + 1$ ;

Find  $\mathbf{x}_k \in \mathbf{x}_0 + \mathcal{V}_k$  such that:

$$\mathbf{x}_k = \arg \min_{\mathbf{x} \in \mathbf{x}_0 + \mathcal{V}_k} \|\mathbf{x}_* - \mathbf{x}\|_A$$

**until**  $k = n$

where  $\mathcal{V}_k$  is the subspace of  $\mathbb{R}^n$  generated by the first  $k$  conjugate direction; i.e.,

$$\mathcal{V}_k = \text{SPAN}\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k\}.$$

# Step: $\mathbf{x}_0 \rightarrow \mathbf{x}_1$

At the first step we consider the subspace  $\mathbf{x}_0 + \text{SPAN}\{\mathbf{p}_1\}$  which consists in vectors of the form

$$\mathbf{x}(\alpha) = \mathbf{x}_0 + \alpha\mathbf{p}_1 \quad \alpha \in \mathbb{R}$$

The minimization problem becomes:

## Minimization step $\mathbf{x}_0 \rightarrow \mathbf{x}_1$

Find  $\mathbf{x}_1 = \mathbf{x}_0 + \alpha_1\mathbf{p}_1$  (i.e., find  $\alpha_1!$ ) such that:

$$\|\mathbf{x}_\star - \mathbf{x}_1\|_{\mathbf{A}} = \min_{\alpha \in \mathbb{R}} \|\mathbf{x}_\star - (\mathbf{x}_0 + \alpha\mathbf{p}_1)\|_{\mathbf{A}},$$

# Solving first step method 1

The minimization problem is the minimum respect to  $\alpha$  of the quadratic:

$$\begin{aligned}\Phi(\alpha) &= \|\mathbf{x}_* - (\mathbf{x}_0 + \alpha\mathbf{p}_1)\|_{\mathbf{A}}^2, \\ &= (\mathbf{x}_* - (\mathbf{x}_0 + \alpha\mathbf{p}_1))^T \mathbf{A} (\mathbf{x}_* - (\mathbf{x}_0 + \alpha\mathbf{p}_1)), \\ &= (\mathbf{e}_0 - \alpha\mathbf{p}_1)^T \mathbf{A} (\mathbf{e}_0 - \alpha\mathbf{p}_1), \\ &= \mathbf{e}_0^T \mathbf{A} \mathbf{e}_0 - 2\alpha \mathbf{p}_1^T \mathbf{A} \mathbf{e}_0 + \alpha^2 \mathbf{p}_1^T \mathbf{A} \mathbf{p}_1.\end{aligned}$$

minimum is found by imposing:

$$\frac{d\Phi(\alpha)}{d\alpha} = -2\mathbf{p}_1^T \mathbf{A} \mathbf{e}_0 + 2\alpha \mathbf{p}_1^T \mathbf{A} \mathbf{p}_1 = 0 \quad \Rightarrow$$

$$\alpha_1 = \frac{\mathbf{p}_1^T \mathbf{A} \mathbf{e}_0}{\mathbf{p}_1^T \mathbf{A} \mathbf{p}_1}$$





## Solving first step method 2

(1/2)

Remember the error expansion:

$$\mathbf{x}_\star - \mathbf{x}_0 = \sigma_1 \mathbf{p}_1 + \sigma_2 \mathbf{p}_2 + \cdots + \sigma_n \mathbf{p}_n.$$

Let  $\mathbf{x}(\alpha) = \mathbf{x}_0 + \alpha \mathbf{p}_1$ , the difference  $\mathbf{x}_\star - \mathbf{x}(\alpha)$  becomes:

$$\mathbf{x}_\star - \mathbf{x}(\alpha) = (\sigma_1 - \alpha) \mathbf{p}_1 + \sigma_2 \mathbf{p}_2 + \cdots + \sigma_n \mathbf{p}_n$$

due to conjugacy the error  $\|\mathbf{x}_\star - \mathbf{x}(\alpha)\|_{\mathbf{A}}$  becomes

$$\begin{aligned} & \|\mathbf{x}_\star - \mathbf{x}(\alpha)\|_{\mathbf{A}}^2 \\ &= \left( (\sigma_1 - \alpha) \mathbf{p}_1 + \sum_{i=2}^n \sigma_i \mathbf{p}_i \right)^T \mathbf{A} \left( (\sigma_1 - \alpha) \mathbf{p}_1 + \sum_{j=2}^n \sigma_j \mathbf{p}_j \right) \\ &= (\sigma_1 - \alpha)^2 \mathbf{p}_1^T \mathbf{A} \mathbf{p}_1 + \sum_{j=2}^n \sigma_j^2 \mathbf{p}_j^T \mathbf{A} \mathbf{p}_j \end{aligned}$$



## Solving first step method 2

(2/2)

Because

$$\|\mathbf{x}_* - \mathbf{x}(\alpha)\|_{\mathbf{A}}^2 = (\sigma_1 - \alpha)^2 \|\mathbf{p}_1\|_{\mathbf{A}}^2 + \sum_{i=2}^n \sigma_i^2 \|\mathbf{p}_i\|_{\mathbf{A}}^2,$$

we have that

$$\|\mathbf{x}_* - \mathbf{x}(\alpha_1)\|_{\mathbf{A}}^2 = \sum_{i=2}^n \sigma_i^2 \|\mathbf{p}_i\|_{\mathbf{A}}^2 \leq \|\mathbf{x}_* - \mathbf{x}(\alpha)\|_{\mathbf{A}}^2 \quad \text{for all } \alpha \neq \sigma_1$$

so that minimum is found by imposing  $\alpha_1 = \sigma_1$ :

$$\alpha_1 = \frac{\mathbf{p}_1^T \mathbf{A} \mathbf{e}_0}{\mathbf{p}_1^T \mathbf{A} \mathbf{p}_1}$$

This argument can be generalized for all  $k > 1$  (see next slides).

# Step, $\mathbf{x}_{k-1} \rightarrow \mathbf{x}_k$

For the step from  $k - 1$  to  $k$  we consider the subspace of  $\mathbb{R}^n$

$$\mathcal{V}_k = \text{SPAN}\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k\}$$

which contains vectors of the form:

$$\mathbf{x}(\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(k)}) = \mathbf{x}_0 + \alpha^{(1)}\mathbf{p}_1 + \alpha^{(2)}\mathbf{p}_2 + \dots + \alpha^{(k)}\mathbf{p}_k$$

The minimization problem becomes:

## Minimization step $\mathbf{x}_{k-1} \rightarrow \mathbf{x}_k$

Find  $\mathbf{x}_k = \mathbf{x}_0 + \alpha_1\mathbf{p}_1 + \alpha_2\mathbf{p}_2 + \dots + \alpha_k\mathbf{p}_k$  (i.e.  $\alpha_1, \alpha_2, \dots, \alpha_k$ ) such that:

$$\|\mathbf{x}_\star - \mathbf{x}_k\|_{\mathbf{A}} = \min_{\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(k)} \in \mathbb{R}} \left\| \mathbf{x}_\star - \mathbf{x}(\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(k)}) \right\|_{\mathbf{A}}$$

Solving  $k$ th Step:  $\mathbf{x}_{k-1} \rightarrow \mathbf{x}_k$ 

(1/2)

Remember the error expansion:

$$\mathbf{x}_\star - \mathbf{x}_0 = \sigma_1 \mathbf{p}_1 + \sigma_2 \mathbf{p}_2 + \cdots + \sigma_n \mathbf{p}_n.$$

Consider a vector of the form

$$\mathbf{x}(\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(k)}) = \mathbf{x}_0 + \alpha^{(1)} \mathbf{p}_1 + \alpha^{(2)} \mathbf{p}_2 + \cdots + \alpha^{(k)} \mathbf{p}_k$$

the error  $\mathbf{x}_\star - \mathbf{x}(\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(k)})$  can be written as

$$\begin{aligned} \mathbf{x}_\star - \mathbf{x}(\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(k)}) &= \mathbf{x}_\star - \mathbf{x}_0 - \sum_{i=1}^k \alpha^{(i)} \mathbf{p}_i, \\ &= \sum_{i=1}^k (\sigma_i - \alpha^{(i)}) \mathbf{p}_i + \sum_{i=k+1}^n \sigma_i \mathbf{p}_i. \end{aligned}$$



Solving  $k$ th Step:  $\mathbf{x}_{k-1} \rightarrow \mathbf{x}_k$ 

(2/2)

using conjugacy of  $\mathbf{p}_i$  we obtain the norm of the error:

$$\begin{aligned} & \left\| \mathbf{x}_* - \mathbf{x}(\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(k)}) \right\|_{\mathbf{A}}^2 \\ &= \sum_{i=1}^k (\sigma_i - \alpha^{(i)})^2 \|\mathbf{p}_i\|_{\mathbf{A}}^2 + \sum_{i=k+1}^n \sigma_i^2 \|\mathbf{p}_i\|_{\mathbf{A}}^2. \end{aligned}$$

So that minimum is found by imposing  $\alpha_i = \sigma_i$ : for  $i = 1, 2, \dots, k$ .

$$\boxed{\alpha_i = \frac{\mathbf{p}_i^T \mathbf{A} \mathbf{e}_0}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i}} \quad i = 1, 2, \dots, k$$



## Successive one dimensional minimization

(1/3)

- notice that  $\alpha_i = \sigma_i$  and that

$$\begin{aligned}\mathbf{x}_k &= \mathbf{x}_0 + \alpha_1 \mathbf{p}_1 + \cdots + \alpha_k \mathbf{p}_k \\ &= \mathbf{x}_{k-1} + \alpha_k \mathbf{p}_k\end{aligned}$$

- so that  $\mathbf{x}_{k-1}$  contains  $k - 1$  coefficients  $\alpha_i$  for the minimization.
- if we consider the one dimensional minimization on the subspace  $\mathbf{x}_{k-1} + \text{SPAN}\{\mathbf{p}_k\}$  we find again  $\mathbf{x}_k$ !



## Successive one dimensional minimization

(2/3)

Consider a vector of the form

$$\mathbf{x}(\alpha) = \mathbf{x}_{k-1} + \alpha \mathbf{p}_k$$

remember that  $\mathbf{x}_{k-1} = \mathbf{x}_0 + \alpha_1 \mathbf{p}_1 + \dots + \alpha_{k-1} \mathbf{p}_{k-1}$  so that the error  $\mathbf{x}_\star - \mathbf{x}(\alpha)$  can be written as

$$\begin{aligned} \mathbf{x}_\star - \mathbf{x}(\alpha) &= \mathbf{x}_\star - \mathbf{x}_0 - \sum_{i=1}^{k-1} \alpha_i \mathbf{p}_i + \alpha \mathbf{p}_k \\ &= \sum_{i=1}^{k-1} (\sigma_i - \alpha_i) \mathbf{p}_i + (\sigma_k - \alpha) \mathbf{p}_k + \sum_{i=k+1}^n \sigma_i \mathbf{p}_i. \end{aligned}$$

due to the equality  $\sigma_i = \alpha_i$  the blue part of the expression is 0.



# Successive one dimensional minimization

(3/3)

Using conjugacy of  $\mathbf{p}_i$  we obtain the norm of the error:

$$\|\mathbf{x}_\star - \mathbf{x}(\alpha)\|_{\mathbf{A}}^2 = (\sigma_k - \alpha)^2 \|\mathbf{p}_k\|_{\mathbf{A}}^2 + \sum_{i=k+1}^n \sigma_i^2 \|\mathbf{p}_i\|_{\mathbf{A}}^2.$$

So that minimum is found by imposing  $\alpha = \sigma_k$ :

$$\alpha_k = \frac{\mathbf{p}_k^T \mathbf{A} \mathbf{e}_0}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}$$

## Remark

*This observation permit to perform the minimization on the  $k$ -dimensional space  $\mathbf{x}_0 + \mathcal{V}_k$  as successive one dimensional minimizations along the conjugate directions  $\mathbf{p}_k$ !*



## Problem (one dimensional successive minimization)

Find  $\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{p}_k$  such that:

$$\|\mathbf{x}_\star - \mathbf{x}_k\|_{\mathbf{A}} = \min_{\alpha \in \mathbb{R}} \|\mathbf{x}_\star - (\mathbf{x}_{k-1} + \alpha \mathbf{p}_k)\|_{\mathbf{A}},$$

The solution is the minimum respect to  $\alpha$  of the quadratic:

$$\begin{aligned} \Phi(\alpha) &= (\mathbf{x}_\star - (\mathbf{x}_{k-1} + \alpha \mathbf{p}_k))^T \mathbf{A} (\mathbf{x}_\star - (\mathbf{x}_{k-1} + \alpha \mathbf{p}_k)), \\ &= (\mathbf{e}_{k-1} - \alpha \mathbf{p}_k)^T \mathbf{A} (\mathbf{e}_{k-1} - \alpha \mathbf{p}_k), \\ &= \mathbf{e}_{k-1}^T \mathbf{A} \mathbf{e}_{k-1} - 2\alpha \mathbf{p}_k^T \mathbf{A} \mathbf{e}_{k-1} + \alpha^2 \mathbf{p}_k^T \mathbf{A} \mathbf{p}_k. \end{aligned}$$

minimum is found by imposing:

$$\frac{d\Phi(\alpha)}{d\alpha} = -2\mathbf{p}_k^T \mathbf{A} \mathbf{e}_{k-1} + 2\alpha \mathbf{p}_k^T \mathbf{A} \mathbf{p}_k = 0 \quad \Rightarrow$$

$$\alpha_k = \frac{\mathbf{p}_k^T \mathbf{A} \mathbf{e}_{k-1}}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}$$



- In the case of minimization on the subspace  $\mathbf{x}_0 + \mathcal{V}_k$  we have:

$$\alpha_k = \mathbf{p}_k^T \mathbf{A} \mathbf{e}_0 / \mathbf{p}_k^T \mathbf{A} \mathbf{p}_k$$

- In the case of one dimensional minimization on the subspace  $\mathbf{x}_{k-1} + \text{SPAN}\{\mathbf{p}_k\}$  we have:

$$\alpha_k = \mathbf{p}_k^T \mathbf{A} \mathbf{e}_{k-1} / \mathbf{p}_k^T \mathbf{A} \mathbf{p}_k$$

- Apparently they are different results, however by using the conjugacy of the vectors  $\mathbf{p}_i$  we have

$$\begin{aligned} \mathbf{p}_k^T \mathbf{A} \mathbf{e}_{k-1} &= \mathbf{p}_k^T \mathbf{A} (\mathbf{x}_\star - \mathbf{x}_{k-1}) \\ &= \mathbf{p}_k^T \mathbf{A} (\mathbf{x}_\star - (\mathbf{x}_0 + \alpha_1 \mathbf{p}_1 + \cdots + \alpha_{k-1} \mathbf{p}_{k-1})) \\ &= \mathbf{p}_k^T \mathbf{A} \mathbf{e}_0 - \alpha_1 \mathbf{p}_k^T \mathbf{A} \mathbf{p}_1 - \cdots - \alpha_{k-1} \mathbf{p}_k^T \mathbf{A} \mathbf{p}_{k-1} \\ &= \mathbf{p}_k^T \mathbf{A} \mathbf{e}_0 \end{aligned}$$



- The **one step minimization** in the space  $\mathbf{x}_0 + \mathcal{V}_n$  and the **successive minimization** in the space  $\mathbf{x}_{k-1} + \text{SPAN}\{\mathbf{p}_k\}$ ,  $k = 1, 2, \dots, n$  are equivalent if  $\mathbf{p}_i$ s are conjugate.
- The successive minimization is useful when  $\mathbf{p}_i$ s are not known in advance but must be computed as the minimization process proceeds.
- The evaluation of  $\alpha_k$  is apparently not computable because  $\mathbf{e}_i$  is not known. However noticing

$$\mathbf{A}\mathbf{e}_k = \mathbf{A}(\mathbf{x}_\star - \mathbf{x}_k) = \mathbf{b} - \mathbf{A}\mathbf{x}_k = \mathbf{r}_k$$

we can write

$$\alpha_k = \mathbf{p}_k^T \mathbf{A}\mathbf{e}_{k-1} / \mathbf{p}_k^T \mathbf{A}\mathbf{p}_k = \mathbf{p}_k^T \mathbf{r}_{k-1} / \mathbf{p}_k^T \mathbf{A}\mathbf{p}_k =$$

- Finally for the residual is valid the recurrence

$$\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k = \mathbf{b} - \mathbf{A}(\mathbf{x}_{k-1} + \alpha_k \mathbf{p}_k) = \mathbf{r}_{k-1} - \alpha_k \mathbf{A}\mathbf{p}_k.$$



# Conjugate direction minimization

## Algorithm (Conjugate direction minimization)

```

 $k \leftarrow 0; \mathbf{x}_0 \text{ assigned};$ 
 $\mathbf{r}_0 \leftarrow \mathbf{b} - \mathbf{A}\mathbf{x}_0;$ 
while not converged do
   $k \leftarrow k + 1;$ 
   $\alpha_k \leftarrow \frac{\mathbf{r}_{k-1}^T \mathbf{p}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k};$ 
   $\mathbf{x}_k \leftarrow \mathbf{x}_{k-1} + \alpha_k \mathbf{p}_k;$ 
   $\mathbf{r}_k \leftarrow \mathbf{r}_{k-1} - \alpha_k \mathbf{A} \mathbf{p}_k;$ 
end while

```

## Observation (Computational cost)

*The conjugate direction minimization requires at each step one matrix–vector product for the evaluation of  $\alpha_k$  and two update **AXPY** for  $\mathbf{x}_k$  and  $\mathbf{r}_k$ .*

# Monotonic behavior of the error

## Remark (Monotonic behavior of the error)

The *energy norm* of the error  $\|e_k\|_A$  is monotonically decreasing in  $k$ . In fact:

$$e_k = \mathbf{x}_* - \mathbf{x}_k = \alpha_{k+1}\mathbf{p}_{k+1} + \dots + \alpha_n\mathbf{p}_n,$$

and by conjugacy

$$\|e_k\|_A^2 = \|\mathbf{x}_* - \mathbf{x}_k\|_A^2 = \sigma_{k+1}^2 \|\mathbf{p}_{k+1}\|_A^2 + \dots + \sigma_n^2 \|\mathbf{p}_n\|_A^2.$$

Finally from this relation we have  $e_n = \mathbf{0}$ .

# Outline

- 1 The Steepest Descent iterative scheme
- 2 Conjugate direction method
- 3 Conjugate Gradient method**
- 4 Conjugate Gradient convergence rate
- 5 Preconditioning the Conjugate Gradient method
- 6 Nonlinear Conjugate Gradient extension

# Conjugate Gradient method

The Conjugate Gradient method combine the **Conjugate Direction** method with an **orthogonalization process** (like Gram-Schmidt) applied to the residual to construct the conjugate directions. In fact, because  $\mathbf{A}$  define a scalar product in the next slide we prove:

- each residue is orthogonal to the previous conjugate directions, and consequently linearly independent from the previous conjugate directions.
- if the residual is not null is can be used to construct a new conjugate direction.



# Orthogonality of the residue $\mathbf{r}_k$ respect $\mathcal{V}_k$

- The residue  $\mathbf{r}_k$  is orthogonal to  $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k$ . In fact, from the error expansion

$$\mathbf{e}_k = \alpha_{k+1}\mathbf{p}_{k+1} + \alpha_{k+2}\mathbf{p}_{k+2} + \dots + \alpha_n\mathbf{p}_n$$

because  $\mathbf{r}_k = \mathbf{A}\mathbf{e}_k$ , for  $i = 1, 2, \dots, k$  we have

$$\begin{aligned} \mathbf{p}_i^T \mathbf{r}_k &= \mathbf{p}_i^T \mathbf{A}\mathbf{e}_k \\ &= \mathbf{p}_i^T \mathbf{A} \sum_{j=k+1}^n \alpha_j \mathbf{p}_j = \sum_{j=k+1}^n \alpha_j \mathbf{p}_i^T \mathbf{A}\mathbf{p}_j \\ &= 0 \end{aligned}$$





## Building new conjugate direction

(1/2)

- The conjugate direction method build **one new** direction at each step.
- If  $\mathbf{r}_k \neq \mathbf{0}$  it can be used to build the new direction  $\mathbf{p}_{k+1}$  by a Gram-Schmidt orthogonalization process

$$\mathbf{p}_{k+1} = \mathbf{r}_k + \beta_1^{(k+1)} \mathbf{p}_1 + \beta_2^{(k+1)} \mathbf{p}_2 + \dots + \beta_k^{(k+1)} \mathbf{p}_k,$$

where the  $k$  coefficients  $\beta_1^{(k+1)}, \beta_2^{(k+1)}, \dots, \beta_k^{(k+1)}$  must satisfy:

$$\mathbf{p}_i^T \mathbf{A} \mathbf{p}_{k+1} = 0, \quad \text{for } i = 1, 2, \dots, k.$$



## Building new conjugate direction

(2/2)

(repeating from previous slide)

$$\mathbf{p}_{k+1} = \mathbf{r}_k + \beta_1^{(k+1)} \mathbf{p}_1 + \beta_2^{(k+1)} \mathbf{p}_2 + \cdots + \beta_k^{(k+1)} \mathbf{p}_k,$$

expanding the expression:

$$\begin{aligned} 0 &= \mathbf{p}_i^T \mathbf{A} \mathbf{p}_{k+1}, \\ &= \mathbf{p}_i^T \mathbf{A} (\mathbf{r}_k + \beta_1^{(k+1)} \mathbf{p}_1 + \beta_2^{(k+1)} \mathbf{p}_2 + \cdots + \beta_k^{(k+1)} \mathbf{p}_k), \\ &= \mathbf{p}_i^T \mathbf{A} \mathbf{r}_k + \beta_i^{(k+1)} \mathbf{p}_i^T \mathbf{A} \mathbf{p}_i, \end{aligned}$$

$$\Rightarrow \boxed{\beta_i^{(k+1)} = -\frac{\mathbf{p}_i^T \mathbf{A} \mathbf{r}_k}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i}} \quad i = 1, 2, \dots, k$$



The choice of the residual  $\mathbf{r}_k \neq \mathbf{0}$  for the construction of the new conjugate direction  $\mathbf{p}_{k+1}$  has **three** important consequences:

- ① simplification of the expression for  $\alpha_k$ ;
- ② Orthogonality of the residual  $\mathbf{r}_k$  from the previous residue  $\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{k-1}$ ;
- ③ **three point formula** and simplification of the coefficients  $\beta_i^{(k+1)}$ .

this facts will be examined in the next slides.

Simplification of the expression for  $\alpha_k$ 

Writing the expression for  $\mathbf{p}_k$  from the orthogonalization process

$$\mathbf{p}_k = \mathbf{r}_{k-1} + \beta_1^{(k+1)} \mathbf{p}_1 + \beta_2^{(k+1)} \mathbf{p}_2 + \dots + \beta_{k-1}^{(k+1)} \mathbf{p}_{k-1},$$

using orthogonality of  $\mathbf{r}_{k-1}$  and the vectors  $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_{k-1}$ , (see slide N.48) we have

$$\begin{aligned} \mathbf{r}_{k-1}^T \mathbf{p}_k &= \mathbf{r}_{k-1}^T (\mathbf{r}_{k-1} + \beta_1^{(k+1)} \mathbf{p}_1 + \beta_2^{(k+1)} \mathbf{p}_2 + \dots + \beta_{k-1}^{(k+1)} \mathbf{p}_{k-1}), \\ &= \mathbf{r}_{k-1}^T \mathbf{r}_{k-1}. \end{aligned}$$

recalling the definition of  $\alpha_k$  it follows:

$$\alpha_k = \frac{\mathbf{e}_{k-1}^T \mathbf{A} \mathbf{p}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k} = \frac{\mathbf{r}_{k-1}^T \mathbf{p}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k} = \boxed{\frac{\mathbf{r}_{k-1}^T \mathbf{r}_{k-1}}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}}$$



# Orthogonality of the residue $\mathbf{r}_k$ from $\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{k-1}$

From the definition of  $\mathbf{p}_{i+1}$  it follows:

$$\mathbf{p}_{i+1} = \mathbf{r}_i + \beta_1^{(i+1)} \mathbf{p}_1 + \beta_2^{(i+1)} \mathbf{p}_2 + \dots + \beta_i^{(i+1)} \mathbf{p}_i,$$

$$\Rightarrow \mathbf{r}_i \in \text{SPAN}\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_i, \mathbf{p}_{i+1}\} = \mathcal{V}_{i+1} \quad (\text{obvious})$$

using orthogonality of  $\mathbf{r}_k$  and the vectors  $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k$ , (see slide N.48) for  $i < k$  we have

$$\begin{aligned} \mathbf{r}_k^T \mathbf{r}_i &= \mathbf{r}_k^T \left( \mathbf{p}_{i+1} - \sum_{j=1}^i \beta_j^{(i+1)} \mathbf{p}_j \right), \\ &= \mathbf{r}_k^T \mathbf{p}_{i+1} - \sum_{j=1}^i \beta_j^{(i+1)} \mathbf{r}_k^T \mathbf{p}_j = 0. \end{aligned}$$



# Three point formula and simplification of $\beta_i^{(k+1)}$

From the relation  $\mathbf{r}_k^T \mathbf{r}_i = \mathbf{r}_k^T (\mathbf{r}_{i-1} - \alpha_i \mathbf{A} \mathbf{p}_i)$  we deduce

$$\mathbf{r}_k^T \mathbf{A} \mathbf{p}_i = \frac{\mathbf{r}_k^T \mathbf{r}_{i-1} - \mathbf{r}_k^T \mathbf{r}_i}{\alpha_i} = \begin{cases} -\mathbf{r}_k^T \mathbf{r}_k / \alpha_k & \text{if } i = k; \\ 0 & \text{if } i < k; \end{cases}$$

remembering that  $\alpha_k = \mathbf{r}_{k-1}^T \mathbf{r}_{k-1} / \mathbf{p}_k^T \mathbf{A} \mathbf{p}_k$  we obtain

$$\beta_i^{(k+1)} = -\frac{\mathbf{r}_k^T \mathbf{A} \mathbf{p}_i}{\mathbf{p}_i^T \mathbf{A} \mathbf{p}_i} = \begin{cases} \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_{k-1}^T \mathbf{r}_{k-1}} & i = k; \\ 0 & i < k; \end{cases}$$

i.e. there is only one non zero coefficient  $\beta_k^{(k+1)}$ , so we write

$\beta_k = \beta_k^{(k+1)}$  and obtain the **three point formula**:

$$\mathbf{p}_{k+1} = \mathbf{r}_k + \beta_k \mathbf{p}_k$$



# Conjugate gradient algorithm

initial step:

$k \leftarrow 0$ ;  $\mathbf{x}_0$  assigned;

$\mathbf{r}_0 \leftarrow \mathbf{b} - \mathbf{A}\mathbf{x}_0$ ;

$\mathbf{p}_1 \leftarrow \mathbf{r}_0$ ;

**while**  $\|\mathbf{r}_k\| > \epsilon$  **do**

$k \leftarrow k + 1$ ;

Conjugate direction method

$$\alpha_k \leftarrow \frac{\mathbf{r}_{k-1}^T \mathbf{r}_{k-1}}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k};$$

$$\mathbf{x}_k \leftarrow \mathbf{x}_{k-1} + \alpha_k \mathbf{p}_k;$$

$$\mathbf{r}_k \leftarrow \mathbf{r}_{k-1} - \alpha_k \mathbf{A} \mathbf{p}_k;$$

Residual orthogonalization

$$\beta_k \leftarrow \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_{k-1}^T \mathbf{r}_{k-1}};$$

$$\mathbf{p}_{k+1} \leftarrow \mathbf{r}_k + \beta_k \mathbf{p}_k;$$

**end while**

# Outline

- 1 The Steepest Descent iterative scheme
- 2 Conjugate direction method
- 3 Conjugate Gradient method
- 4 Conjugate Gradient convergence rate**
- 5 Preconditioning the Conjugate Gradient method
- 6 Nonlinear Conjugate Gradient extension



# Polynomial residual expansions

(1/6)

## Lemma

*The residuals and conjugate directions for the Conjugate Gradient iterative scheme of slide 55 can be written as*

$$\mathbf{r}_k = P_k(\mathbf{A})\mathbf{r}_0 \quad k = 0, 1, \dots, n$$

$$\mathbf{p}_k = Q_{k-1}(\mathbf{A})\mathbf{r}_0 \quad k = 1, 2, \dots, n$$

*where  $P_k(x)$  and  $Q_k(x)$  are  $k$ -degree polynomial such that  $P_k(0) = 1$  for all  $k$ .*

## Proof.

(1/2).

The proof is by induction.

**Base  $k = 0$ :**  $\mathbf{p}_1 = \mathbf{r}_0$   
so that  $P_0(x) = 1$  and  $Q_0(x) = 1$ .

## Polynomial residual expansions

(2/6)

Proof.

(2/2).

Let the expansion valid for  $k - 1$ . Consider the recursion for the residual:

$$\begin{aligned}\mathbf{r}_k &= \mathbf{r}_{k-1} - \alpha_k \mathbf{A} \mathbf{p}_k \\ &= P_{k-1}(\mathbf{A}) \mathbf{r}_0 + \alpha_k \mathbf{A} Q_{k-1}(\mathbf{A}) \mathbf{r}_0 \\ &= (P_{k-1}(\mathbf{A}) + \alpha_k \mathbf{A} Q_{k-1}(\mathbf{A})) \mathbf{r}_0\end{aligned}$$

then  $P_k(x) = P_{k-1}(x) + \alpha_k x Q_{k-1}(x)$  and  $P_k(0) = P_{k-1}(0) = 1$ . Consider the recursion for the conjugate direction

$$\begin{aligned}\mathbf{p}_{k+1} &= P_k(\mathbf{A}) \mathbf{r}_0 + \beta_k Q_{k-1}(\mathbf{A}) \mathbf{r}_0 \\ &= (P_k(\mathbf{A}) + \beta_k Q_{k-1}(\mathbf{A})) \mathbf{r}_0\end{aligned}$$

then  $Q_k(x) = P_k(x) + \beta_k Q_{k-1}(x)$ .

□



## Polynomial residual expansions

(3/6)

## Corollary

$$\mathbf{e}_k = P_k(\mathbf{A})\mathbf{e}_0.$$

## Proof.

$$\begin{aligned}\mathbf{e}_k &= \mathbf{x}_\star - \mathbf{x}_k = \mathbf{A}^{-1}\mathbf{r}_k \\ &= \mathbf{A}^{-1}P_k(\mathbf{A})\mathbf{r}_0 \\ &= P_k(\mathbf{A})\mathbf{A}^{-1}\mathbf{r}_0 \\ &= P_k(\mathbf{A})(\mathbf{x}_\star - \mathbf{x}_0) \\ &= P_k(\mathbf{A})\mathbf{e}_0.\end{aligned}$$



## Polynomial residual expansions

(4/6)

## Lemma

For the Conjugate Gradient iterative scheme of slide n.55 we have:

$$\mathcal{V}_k = \{p(\mathbf{A})\mathbf{e}_0 \mid p \in \mathbb{P}^k, p(0) = 0\}$$

## Proof.

Using expansion of slide n.57 and  $\mathbf{r}_0 = \mathbf{A}\mathbf{e}_0$  we have:

$$\begin{aligned} \mathcal{V}_k &= \text{SPAN}\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k\} \\ &= \left\{ \sum_{i=0}^{k-1} \beta_i Q_i(\mathbf{A})\mathbf{r}_0 \mid (\beta_0, \dots, \beta_{k-1}) \in \mathbb{R}^{k-1} \right\} \\ &= \{q(\mathbf{A})\mathbf{A}\mathbf{e}_0 \mid p \in \mathbb{P}^{k-1}\} = \{p(\mathbf{A})\mathbf{e}_0 \mid p \in \mathbb{P}^k, p(0) = 0\} \end{aligned}$$



## Polynomial residual expansions

(5/6)

By using the equality

$$\mathcal{V}_k = \{p(\mathbf{A})\mathbf{e}_0 \mid p \in \mathbb{P}^k, p(0) = 0\}$$

The optimality of CG step can be written as

$$\|\mathbf{x}_\star - \mathbf{x}_k\|_{\mathbf{A}} \leq \|\mathbf{x}_\star - \mathbf{x}\|_{\mathbf{A}}, \quad \forall \mathbf{x} \in \mathbf{x}_0 + \mathcal{V}_k$$

$$\|\mathbf{x}_\star - \mathbf{x}_k\|_{\mathbf{A}} \leq \|\mathbf{x}_\star - (\mathbf{x}_0 + p(\mathbf{A})\mathbf{e}_0)\|_{\mathbf{A}}, \quad \forall p \in \mathbb{P}^k, p(0) = 0$$

$$\|\mathbf{x}_\star - \mathbf{x}_k\|_{\mathbf{A}} \leq \|P(\mathbf{A})\mathbf{e}_0\|_{\mathbf{A}}, \quad \forall P \in \mathbb{P}^k, P(0) = 1$$

And using the results of slide 60 and 59 we can write

$$\mathbf{e}_k = P_k(\mathbf{A})\mathbf{e}_0,$$

$$\|\mathbf{e}_k\|_{\mathbf{A}} = \|P_k(\mathbf{A})\mathbf{e}_0\|_{\mathbf{A}} \leq \|P(\mathbf{A})\mathbf{e}_0\|_{\mathbf{A}} \quad \forall P \in \mathbb{P}^k, P(0) = 1$$



# Polynomial residual expansions

(6/6)

From previous equations we have the characterization of CG error

$$\|e_k\|_{\mathbf{A}} = \inf_{P \in \mathbb{P}^k, P(0)=1} \|P(\mathbf{A})e_0\|_{\mathbf{A}}$$

Thus, an estimate of the form

$$\|e_k\|_{\mathbf{A}} \leq C_k \|e_0\|_{\mathbf{A}}$$

can be obtained by using estimate on the polynomial of the form

$$\left\{ P \in \mathbb{P}^k, P(0) = 1 \right\}$$



# Convergence rate calculation

## Lemma

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  an SPD matrix, and  $p \in \mathbb{P}^k$  a polynomial, then

$$\|p(\mathbf{A})\mathbf{x}\|_{\mathbf{A}} \leq \|p(\mathbf{A})\|_2 \|\mathbf{x}\|_{\mathbf{A}}$$

## Proof.

(1/2).

The matrix  $\mathbf{A}$  is SPD so that we can write

$$\mathbf{A} = \mathbf{U}^T \mathbf{\Lambda} \mathbf{U}, \quad \mathbf{\Lambda} = \text{DIAG}\{\lambda_1, \lambda_2, \dots, \lambda_n\}$$

where  $\mathbf{U}$  is an orthogonal matrix (i.e.  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ ) and  $\mathbf{\Lambda} \geq \mathbf{0}$  is diagonal. We can define the SPD matrix  $\mathbf{A}^{1/2}$  as follows

$$\mathbf{A}^{1/2} = \mathbf{U}^T \mathbf{\Lambda}^{1/2} \mathbf{U}, \quad \mathbf{\Lambda}^{1/2} = \text{DIAG}\{\lambda_1^{1/2}, \lambda_2^{1/2}, \dots, \lambda_n^{1/2}\}$$

and obviously  $\mathbf{A}^{1/2} \mathbf{A}^{1/2} = \mathbf{A}$ .



Proof.

(2/2).

Notice that

$$\|\mathbf{x}\|_{\mathbf{A}}^2 = \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{A}^{1/2} \mathbf{A}^{1/2} \mathbf{x} = \left\| \mathbf{A}^{1/2} \mathbf{x} \right\|_2^2$$

so that

$$\begin{aligned} \|p(\mathbf{A})\mathbf{x}\|_{\mathbf{A}} &= \left\| \mathbf{A}^{1/2} p(\mathbf{A}) \mathbf{x} \right\|_2 \\ &= \left\| p(\mathbf{A}) \mathbf{A}^{1/2} \mathbf{x} \right\|_2 \\ &\leq \|p(\mathbf{A})\|_2 \left\| \mathbf{A}^{1/2} \mathbf{x} \right\|_2 \\ &= \|p(\mathbf{A})\|_2 \|\mathbf{x}\|_{\mathbf{A}} \end{aligned}$$





## Lemma

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  an SPD matrix, and  $p \in \mathbb{P}^k$  a polynomial, then

$$\|p(\mathbf{A})\|_2 = \max_{\lambda \in \sigma(\mathbf{A})} |p(\lambda)|$$

## Proof.

The matrix  $p(\mathbf{A})$  is symmetric, and for a generic symmetric matrix  $\mathbf{B}$  we have

$$\|\mathbf{B}\|_2 = \max_{\lambda \in \sigma(\mathbf{B})} |\lambda|$$

observing that if  $\lambda$  is an eigenvalue of  $\mathbf{A}$  then  $p(\lambda)$  is an eigenvalue of  $p(\mathbf{A})$  the thesis easily follows.  $\square$

- Starting the error estimate

$$\|e_k\|_{\mathbf{A}} \leq \inf_{P \in \mathbb{P}^k, P(0)=1} \|P(\mathbf{A})e_0\|_{\mathbf{A}}$$

- Combining the last two lemma we easily obtain the estimate

$$\|e_k\|_{\mathbf{A}} \leq \inf_{P \in \mathbb{P}^k, P(0)=1} \left[ \max_{\lambda \in \sigma(\mathbf{A})} |P(\lambda)| \right] \|e_0\|_{\mathbf{A}}$$

- The convergence rate is estimated by bounding the constant

$$\inf_{P \in \mathbb{P}^k, P(0)=1} \left[ \max_{\lambda \in \sigma(\mathbf{A})} |P(\lambda)| \right]$$



# Finite termination of Conjugate Gradient

## Theorem (Finite termination of Conjugate Gradient)

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  an SPD matrix, the the *Conjugate Gradient* applied to the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  terminate finding the exact solution in at most  $n$ -step.

## Proof.

From the estimate

$$\|\mathbf{e}_k\|_{\mathbf{A}} \leq \inf_{P \in \mathbb{P}^k, P(0)=1} \left[ \max_{\lambda \in \sigma(\mathbf{A})} |P(\lambda)| \right] \|\mathbf{e}_0\|_{\mathbf{A}}$$

choosing 
$$P(x) = \prod_{\lambda \in \sigma(\mathbf{A})} (x - \lambda) / \prod_{\lambda \in \sigma(\mathbf{A})} (0 - \lambda)$$

we have  $\max_{\lambda \in \sigma(\mathbf{A})} |P(\lambda)| = 0$  and  $\|\mathbf{e}_n\|_{\mathbf{A}} = 0$ . □

# Convergence rate of Conjugate Gradient

- 1 The constant

$$\inf_{P \in \mathbb{P}^k, P(0)=1} \left[ \max_{\lambda \in \sigma(\mathbf{A})} |P(\lambda)| \right]$$

is not easy to evaluate,

- 2 The following bound, is useful

$$\max_{\lambda \in \sigma(\mathbf{A})} |P(\lambda)| \leq \max_{\lambda \in [\lambda_1, \lambda_n]} |P(\lambda)|$$

- 3 in particular the final estimate will be obtained by

$$\inf_{P \in \mathbb{P}^k, P(0)=1} \left[ \max_{\lambda \in \sigma(\mathbf{A})} |P(\lambda)| \right] \leq \max_{\lambda \in [\lambda_1, \lambda_n]} |\bar{P}_k(\lambda)|$$

where  $\bar{P}_k(x)$  is an opportune  $k$ -degree polynomial for which  $\bar{P}_k(0) = 1$  and it is easy to evaluate  $\max_{\lambda \in [\lambda_1, \lambda_n]} |\bar{P}_k(\lambda)|$ .

# Chebyshev Polynomials

(1/4)

- 1 The **Chebyshev Polynomials of the First Kind** are the right polynomial for this estimate. This polynomial have the following definition in the interval  $[-1, 1]$ :

$$T_k(x) = \cos(k \arccos(x))$$

- 2 Another equivalent definition valid in the interval  $(-\infty, \infty)$  is the following

$$T_k(x) = \frac{1}{2} \left[ \left( x + \sqrt{x^2 - 1} \right)^k + \left( x - \sqrt{x^2 - 1} \right)^k \right]$$

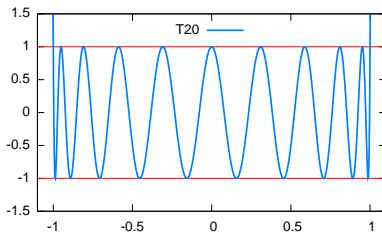
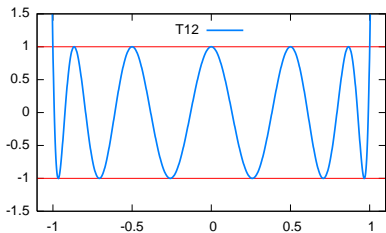
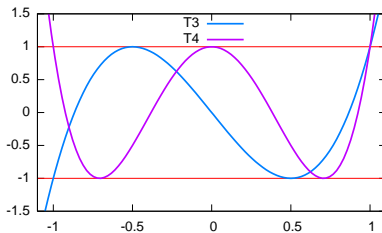
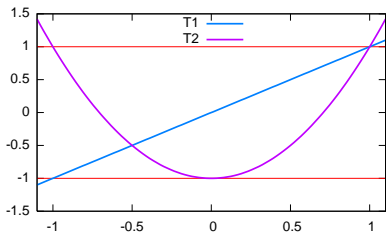
- 3 In spite of these definition,  $T_k(x)$  is effectively a polynomial.



# Chebyshev Polynomials

(2/4)

Some example of Chebyshev Polynomials.



## Chebyshev Polynomials

(3/4)

- ① It is easy to show that  $T_k(x)$  is a polynomial by the use of

$$\cos(\alpha + \beta) = \cos \alpha \cos \beta - \sin \alpha \sin \beta$$

$$\cos(\alpha + \beta) + \cos(\alpha - \beta) = 2 \cos \alpha \cos \beta$$

let  $\theta = \arccos(x)$ :

①  $T_0(x) = \cos(0\theta) = 1;$

②  $T_1(x) = \cos(1\theta) = x;$

③  $T_2(x) = \cos(2\theta) = \cos(\theta)^2 - \sin(\theta)^2 = 2 \cos(\theta)^2 - 1 = 2x^2 - 1;$

④  $T_{k+1}(x) + T_{k-1}(x) = \cos((k+1)\theta) + \cos((k-1)\theta)$   
 $= 2 \cos(k\theta) \cos(\theta) = 2x T_k(x)$

- ② In general we have the following recurrence:

①  $T_0(x) = 1;$

②  $T_1(x) = x;$

③  $T_{k+1}(x) = 2x T_k(x) - T_{k-1}(x).$



# Chebyshev Polynomials

(4/4)

- Solving the recurrence:
  - ①  $T_0(x) = 1;$
  - ②  $T_1(x) = x;$
  - ③  $T_{k+1}(x) = 2xT_k(x) - T_{k-1}(x).$
- We obtain the explicit form of the Chebyshev Polynomials

$$T_k(x) = \frac{1}{2} \left[ \left( x + \sqrt{x^2 - 1} \right)^k + \left( x - \sqrt{x^2 - 1} \right)^k \right]$$

- The translated and scaled polynomial is useful in the study of the conjugate gradient method:

$$T_k(x; a, b) = T_k\left(\frac{a + b - 2x}{b - a}\right)$$

where we have  $|T_k(x; a, b)| \leq 1$  for all  $x \in [a, b]$ .



# Convergence rate of Conjugate Gradient method

## Theorem (Convergence rate of Conjugate Gradient method)

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  an SPD matrix then the **Conjugate Gradient** method converge to the solution  $\mathbf{x}_* = \mathbf{A}^{-1}\mathbf{b}$  with at least linear  $r$ -rate in the norm  $\|\cdot\|_{\mathbf{A}}$ . Moreover we have the error estimate

$$\|\mathbf{e}_k\|_{\mathbf{A}} \lesssim 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \|\mathbf{e}_0\|_{\mathbf{A}}$$

$\kappa = M/m$  is the **condition number** where  $m = \lambda_1$  is the smallest eigenvalue of  $\mathbf{A}$  and  $M = \lambda_n$  is the biggest eigenvalue of  $\mathbf{A}$ .

The expression  $a_k \lesssim b_k$  means that for all  $\epsilon > 0$  there exists  $k_0 > 0$  such that:

$$a_k \leq (1 - \epsilon)b_k, \quad \forall k > k_0$$

## Proof.

From the estimate

$$\|e_k\|_{\mathbf{A}} \leq \max_{\lambda \in [m, M]} |P(\lambda)| \|e_0\|_{\mathbf{A}}, \quad P \in \mathbb{P}^k, P(0) = 1$$

choosing  $P(x) = T_k(x; m, M)/T_k(0; m, M)$  from the fact that  $|T_k(x; m, M)| \leq 1$  for  $x \in [m, M]$  we have

$$\|e_k\|_{\mathbf{A}} \leq T_k(0; m, M)^{-1} \|e_0\|_{\mathbf{A}} = T_k\left(\frac{M+m}{M-m}\right)^{-1} \|e_0\|_{\mathbf{A}}$$

observe that  $\frac{M+m}{M-m} = \frac{\kappa+1}{\kappa-1}$  and

$$T_k\left(\frac{\kappa+1}{\kappa-1}\right)^{-1} = 2 \left[ \left(\frac{\sqrt{\kappa}+1}{\sqrt{\kappa}-1}\right)^k + \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^k \right]^{-1}$$

finally notice that  $\left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^k \rightarrow 0$  as  $k \rightarrow \infty$ . □



# Outline

- 1 The Steepest Descent iterative scheme
- 2 Conjugate direction method
- 3 Conjugate Gradient method
- 4 Conjugate Gradient convergence rate
- 5 Preconditioning the Conjugate Gradient method**
- 6 Nonlinear Conjugate Gradient extension

# Preconditioning

## Problem (Preconditioned linear system)

Given  $\mathbf{A}, \mathbf{P} \in \mathbb{R}^{n \times n}$ , with  $\mathbf{A}$  an SPD matrix and  $\mathbf{P}$  non singular matrix and  $\mathbf{b} \in \mathbb{R}^n$ .

$$\text{Find } \mathbf{x}_* \in \mathbb{R}^n \text{ such that: } \mathbf{P}^{-T} \mathbf{A} \mathbf{x}_* = \mathbf{P}^{-T} \mathbf{b}.$$

A **good** choice for  $\mathbf{P}$  should be such that  $\mathbf{M} = \mathbf{P}^T \mathbf{P} \approx \mathbf{A}$ , where  $\approx$  denotes that  $\mathbf{M}$  is an approximation of  $\mathbf{A}$  in **some sense to precise later**.

Notice that:

- $\mathbf{P}$  non singular imply:

$$\mathbf{P}^{-T}(\mathbf{b} - \mathbf{A}\mathbf{x}) = \mathbf{0} \quad \iff \quad \mathbf{b} - \mathbf{A}\mathbf{x} = \mathbf{0};$$

- $\mathbf{A}$  SPD imply  $\tilde{\mathbf{A}} = \mathbf{P}^{-T} \mathbf{A} \mathbf{P}^{-1}$  is also SPD (obvious proof).

Now we reformulate the preconditioned system:

### Problem (Preconditioned linear system)

Given  $\mathbf{A}, \mathbf{P} \in \mathbb{R}^{n \times n}$ , with  $\mathbf{A}$  an SPD matrix and  $\mathbf{P}$  non singular matrix and  $\mathbf{b} \in \mathbb{R}^n$  the preconditioned problem is the following:

$$\text{Find } \widetilde{\mathbf{x}}_{\star} \in \mathbb{R}^n \text{ such that: } \quad \widetilde{\mathbf{A}}\widetilde{\mathbf{x}}_{\star} = \widetilde{\mathbf{b}}$$

where

$$\widetilde{\mathbf{A}} = \mathbf{P}^{-T}\mathbf{A}\mathbf{P}^{-1} \quad \widetilde{\mathbf{b}} = \mathbf{P}^{-T}\mathbf{b}$$

notice that if  $\mathbf{x}_{\star}$  is the solution of the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  then  $\widetilde{\mathbf{x}}_{\star} = \mathbf{P}\mathbf{x}_{\star}$  is the solution of the linear system  $\widetilde{\mathbf{A}}\mathbf{x} = \widetilde{\mathbf{b}}$ .



## PCG: preliminary version

initial step:

$k \leftarrow 0$ ;  $\mathbf{x}_0$  assigned;

$\tilde{\mathbf{x}}_0 \leftarrow \mathbf{P}\mathbf{x}_0$ ;  $\tilde{\mathbf{r}}_0 \leftarrow \tilde{\mathbf{b}} - \tilde{\mathbf{A}}\tilde{\mathbf{x}}_0$ ;  $\tilde{\mathbf{p}}_1 \leftarrow \tilde{\mathbf{r}}_0$ ;

**while**  $\|\tilde{\mathbf{r}}_k\| > \epsilon$  **do**

$k \leftarrow k + 1$ ;

Conjugate direction method

$$\tilde{\alpha}_k \leftarrow \frac{\tilde{\mathbf{r}}_{k-1}^T \tilde{\mathbf{r}}_{k-1}}{\tilde{\mathbf{p}}_k^T \tilde{\mathbf{A}} \tilde{\mathbf{p}}_k};$$

$$\tilde{\mathbf{x}}_k \leftarrow \tilde{\mathbf{x}}_{k-1} + \tilde{\alpha}_k \tilde{\mathbf{p}}_k;$$

$$\tilde{\mathbf{r}}_k \leftarrow \tilde{\mathbf{r}}_{k-1} - \tilde{\alpha}_k \tilde{\mathbf{A}} \tilde{\mathbf{p}}_k;$$

Residual orthogonalization

$$\tilde{\beta}_k \leftarrow \frac{\tilde{\mathbf{r}}_k^T \tilde{\mathbf{r}}_k}{\tilde{\mathbf{r}}_{k-1}^T \tilde{\mathbf{r}}_{k-1}};$$

$$\tilde{\mathbf{p}}_{k+1} \leftarrow \tilde{\mathbf{r}}_k + \tilde{\beta}_k \tilde{\mathbf{p}}_k;$$

**end while**

final step

$$\mathbf{P}^{-1} \tilde{\mathbf{x}}_k;$$



Conjugate gradient algorithm applied to  $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$  require the evaluation of thing like:

$$\tilde{\mathbf{A}}\tilde{\mathbf{p}}_k = \mathbf{P}^{-T} \mathbf{A} \mathbf{P}^{-1} \tilde{\mathbf{p}}_k.$$

this can be done **without evaluate directly the matrix  $\tilde{\mathbf{A}}$** , by the following operations:

- ① solve  $\mathbf{P} \mathbf{s}'_k = \tilde{\mathbf{p}}_k$  for  $\mathbf{s}'_k = \mathbf{P}^{-1} \tilde{\mathbf{p}}_k$ ;
- ② evaluate  $\mathbf{s}''_k = \mathbf{A} \mathbf{s}'_k$ ;
- ③ solve  $\mathbf{P}^T \mathbf{s}'''_k = \mathbf{s}''_k$  for  $\mathbf{s}'''_k = \mathbf{P}^{-T} \mathbf{s}''_k$ .

Step 1 and 3 require the solution of two auxiliary linear system. This is not a big problem if  $\mathbf{P}$  and  $\mathbf{P}^T$  are triangular matrices (see e.g. **incomplete Cholesky**).



**However...** we can reformulate the algorithm using only the matrices  $\mathbf{A}$  and  $\mathbf{P}$ !

### Definition

For all  $k \geq 1$ , we introduce the vector  $\mathbf{q}_k = \mathbf{P}^{-1}\tilde{\mathbf{p}}_k$ .

### Observation

If the vectors  $\tilde{\mathbf{p}}_1, \tilde{\mathbf{p}}_2, \dots, \tilde{\mathbf{p}}_k$  for all  $1 \leq k \leq n$  are  $\tilde{\mathbf{A}}$ -conjugate, then the corresponding vectors  $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k$  are  $\mathbf{A}$ -conjugate. In fact:

$$\mathbf{q}_j^T \mathbf{A} \mathbf{q}_i = \underbrace{\tilde{\mathbf{p}}_j^T \mathbf{P}^{-T}}_{=\mathbf{q}_j^T} \mathbf{A} \underbrace{\mathbf{P}^{-1} \tilde{\mathbf{p}}_i}_{=\mathbf{q}_i} = \tilde{\mathbf{p}}_j^T \underbrace{\tilde{\mathbf{A}}}_{=\mathbf{P}^{-T} \mathbf{A} \mathbf{P}^{-1}} \tilde{\mathbf{p}}_i = 0, \quad \text{if } i \neq j,$$

that is a consequence of  $\tilde{\mathbf{A}}$ -conjugation of vectors  $\tilde{\mathbf{p}}_i$ .



## Definition

For all  $k \geq 1$ , we introduce the vectors

$$\mathbf{x}_k = \mathbf{x}_{k-1} + \tilde{\alpha}_k \mathbf{q}_k.$$

## Observation

If we assume, **by construction**,  $\tilde{\mathbf{x}}_0 = \mathbf{P}\mathbf{x}_0$ , then we have

$$\tilde{\mathbf{x}}_k = \mathbf{P}\mathbf{x}_k, \quad \text{for all } k \text{ with } 1 \leq k \leq n.$$

In fact, if  $\tilde{\mathbf{x}}_{k-1} = \mathbf{P}\mathbf{x}_{k-1}$  (inductive hypothesis), then

$$\begin{aligned} \tilde{\mathbf{x}}_k &= \tilde{\mathbf{x}}_{k-1} + \tilde{\alpha}_k \tilde{\mathbf{p}}_k && \text{[preconditioned CG]} \\ &= \mathbf{P}\mathbf{x}_{k-1} + \tilde{\alpha}_k \mathbf{P}\mathbf{q}_k && \text{[inductive Hyp. defs of } \mathbf{q}_k \text{]} \\ &= \mathbf{P}(\mathbf{x}_{k-1} + \tilde{\alpha}_k \mathbf{q}_k) && \text{[obvious]} \\ &= \mathbf{P}\mathbf{x}_k && \text{[defs. of } \mathbf{x}_k \text{]} \end{aligned}$$



## Observation

Because  $\tilde{\mathbf{x}}_k = \mathbf{P}\mathbf{x}_k$  for all  $k \geq 0$ , we have the recurrence between the corresponding residue  $\tilde{\mathbf{r}}_k = \tilde{\mathbf{b}} - \tilde{\mathbf{A}}\tilde{\mathbf{x}}_k$  and  $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$ :

$$\tilde{\mathbf{r}}_k = \mathbf{P}^{-T}\mathbf{r}_k.$$

In fact,

$$\begin{aligned} \tilde{\mathbf{r}}_k &= \tilde{\mathbf{b}} - \tilde{\mathbf{A}}\tilde{\mathbf{x}}_k, && [\text{defs. of } \tilde{\mathbf{r}}_k] \\ &= \mathbf{P}^{-T}\mathbf{b} - \mathbf{P}^{-T}\mathbf{A}\mathbf{P}^{-1}\mathbf{P}\mathbf{x}_k, && [\text{defs. of } \tilde{\mathbf{b}}, \tilde{\mathbf{A}}, \tilde{\mathbf{x}}_k] \\ &= \mathbf{P}^{-T}(\mathbf{b} - \mathbf{A}\mathbf{x}_k), && [\text{obvious}] \\ &= \mathbf{P}^{-T}\mathbf{r}_k. && [\text{defs. of } \mathbf{r}_k] \end{aligned}$$

## Definition

For all  $k$ , with  $1 \leq k \leq n$ , the vector  $\mathbf{z}_k$  is the solution of the linear system

$$\mathbf{M}\mathbf{z}_k = \mathbf{r}_k.$$

where  $\mathbf{M} = \mathbf{P}^T \mathbf{P}$ . Formally,

$$\mathbf{z}_k = \mathbf{M}^{-1} \mathbf{r}_k = \mathbf{P}^{-1} \mathbf{P}^{-T} \mathbf{r}_k.$$

Using the vectors  $\{\mathbf{z}_k\}$ ,

- we can express  $\tilde{\alpha}_k$  and  $\tilde{\beta}_k$  in terms of  $\mathbf{A}$ , the residual  $\mathbf{r}_k$ , and conjugate direction  $\mathbf{q}_k$ ;
- we can build a recurrence relation for the  $\mathbf{A}$ -conjugate directions  $\mathbf{q}_k$ .

## Observation

$$\begin{aligned}\tilde{\alpha}_k &= \frac{\tilde{\mathbf{r}}_{k-1}^T \tilde{\mathbf{r}}_{k-1}}{\tilde{\mathbf{p}}_k^T \tilde{\mathbf{A}} \tilde{\mathbf{p}}_k} = \frac{\mathbf{r}_{k-1}^T \mathbf{P}^{-1} \mathbf{P}^{-T} \mathbf{r}_{k-1}}{\mathbf{q}_k^T \mathbf{P}^T \mathbf{P}^{-T} \mathbf{A} \mathbf{P}^{-1} \mathbf{P} \mathbf{q}_k} = \frac{\mathbf{r}_{k-1}^T \mathbf{M}^{-1} \mathbf{r}_{k-1}}{\mathbf{q}_k^T \mathbf{A} \mathbf{q}_k}, \\ &= \boxed{\frac{\mathbf{r}_{k-1}^T \mathbf{z}_{k-1}}{\mathbf{q}_k^T \mathbf{A} \mathbf{q}_k}}.\end{aligned}$$

## Observation

$$\begin{aligned}\tilde{\beta}_k &= \frac{\tilde{\mathbf{r}}_k^T \tilde{\mathbf{r}}_k}{\tilde{\mathbf{r}}_{k-1}^T \tilde{\mathbf{r}}_{k-1}} = \frac{\mathbf{r}_k^T \mathbf{P}^{-1} \mathbf{P}^{-T} \mathbf{r}_k}{\mathbf{r}_{k-1}^T \mathbf{P}^{-1} \mathbf{P}^{-T} \mathbf{r}_{k-1}} = \frac{\mathbf{r}_k^T \mathbf{M}^{-1} \mathbf{r}_k}{\mathbf{r}_{k-1}^T \mathbf{M}^{-1} \mathbf{r}_{k-1}}, \\ &= \boxed{\frac{\mathbf{r}_k^T \mathbf{z}_k}{\mathbf{r}_{k-1}^T \mathbf{z}_{k-1}}}.\end{aligned}$$



## Observation

Using the vector  $\mathbf{z}_k = \mathbf{M}^{-1}\mathbf{r}_k$ , the following recurrence is true

$$\mathbf{q}_{k+1} = \mathbf{z}_k + \tilde{\beta}_k \mathbf{q}_k$$

In fact:

$$\tilde{\mathbf{p}}_{k+1} = \tilde{\mathbf{r}}_k + \tilde{\beta}_k \tilde{\mathbf{p}}_k \quad [\text{preconditioned CG}]$$

$$\mathbf{P}^{-1} \tilde{\mathbf{p}}_{k+1} = \mathbf{P}^{-1} \tilde{\mathbf{r}}_k + \tilde{\beta}_k \mathbf{P}^{-1} \tilde{\mathbf{p}}_k \quad [\text{left mult } \mathbf{P}^{-1}]$$

$$\mathbf{P}^{-1} \tilde{\mathbf{p}}_{k+1} = \mathbf{P}^{-1} \mathbf{P}^{-T} \mathbf{r}_k + \tilde{\beta}_k \mathbf{P}^{-1} \tilde{\mathbf{p}}_k \quad [\mathbf{r}_{k+1} = \mathbf{P}^{-T} \mathbf{r}_{k+1}]$$

$$\mathbf{P}^{-1} \tilde{\mathbf{p}}_{k+1} = \mathbf{M}^{-1} \mathbf{r}_k + \tilde{\beta}_k \mathbf{P}^{-1} \tilde{\mathbf{p}}_k \quad [\mathbf{M}^{-1} = \mathbf{P}^{-1} \mathbf{P}^{-T}]$$

$$\mathbf{q}_{k+1} = \mathbf{z}_k + \tilde{\beta}_k \mathbf{q}_k \quad [\mathbf{q}_k = \mathbf{P}^{-1} \tilde{\mathbf{p}}_k]$$



# PCG: final version

initial step:

$k \leftarrow 0$ ;  $\mathbf{x}_0$  assigned;

$\mathbf{r}_0 \leftarrow \mathbf{b} - \mathbf{A}\mathbf{x}_0$ ;  $\mathbf{q}_1 \leftarrow \mathbf{r}_0$ ;

**while**  $\|\mathbf{z}_k\| > \epsilon$  **do**

$k \leftarrow k + 1$ ;

Conjugate direction method

$$\tilde{\alpha}_k \leftarrow \frac{\mathbf{r}_{k-1}^T \mathbf{z}_{k-1}}{\mathbf{q}_k^T \mathbf{A} \mathbf{q}_k};$$

$$\mathbf{x}_k \leftarrow \mathbf{x}_{k-1} + \tilde{\alpha}_k \mathbf{q}_k;$$

$$\mathbf{r}_k \leftarrow \mathbf{r}_{k-1} - \tilde{\alpha}_k \mathbf{A} \mathbf{q}_k;$$

Preconditioning

$$\mathbf{z}_k = \mathbf{M}^{-1} \mathbf{r}_k;$$

Residual orthogonalization

$$\tilde{\beta}_k \leftarrow \frac{\mathbf{r}_k^T \mathbf{z}_k}{\mathbf{r}_{k-1}^T \mathbf{z}_{k-1}};$$

$$\mathbf{q}_{k+1} \leftarrow \mathbf{z}_k + \tilde{\beta}_k \mathbf{q}_k;$$

**end while**

# Outline

- 1 The Steepest Descent iterative scheme
- 2 Conjugate direction method
- 3 Conjugate Gradient method
- 4 Conjugate Gradient convergence rate
- 5 Preconditioning the Conjugate Gradient method
- 6 Nonlinear Conjugate Gradient extension

# Nonlinear Conjugate Gradient extension

- 1 The conjugate gradient algorithm can be extended for nonlinear minimization.
- 2 Fletcher and Reeves extend CG for the minimization of a general non linear function  $f(\mathbf{x})$  as follows:
  - 1 Substitute the evaluation of  $\alpha_k$  by an line search
  - 2 Substitute the residual  $\mathbf{r}_k$  with the gradient  $\nabla f(\mathbf{x}_k)$
- 3 We also translate the index for the search direction  $\mathbf{p}_k$  to be more consistent with the gradients. The resulting algorithm is in the next slide



# Fletcher and Reeves Nonlinear Conjugate Gradient

**initial step:**

$k \leftarrow 0$ ;  $\mathbf{x}_0$  assigned;

$f_0 \leftarrow f(\mathbf{x}_0)$ ;  $\mathbf{g}_0 \leftarrow \nabla f(\mathbf{x}_0)^T$ ;

$\mathbf{p}_0 \leftarrow -\mathbf{g}_0$ ;

**while**  $\|\mathbf{g}_k\| > \epsilon$  **do**

$k \leftarrow k + 1$ ;

**Conjugate direction method**

Compute  $\alpha_k$  by line-search;

$\mathbf{x}_k \leftarrow \mathbf{x}_{k-1} + \alpha_k \mathbf{p}_{k-1}$ ;

$\mathbf{g}_k \leftarrow \nabla f(\mathbf{x}_k)^T$ ;

**Residual orthogonalization**

$$\beta_k^{FR} \leftarrow \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_{k-1}^T \mathbf{g}_{k-1}};$$

$$\mathbf{p}_k \leftarrow -\mathbf{g}_k + \beta_k^{FR} \mathbf{p}_{k-1};$$

**end while**

- ① To ensure convergence and apply Zoutendijk global convergence theorem we need to ensure that  $\mathbf{p}_k$  is a descent direction.
- ②  $\mathbf{p}_0$  is a descent direction by construction, for  $\mathbf{p}_k$  we have

$$\mathbf{g}_k^T \mathbf{p}_k = -\|\mathbf{g}_k\|^2 + \beta_k^{FR} \mathbf{g}_k^T \mathbf{p}_{k-1}$$

if the line-search is **exact** than  $\mathbf{g}_k^T \mathbf{p}_{k-1} = 0$  because  $\mathbf{p}_{k-1}$  is the direction of the line-search. So by induction  $\mathbf{p}_k$  is a descent direction.

- ③ Exact line-search is expensive, however if we use inexact line-search with **strong Wolfe** conditions
  - ① **sufficient decrease**:  $f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) \leq f(\mathbf{x}_k) + c_1 \alpha_k \nabla f(\mathbf{x}_k)^T \mathbf{p}_k$ ;
  - ② **curvature condition**:  $|\nabla f(\mathbf{x}_k + \alpha_k \mathbf{p}_k)^T \mathbf{p}_k| \leq c_2 |\nabla f(\mathbf{x}_k)^T \mathbf{p}_k|$ .

with  $0 < c_1 < c_2 < 1/2$  then we can prove that  $\mathbf{p}_k$  is a descent direction.



The previous consideration permits to say that Fletcher and Reeves nonlinear conjugate gradient method with strong Wolfe line-search is globally convergent<sup>1</sup>

To prove globally convergence we need the following lemma:

### Lemma (descent direction bound)

*Suppose we apply Fletcher and Reeves nonlinear conjugate gradient method to  $f(x)$  with strong Wolfe line-search with  $0 < c_2 < 1/2$ . The the method generates descent direction  $\mathbf{p}_k$  that satisfy the following inequality*

$$-\frac{1}{1-c_2} \leq \frac{\mathbf{g}_k^T \mathbf{p}_k}{\|\mathbf{g}_k\|^2} \leq -\frac{1-2c_2}{1-c_2}, \quad k = 0, 1, 2, \dots$$

---

<sup>1</sup>globally here means that Zoutendijk like theorem apply

Proof.

(1/3).

The proof is by induction. First notice that the function

$$t(\xi) = \frac{2\xi - 1}{1 - \xi}$$

is monotonically increasing on the interval  $[0, 1/2]$  and that  $t(0) = -1$  and  $t(1/2) = 0$ . Hence, because of  $c_2 \in (0, 1/2)$  we have:

$$-1 < \frac{2c_2 - 1}{1 - c_2} < 0. \quad (\star)$$

**base of induction  $k = 0$ :** For  $k = 0$  we have  $p_0 = -g_0$  so that  $g_0^T p_0 / \|g_0\|^2 = -1$ . From  $(\star)$  the lemma inequality is trivially satisfied.

Proof.

(2/3).

Using update direction formula's of the algorithm:

$$\beta_k^{FR} = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_{k-1}^T \mathbf{g}_{k-1}} \quad \mathbf{p}_k = -\mathbf{g}_k + \beta_k^{FR} \mathbf{p}_{k-1}$$

we can write

$$\frac{\mathbf{g}_k^T \mathbf{p}_k}{\|\mathbf{g}_k\|^2} = -1 + \beta_k^{FR} \frac{\mathbf{g}_k^T \mathbf{p}_{k-1}}{\|\mathbf{g}_k\|^2} = -1 + \frac{\mathbf{g}_k^T \mathbf{p}_{k-1}}{\|\mathbf{g}_{k-1}\|^2}$$

and by using second strong Wolfe condition:

$$-1 + c_2 \frac{\mathbf{g}_{k-1}^T \mathbf{p}_{k-1}}{\|\mathbf{g}_{k-1}\|^2} \leq \frac{\mathbf{g}_k^T \mathbf{p}_k}{\|\mathbf{g}_k\|^2} \leq -1 - c_2 \frac{\mathbf{g}_{k-1}^T \mathbf{p}_{k-1}}{\|\mathbf{g}_{k-1}\|^2}$$



Proof.

(3/3).

by induction we have

$$\frac{1}{1 - c_2} \geq -\frac{\mathbf{g}_{k-1}^T \mathbf{p}_{k-1}}{\|\mathbf{g}_{k-1}\|^2} > 0$$

so that

$$\frac{\mathbf{g}_k^T \mathbf{p}_k}{\|\mathbf{g}_k\|^2} \leq -1 - c_2 \frac{\mathbf{g}_{k-1}^T \mathbf{p}_{k-1}}{\|\mathbf{g}_{k-1}\|^2} \leq -1 + c_2 \frac{1}{1 - c_2} = \frac{2c_2 - 1}{1 - c_2}$$

and

$$\frac{\mathbf{g}_k^T \mathbf{p}_k}{\|\mathbf{g}_k\|^2} \geq -1 + c_2 \frac{\mathbf{g}_{k-1}^T \mathbf{p}_{k-1}}{\|\mathbf{g}_{k-1}\|^2} \geq -1 - c_2 \frac{1}{1 - c_2} = -\frac{1}{1 - c_2}$$



- 1 The inequality of the the previous lemma can be written as:

$$\frac{1}{1 - c_2} \frac{\|\mathbf{g}_k\|}{\|\mathbf{p}_k\|} \geq -\frac{\mathbf{g}_k^T \mathbf{p}_k}{\|\mathbf{g}_k\| \|\mathbf{p}_k\|} \geq \frac{1 - 2c_2}{1 - c_2} \frac{\|\mathbf{g}_k\|}{\|\mathbf{p}_k\|} > 0$$

- 2 Remembering the Zoutendijk theorem we have

$$\sum_{k=1}^{\infty} (\cos \theta_k)^2 \|\mathbf{g}_k\|^2 < \infty, \quad \text{where} \quad \cos \theta_k = -\frac{\mathbf{g}_k^T \mathbf{p}_k}{\|\mathbf{g}_k\| \|\mathbf{p}_k\|}$$

- 3 so that if  $\|\mathbf{g}_k\| / \|\mathbf{p}_k\|$  is bounded from below we have that  $\cos \theta_k \geq \delta$  for all  $k$  and then from Zoutendijk theorem the scheme converge.
- 4 Unfortunately this bound cant be proved so that Zoutendijk theorem cant be applied directly. However it is possible to prove a weaker results, i.e. that  $\liminf_{k \rightarrow \infty} \|\mathbf{g}_k\| = 0!$



# Convergence of Fletcher and Reeves method

## Assumption (Regularity assumption)

We assume  $f \in C^1(\mathbb{R}^n)$  with Lipschitz continuous gradient, i.e. there exists  $\gamma > 0$  such that

$$\|\nabla f(\mathbf{x})^T - \nabla f(\mathbf{y})^T\| \leq \gamma \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$$





## Theorem (Convergence of Fletcher and Reeves method)

Suppose the method of *Fletcher and Reeves* is implemented with strong Wolfe line-search with  $0 < c_1 < c_2 < 1/2$ . If  $f(\mathbf{x})$  and  $\mathbf{x}_0$  satisfy the previous regularity assumptions, then

$$\liminf_{k \rightarrow \infty} \|\mathbf{g}_k\| = 0$$

Proof.

(1/4).

From previous Lemma we have

$$\cos \theta_k \geq \frac{1}{1 - c_2} \frac{\|\mathbf{g}_k\|}{\|\mathbf{p}_k\|} \quad k = 1, 2, \dots$$

substituting in Zoutendijk condition we have  $\sum_{k=1}^{\infty} \frac{\|\mathbf{g}_k\|^4}{\|\mathbf{p}_k\|^2} < \infty$ .

**The proof is by contradiction.** in fact if theorem is not true than the series diverge. Next we want to bound  $\|\mathbf{p}_k\|$ .



Proof. (bounding  $\|\mathbf{p}_k\|$ )

(2/4).

Using second Wolfe condition and previous Lemma

$$|\mathbf{g}_k^T \mathbf{p}_{k-1}| \leq -c_2 \mathbf{g}_k^T \mathbf{p}_{k-1} \leq \frac{c_2}{1 - c_2} \|\mathbf{g}_{k-1}\|^2$$

using  $\mathbf{p}_k = -\mathbf{g}_k + \beta_k^{FR} \mathbf{p}_{k-1}$  we have

$$\begin{aligned} \|\mathbf{p}_k\|^2 &\leq \|\mathbf{g}_k\|^2 + 2\beta_k^{FR} |\mathbf{g}_k^T \mathbf{p}_{k-1}| + (\beta_k^{FR})^2 \|\mathbf{p}_{k-1}\|^2 \\ &\leq \|\mathbf{g}_k\|^2 + \frac{2c_2}{1 - c_2} \beta_k^{FR} \|\mathbf{g}_{k-1}\|^2 + (\beta_k^{FR})^2 \|\mathbf{p}_{k-1}\|^2 \end{aligned}$$

recall that  $\beta_k^{FR} = \|\mathbf{g}_k\|^2 / \|\mathbf{g}_{k-1}\|^2$  then

$$\|\mathbf{p}_k\|^2 \leq \frac{1 + c_2}{1 - c_2} \|\mathbf{g}_k\|^2 + (\beta_k^{FR})^2 \|\mathbf{p}_{k-1}\|^2$$



## Proof. (bounding $\|\mathbf{p}_k\|$ ) (3/4).

setting  $c_3 = \frac{1+c_2}{1-c_2}$  and using repeatedly the last inequality we obtain:

$$\begin{aligned}
 \|\mathbf{p}_k\|^2 &\leq c_3 \|\mathbf{g}_k\|^2 + (\beta_k^{FR})^2 (c_3 \|\mathbf{g}_{k-1}\|^2 + (\beta_{k-1}^{FR})^2 \|\mathbf{p}_{k-2}\|^2) \\
 &= c_3 \|\mathbf{g}_k\|^4 \left( \|\mathbf{g}_k\|^{-2} + \|\mathbf{g}_{k-1}\|^{-2} \right) + \frac{\|\mathbf{g}_k\|^4}{\|\mathbf{g}_{k-2}\|^4} \|\mathbf{p}_{k-2}\|^2 \\
 &\leq c_3 \|\mathbf{g}_k\|^4 \left( \|\mathbf{g}_k\|^{-2} + \|\mathbf{g}_{k-1}\|^{-2} + \|\mathbf{g}_{k-2}\|^{-2} \right) \\
 &\quad + \frac{\|\mathbf{g}_k\|^4}{\|\mathbf{g}_{k-3}\|^4} \|\mathbf{p}_{k-3}\|^2 \\
 &\leq c_3 \|\mathbf{g}_k\|^4 \sum_{j=1}^k \|\mathbf{g}_j\|^{-2}
 \end{aligned}$$

Proof.

(4/4).

Suppose now **by contradiction** there exists  $\delta > 0$  such that  $\|\mathbf{g}_k\| \geq \delta^a$  by using the regularity assumptions we have

$$\|\mathbf{p}_k\|^2 \leq c_3 \|\mathbf{g}_k\|^4 \sum_{j=1}^k \|\mathbf{g}_j\|^{-2} \leq c_3 \|\mathbf{g}_k\|^4 \delta^{-2} k$$

Substituting in Zoutendijk condition we have

$$\infty > \sum_{k=1}^{\infty} \frac{\|\mathbf{g}_k\|^4}{\|\mathbf{p}_k\|^2} \geq \frac{\delta^2}{c_4} \sum_{k=1}^{\infty} \frac{1}{k} = \infty$$

this contradict assumption. □

---

<sup>a</sup>the correct assumption is that there exists  $k_0$  such that  $\|\mathbf{g}_k\| \geq \delta$  for  $k \geq k_0$  but this complicate a little bit the following inequality without introducing new idea.

# Weakness of Fletcher and Reeves method

- Suppose that  $\mathbf{p}_k$  is a **bad** search direction, i.e.  $\cos \theta_k \approx 0$ .
- From the **descent direction bound** Lemma (see slide 91) we have

$$\frac{1}{1 - c_2} \frac{\|\mathbf{g}_k\|}{\|\mathbf{p}_k\|} \geq \cos \theta_k \geq \frac{1 - 2c_2}{1 - c_2} \frac{\|\mathbf{g}_k\|}{\|\mathbf{p}_k\|} > 0$$

- so that to have  $\cos \theta_k \approx 0$  we need  $\|\mathbf{p}_k\| \gg \|\mathbf{g}_k\|$ .
- since  $\mathbf{p}_k$  is a bad direction near orthogonal to  $\mathbf{g}_k$  it is likely that the step is small and  $\mathbf{x}_{k+1} \approx \mathbf{x}_k$ . If so we have also  $\mathbf{g}_{k+1} \approx \mathbf{g}_k$  and  $\beta_{k+1}^{FR} \approx 1$ .
- but remember that  $\mathbf{p}_{k+1} \leftarrow -\mathbf{g}_{k+1} + \beta_{k+1}^{FR} \mathbf{p}_k$ , so that  $\mathbf{p}_{k+1} \approx \mathbf{p}_k$ .
- This means that a **long sequence of unproductive iterates** will follow.

# Polack and Ribière Nonlinear Conjugate Gradient

- 1 The previous problem can be elided if we restart anew when the iterate stagnate.
- 2 Restarting is obtained by simply set  $\beta_k^{FR} = 0$ .
- 3 A more elegant solution can be obtained with a new definition of  $\beta_k$  due to Polack and Ribière is the following:

$$\beta_k^{PR} = \frac{\mathbf{g}_k^T (\mathbf{g}_k - \mathbf{g}_{k-1})}{\mathbf{g}_{k-1}^T \mathbf{g}_{k-1}}$$

- 4 This definition of  $\beta_k^{PR}$  is identical of  $\beta_k^{FR}$  in the case of quadratic function because  $\mathbf{g}_k^T \mathbf{g}_{k-1} = 0$ . The definition differs in non linear case and in particular when there is stagnation i.e.  $\mathbf{g}_k \approx \mathbf{g}_{k-1}$  we have  $\beta_k^{PR} \approx 0$ , i.e. we have an automatic restart.



# Polack and Ribière Nonlinear Conjugate Gradient

initial step:

$k \leftarrow 0$ ;  $\mathbf{x}_0$  assigned;

$f_0 \leftarrow f(\mathbf{x}_0)$ ;  $\mathbf{g}_0 \leftarrow \nabla f(\mathbf{x}_0)^T$ ;

$\mathbf{p}_0 \leftarrow -\mathbf{g}_0$ ;

**while**  $\|\mathbf{g}_k\| > \epsilon$  **do**

$k \leftarrow k + 1$ ;

**Conjugate direction method**

Compute  $\alpha_k$  by line-search;

$\mathbf{x}_k \leftarrow \mathbf{x}_{k-1} + \alpha_k \mathbf{p}_{k-1}$ ;

$\mathbf{g}_k \leftarrow \nabla f(\mathbf{x}_k)^T$ ;

**Residual orthogonalization**

$$\beta_k^{PR} \leftarrow \frac{\mathbf{g}_k^T (\mathbf{g}_k - \mathbf{g}_{k-1})}{\mathbf{g}_{k-1}^T \mathbf{g}_{k-1}};$$

$$\mathbf{p}_k \leftarrow -\mathbf{g}_k + \beta_k^{PR} \mathbf{p}_{k-1}$$

**end while**

# Weakness of Polack and Ribière method

(1/2)

- Although the modification is minimal, for the Polack and Ribière method with strong Wolfe line-search it can happen that  $\mathbf{p}_k$  is not a descent direction.
- If  $\mathbf{p}_k$  is not a descent direction we can restart i.e. set  $\beta_k^{PR} = 0$  or modify  $\beta_k^{PR}$  as follows

$$\beta_k^{PR+} = \max\{\beta_k^{PR}, 0\}$$

this new coefficient with a modified Wolfe line-search ensure that  $\mathbf{p}_k$  is a descent direction.





# Weakness of Polack and Ribière method

(2/2)

- Polack and Ribière choice on the average perform better than Fletcher and Reeves but there is **not** convergence results!
- Although there is not convergence results there is a negative results due to Powell:

## Theorem

*Consider the Polack and Ribière method with exact line-search. There exists a twice continuously differentiable function  $f : \mathbb{R}^3 \mapsto \mathbb{R}$  and a starting point  $x_0$  such that the sequence of gradients  $\{ \|g_k\| \}$  is bounded away from zero.*

- However in spite of this results Polack and Ribière is the first choice among conjugate direction methods.

# Other choices




- There are many other modification of the coefficient  $\beta_k$  that collapse to the same coefficient in the case o quadratic function. One important choice is the Hestenes and Stiefel choice

$$\beta_k^{HS} = \frac{\mathbf{g}_k^T (\mathbf{g}_k - \mathbf{g}_{k-1})}{(\mathbf{g}_k^T - \mathbf{g}_{k-1}^T) \mathbf{p}_{k-1}}$$

- For this choice there is similar convergence results of Fletcher and Reeves and similar performance.



# References

-  J. E. Dennis, Jr. and Robert B. Schnabel  
Numerical Methods for Unconstrained Optimization and  
Nonlinear Equations  
SIAM, Classics in Applied Mathematics, **16**, 1996.
-  J. Nocedal and S. J. Wright  
Numerical Optimization  
Springer Series in Operation Research, 1999.
-  J. Stoer and R. Bulirsch  
Introduction to numerical analysis  
Springer-Verlag, Texts in Applied Mathematics, **12**, 2002.