

## §5. The Conjugate Gradient method

Iterative solver for symm pos def systems

$$Ax = b \quad A \in \mathbb{R}^{n \times n} \text{ symm pos def}$$

$$\Leftrightarrow \min_{x \in \mathbb{R}^n} \frac{1}{2} x^T A x - x^T b = f(x)$$
$$\nabla f(x) = Ax - b$$
$$\nabla^2 f(x) = A$$

### Gradient method:

$$x_{k+1} = x_k - \alpha_k \nabla f_k$$

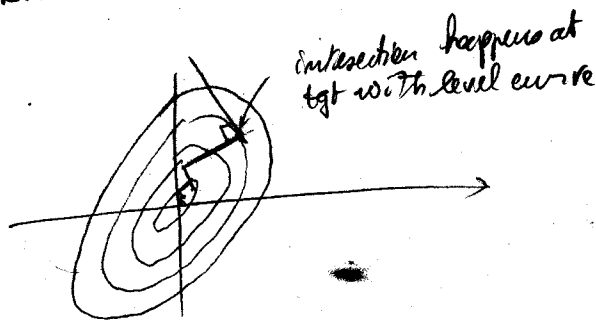
where  $\alpha_k$  solves  $\min_{\alpha} \varphi(\alpha) = f(x_k - \alpha \nabla f_k)$  i.e.:

$$\begin{aligned} \varphi(\alpha) &= -\nabla f_k^T (x_k - \alpha \nabla f_k) + \frac{1}{2} (x_k - \alpha \nabla f_k)^T A (x_k - \alpha \nabla f_k) \\ &= -\nabla f_k^T A x_k + \alpha \nabla f_k^T A x_k + \frac{1}{2} \|\nabla f_k\|^2 - \alpha \nabla f_k^T A \nabla f_k \\ &= \frac{1}{2} \|\nabla f_k\|^2 + \alpha \nabla f_k^T A \nabla f_k \end{aligned}$$

$$\varphi'(\alpha) > 0 \quad \rightarrow \quad \alpha_k = \frac{\|\nabla f_k\|^2}{\nabla f_k^T A \nabla f_k}$$

Possible to show  $\nabla f_k^T \nabla f_{k+1} = 0$ . Method can take long to converge

"Zig-zagging" 2D:



### The linear conjugate gradient method

$x_{k+1} = x_k + \alpha_k p_k$ , where  $\alpha_k p_k$  is determined s.t.  $x_{k+1}$  solves

$$\min_{x \in x_0 + \text{span}\{p_0, \dots, p_{k-1}, r_k\}} \frac{1}{2} x^T A x - b^T x \quad (1)$$

where  $r_k = -\nabla f_k = b - Ax_k$

(1)  $\Leftrightarrow$  (2)

$$\min \frac{1}{2} \hat{x}^T A \hat{x} - \hat{x}^T r_0$$

$$\hat{x} \in \text{span}\{p_0, \dots, p_{k-1}, r_k\}$$

$$\text{where } \hat{x} = x - x_0$$

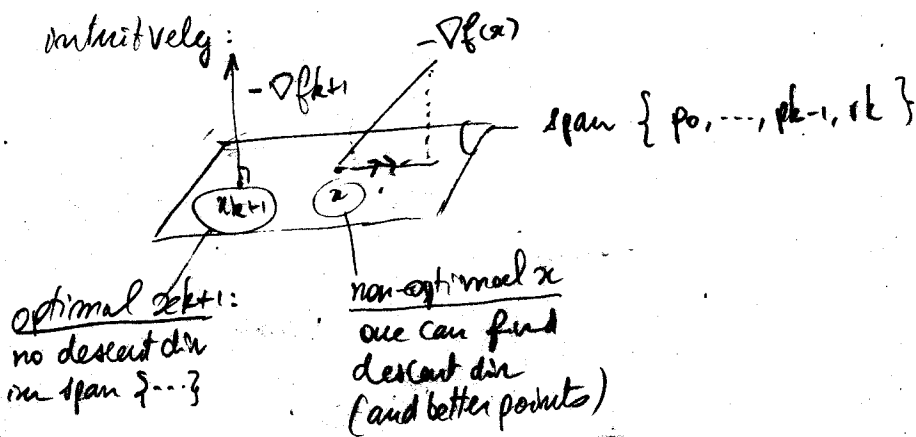
(22)

It can be shown that  $\hat{x}_{k+1}$  solves (2) iff

$$(A \hat{x}_{k+1} - r_0)^T v = 0 \quad \forall v \in \text{span}\{p_0, \dots, p_{k-1}, r_k\}$$

$$\Leftrightarrow \underbrace{(A \hat{x}_{k+1} - b)^T v = 0}_{-\nabla f_{k+1}}$$

(~ Galerkin)



In previous step:  $\hat{x}_k$  solves:

$$\min \frac{1}{2} \hat{x}^T A \hat{x} - \hat{x}^T r_0$$

$$\hat{x} \in \text{span}\{p_0, \dots, p_{k-1}\}$$

$$\Rightarrow P_j^T (b - A x_k) = 0$$
  
$$j = 0 \dots k-1$$

Since  $x_{k+1} = x_k + \alpha_k p_k$ :

$$0 = (A x_{k+1} - b)^T p_j = (A x_k - b)^T p_j + \alpha_k p_k^T A p_j$$

$\Rightarrow p_k$  is  $A$ -orthogonal (constant inner prod  $(u, v)_A = u^T A v$ )  
to previous  $k$  directions

$\rightarrow$  get  $p_k$  using Gram-Schmidt orthogonalization,  
which will greatly simplify:

$$\begin{cases} p_0 = r_0 \\ \vdots \\ p_k = r_k - \sum_{j=0}^{k-1} \frac{r_k^T A p_j}{p_j^T A p_j} p_j \end{cases} \quad (GS)$$

Now that we know in which direction to go, we can use gradient method to find

$$\boxed{\alpha_k = \underset{\alpha}{\operatorname{argmin}} \phi(\alpha) = f(x_k + \alpha p_k)} \\ = \frac{r_k^T p_k}{p_k^T A p_k}$$

problem w/ two preliminary versions } # of operations grows linearly } with storage

A closer look to the subspace used gives:

$$\begin{aligned} \operatorname{span}\{p_0, \dots, p_{k-1}, r_k\} &= \operatorname{span}\{p_0, \dots, p_{k-1}, p_k\} \\ &= \operatorname{span}\{r_0, \dots, r_k\} \\ &= \mathcal{K}_{k+1}(A, r_0) \\ &= \operatorname{span}\{r_0, A r_0, A^2 r_0, \dots, A^k r_0\} \\ &= \text{Krylov subspace} \end{aligned}$$

why?  $r_0 = b - A x_0 \in \mathcal{K}_1(A, r_0)$   
 $r_1 = b - A x_1 = b - A(x_0 + \alpha(b - A x_0)) \in \mathcal{K}_2(A, r_0)$   
 etc...

Why is this useful?

optimality conditions:  $r_{k+1}^T v = 0 \quad \forall v \in \mathcal{K}_{k+1}(A, r_0)$

$r_k^T v = 0 \quad \forall v \in \mathcal{K}_k(A, r_0)$

take  $p_i \in \mathcal{K}_{k-1}(A, r_0) \Rightarrow A p_i \in \mathcal{K}_k(A, r_0)$

$\Rightarrow r_k^T A p_i = 0 \quad i = 0 \dots k-2$

thus (GS) reduces to:

$$p_k = r_k - \frac{r_k^T A p_{k-1}}{p_{k-1}^T A p_{k-1}} p_{k-1} \\ = \beta_k$$

Storage:  $x_k, r_k, p_k, A p_k$   
 requires only knowledge of matrix-vector product

Can show:

$$\alpha_k = \frac{\|r_k\|^2}{p_k^T A p_k}$$

$$\beta_k = \frac{\|r_{k+1}\|^2}{\|r_k\|^2}$$

→ classical CG formulation (saves some MV prod)  
see Algo 5.2 p 112.

What happens if A is indefinite? ⇒

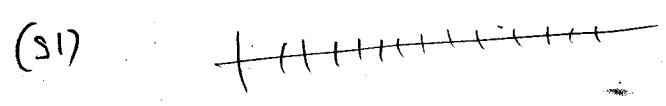
→ ∃ k for which  $p_k^T A p_k < 0$

→ "negative curvature" direction → stop iterations

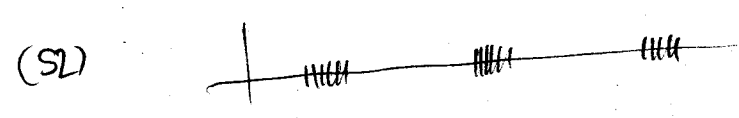
$\min x^T A x - x^T b$  has no sol

Convergence results: convergence properties determined by eigenvalues of A.

the fewer "clusters" the better.



many iterations



few (~3) iterations.

preconditioning

$$A x = b \Leftrightarrow M A x = M b$$

M = preconditioner  
easy to compute  
and somehow transform  
→ spectrum of A from  
(S1) to (S2).