On preconditioned Krylov subspace methods in numerical PDEs.

Zdeněk Strakoš

Nečas Center for Mathematical Modeling Charles University in Prague and Czech Academy of Sciences http://www.karlin.mff.cuni.cz/~strakos

Ascona, CSF, September 2013



R. C. Kirby, SIREV (2010):

"We examine condition numbers, preconditioners and iterative methods for FEM discretization of coercive PDEs in the context of the solvability result, the Lax-Milgram lemma.

Moreover, useful insight is gained as to the relationship between Hilbert space and matrix condition numbers, and translating Hilbert space fixed point iterations into matrix computations provides new ways of motivating and explaining some classic iteration schemes."

Krylov subspace methods, fixed point iteration and pre-conditioning?



K. A. Mardal and R. Winther, NLAA (2011):

"The main focus will be on an abstract approach to the construction of preconditioners for symmetric linear systems in a Hilbert space setting [...] The discussion of preconditioned Krylov space methods for the continuous systems will be a starting point for a corresponding discrete theory.

By using this characterization it can be established that the conjugate gradient method converges [...] with a rate which can be bounded by the condition number [...] However, if the operator has a few eigenvalues far away from the rest of the spectrum, then the estimate is not sharp. In fact, a few 'bad eigenvalues' will have almost no effect on the asymptotic convergence of the method."

Axelsson (1976), quote Jennings (1977)



FIG. 4. A Chebyshev polynomial modified by a simple third order auxiliary polynomial having zeros at λ_1 , λ_2 and λ_n .

p. 72: ... it may be inferred that rounding errors ... affects the convergence rate when large outlying eigenvalues are present.



O. Axelsson and J. Karátson, Numer. Alg. (2009):

"To preserve sparsity, the arising system is normally solved using an iterative solution method, commonly a preconditioned conjugate gradient [PCG] method [...] the rate of convergence depends in general on a generalized condition number of the preconditioned operator [...]

- if the two operators (original and preconditioner) are equivalent then the corresponding PCG method provides mesh independent linear convergence [...]
- if the two operators (original and preconditioner) are compact-equivalent then the corresponding PCG method provides mesh independent superlinear convergence."



- Computational cost of finding sufficiently accurate approximation to the exact solution of the algebraic problem heavily depends on
 - * the underlying real world problem,
 - * the mathematical model including the concept of solution,
 - * on its discretization.
- Construction and analysis of computational algorithms should respect that. We should always think in terms of approximations. Exact solutions can *in principle* be uncomputable (eigenvalues).
- Evaluation of accuracy and of computational cost in numerical PDEs must take into account algebraic errors, including rounding errors.



Outline

- 1. Functional formulation and preconditioning in PDEs
- 2. Numerical discretization: consistency, stability, and convergence
- 3. CG, conditioning, *a-priori* algebraic error estimates
- 4. Spectral theory and the moment problem formulation
- 5. How to measure (*a-posteriori*) errors?
- 6. Preconditioning as transformation of the basis
- 7. Reaching an arbitrary accuracy?
- 8. Conclusions



Let V be a real infinite dimensional Hilbert space with the inner product

 $(\cdot, \cdot)_V : V \times V \to R$, the associated norm $\| \cdot \|_V$,

 $V^{\#}$ be the dual space of bounded (continuous) linear functionals on V with the duality pairing

$$\langle \cdot, \cdot \rangle : V^{\#} \times V \to R.$$

For each $f \in V^{\#}$ there exists a unique $\tau f \in V$ such that

 $\langle f, v \rangle = (\tau f, v)_V \text{ for all } v \in V.$

In this way the inner product $(\cdot, \cdot)_V$ determines the Riesz map

 $\tau: V^{\#} \to V.$



Let $a(\cdot, \cdot): V \times V \to R$ be a bounded and V-elliptic bilinear form. For a fixed $u \in V$ we can see $\mathcal{A}u \equiv a(u, \cdot) \in V^{\#}$, i.e.,

 $\langle \mathcal{A} u, v \rangle = a(u,v) \quad \text{for all } v \in V \,.$

This defines the bounded and $(\alpha -)$ coercive operator

$$\mathcal{A}: V \to V^{\#}, \quad \inf_{u \in V, \, \|u\|_{V}=1} \langle \mathcal{A}u, u \rangle = \alpha > 0, \, \|\mathcal{A}\| = C.$$

Using the Lax-Milgram theorem, the PDE problem is well-posed: For any $b \in V^{\#}$ there exist a unique solution $x \in V$ of

 $a(x,v) = \langle b,v \rangle$ for all $v \in V$.

and x depends continuously on the data b ,

$$||x||_V \leq \frac{1}{\alpha} ||b||_{V^{\#}}.$$



 $\langle \mathcal{A}x - b, v \rangle = 0$ for all $v \in V$

gives the (functional) equation in (the data space) $V^{\#}$,

$$\mathcal{A}x = b, \quad , \quad \mathcal{A}: V \to V^{\#}, \quad x \in V, \quad b \in V^{\#}.$$

Using the Riesz map,

$$(\tau \mathcal{A}x - \tau b, v)_V = 0$$
 for all $v \in V$.

Clearly, application of the Riesz map τ can be interpreted as transformation of the original problem Ax = b in the data space $V^{\#}$ into the equation in the solution space V,

$$\tau \mathcal{A} x = \tau b, \qquad \tau \mathcal{A} : V \to V, \quad x \in V, \quad \tau b \in V,$$

which is commonly (and inaccurately) called preconditioning.



Arnold, Falk, and Winther (1997, 1997); Steinbach and Wendland (1998); Mc Lean and Tran (1997); Christiansen and Nédélec (2000, 2000); Powell and Silvester (2003); Elman, Silvester, and Wathen (2005); Axelsson and Karátson (2009); Mardal and Winther (2011); Kirby (2011); Zulehner (2011); Preconditioning Conference 2013, Oxford; ...

Inner product \longrightarrow Riesz map \longrightarrow Preconditioning \longrightarrow Spectral bounds

However, there is a point to consider. What is the appropriate inner product ? A standard way is to focus on the mesh (model) parameters independence of the condition number-based convergence bounds.

Operator preconditioning \longrightarrow PDEs.

Algebraic preconditioning \longrightarrow Matrices.

Preconditioning \longrightarrow what does it tell us about **discretization?**



Consider an *N*-dimensional discrete solution subspace $V_h \subset V$ with the duality pairing and the inner product as above. Then the restriction to V_h gives an approximation $x_h \in V_h$ to $x \in V$,

$$a(x_h, v) = \langle b, v \rangle$$
 for all $v \in V_h$.

As above, the bilinear form $a(\cdot, \cdot): V_h \times V_h \to R$ defines the operator $\mathcal{A}_h: V_h \to V_h^{\#}$ such that

$$\langle \mathcal{A}_h x_h - b, v \rangle = 0$$
 for all $v \in V_h$.

With restricting *b* to V_h , i.e. $\langle b_h, v \rangle \equiv \langle b, v \rangle$ for all $v \in V_h$, we get the equation in the discrete data space $V_h^{\#}$,

$$\mathcal{A}_h x_h = b_h, \qquad x_h \in V_h, \quad b_h \in V_h^{\#}, \quad \mathcal{A}_h : V_h \to V_h^{\#}.$$



Let X_h be a representation of the solution x in V_h .

- Consistency error norm $\|A_h X_h b_h\|_{V_h^{\#}}$. The discretization scheme is consistent if the consistency error norm tends to 0 with h.
- Stability = continuity of the discrete mapping $\mathcal{A}_h^{-1}: V_h^{\#} \to V_h$. The discretization scheme is stable if the stability constant $\|\mathcal{A}_h^{-1}\|_{V_h^{\#}, V_h}$ is bounded uniformly in h. Here

$$\|\mathcal{A}_h^{-1}\|_{V_h^{\#},V_h} \leq \frac{1}{\alpha}.$$

• The discretization scheme is convergent if the error norm $||X_h - x_h||_{V_h}$ tends to 0 with h.



Using

$$\|X_h - x_h\|_{V_h} \leq \|\mathcal{A}_h^{-1}\|_{V_h^{\#}, V_h} \|\mathcal{A}_h X_h - b_h\|_{V_h^{\#}},$$

a discretization scheme which is consistent and stable is convergent.

Instructive (and more general) exposition in Arnold (2012); see also Arnold, Falk and Winther (2010); ... From the computational point of view, one issue is, however, missing. Here it is assumed that x_h satisfies

$$\mathcal{A}_h x_h = b_h \, .$$

In numerical algebra this reminds of a bound for the forward error using a residual backward error and a conditioning of the problem.

Incorporating algebraic errors?



Using the Riesz map, $\ \tau \mathcal{A}: V \to V$. We can therefore form for $\ g \in V$ the Krylov sequence

$$g, \tau \mathcal{A}g, (\tau \mathcal{A})^2 g, \ldots$$
 in V

and define Krylov subspace methods in the Hilbert space operator setting (here we will do CG). Our goal is to construct a method for solving the functional equation

$$\mathcal{A}x = b, \quad x \in V, \quad b \in V^{\#}$$

such that with $r_0 = b - Ax_0 \in V^{\#}$ the approximations x_n to the solution x, n = 1, 2, ... belong to the Krylov manifolds in V

$$x_n \in x_0 + K_n(\tau \mathcal{A}, \tau r_0) \equiv$$

$$x_0 + \operatorname{span}\{\tau r_0, \tau \mathcal{A}(\tau r_0), (\tau \mathcal{A})^2(\tau r_0), \dots, (\tau \mathcal{A})^{n-1}(\tau r_0)\}.$$



3 Self-adjoint \mathcal{A} wrt the duality pairing

CG in Hilbert spaces : $r_0 = b - Ax_0 \in V^{\#}$, $p_0 = \tau r_0 \in V$ For $n = 1, 2, ..., n_{max}$

$$\begin{aligned} \alpha_{n-1} &= \frac{\langle r_{n-1}, \tau r_{n-1} \rangle}{\langle \mathcal{A}p_{n-1}, p_{n-1} \rangle} = \frac{(\tau r_{n-1}, \tau r_{n-1})_V}{(\tau \mathcal{A}p_{n-1}, p_{n-1})_V} \\ x_n &= x_{n-1} + \alpha_{n-1}p_{n-1} , \qquad \text{stop when the stopping criterion is satisfied} \\ r_n &= r_{n-1} - \alpha_{n-1}\mathcal{A}p_{n-1} \\ \beta_n &= \frac{\langle r_n, \tau r_n \rangle}{\langle r_{n-1}, \tau r_{n-1} \rangle} = \frac{(\tau r_n, \tau r_n)_V}{(\tau r_{n-1}, \tau r_{n-1})_V} \\ p_n &= \tau r_n + \beta_n p_{n-1} \end{aligned}$$

End

Hayes (1954); ...; Glowinski (2003); Axelsson and Karatson (2009); Mardal and Winther (2011); Günnel, Herzog and Sachs (2013)



Let $\Phi_h = (\phi_1^{(h)}, \dots, \phi_N^{(h)})$ be the basis of V_h , $\Phi_h^{\#} = (\phi_1^{(h)\#}, \dots, \phi_N^{(h)\#})$ the canonical basis of its dual $V_h^{\#}$. Using the coordinates in Φ_h and in $\Phi_h^{\#}$,

$$\begin{split} \langle f, v \rangle &\to \mathbf{v}^* \mathbf{f} ,\\ (u, v)_V &\to \mathbf{v}^* \mathbf{M} \mathbf{u}, \quad (\mathbf{M}_{ij}) = \left((\phi_j, \phi_i)_V \right)_{i,j=1,\dots,N} ,\\ \tau &\to \mathbf{M}^{-1} , \quad \text{the inverse of the Gram matrix of } \Phi_h \; \text{wrt} \; (.,.)_V \\ \mathcal{A}_h &\to \mathbf{A}, \quad (\mathbf{A}_{ij}) = \left(a(\phi_j, \phi_i) \right)_{i,j=1,\dots,N} = \left(\langle \mathcal{A} \phi_j, \phi_i \rangle \right)_{i,j=1,\dots,N} ,\\ b \to \mathbf{b} , \end{split}$$

we get with $x_n = \Phi_h \mathbf{x}_n, \ p_n = \Phi_h \mathbf{p}_n, \ r_n = \Phi_h^{\#} \mathbf{r}_n$



$$\mathbf{r}_0 = \mathbf{b} - \mathbf{A} \mathbf{x}_0, \quad \text{solve} \quad \mathbf{M} \mathbf{z}_0 = \mathbf{r}_0, \ \mathbf{p}_0 = \mathbf{z}_0$$

For $n = 1, \ldots, n_{\max}$

$$\begin{split} \alpha_{n-1} &= \frac{\mathbf{z}_{n-1}^* \mathbf{r}_{n-1}}{\mathbf{p}_{n-1}^* \mathbf{A} \mathbf{p}_{n-1}} \\ \mathbf{x}_n &= \mathbf{x}_{n-1} + \alpha_{n-1} \mathbf{p}_{n-1}, \text{ stop when the stopping criterion is satisfied} \\ \mathbf{r}_n &= \mathbf{r}_{n-1} - \alpha_{n-1} \mathbf{A} \mathbf{p}_{n-1} \\ \mathbf{z}_n &= \mathbf{M}^{-1} \mathbf{r}_n, \text{ solve for } \mathbf{z}_n \\ \beta_n &= \frac{\mathbf{z}_n^* \mathbf{r}_n}{\mathbf{z}_{n-1}^* \mathbf{r}_{n-1}} \\ \mathbf{p}_n &= \mathbf{z}_n + \beta_n \mathbf{p}_{n-1} \end{split}$$

End



3 Philosophy of a-priori robust bounds

Theorem

$$\kappa(\mathbf{M}^{-1}\mathbf{A}) \leq \frac{C}{\alpha} = \frac{\|\mathcal{A}\|}{\inf_{u \in V, \|u\|_{V}=1} \langle \mathcal{A}u, u \rangle}$$

"Knowledge of robust estimates not only contributes to the question of well-posedness, but also to discretization error estimates and the construction of efficient solvers for the discretized problem. In the dicretized case, having robust estimates [...] translates to having a [...] preconditioner for the linear operator [...] with robust estimates on the condition number. This would immediately imply that Krylov subspace methods like the minimum residual method [...] converge with convergence rates independent on [...] h."

Zullehner, SIAM J. Matrix Anal. Appl. (2011)



3 Liesen, S (2012); Gergelits, S (2013)



Short recurrences always mean in practical computations loss of (bi-)orthogonality due to rounding errors! Principal consequences are not resolved by *the common assumption* that this phenomenon does not take place.



Hiptmair, CMA (2006):

"Operator preconditioning is a very general recipe [...]. It is simple to apply, but may not be particularly efficient, because in case of the [*condition number*] bound of Theorem [...] is too large, the operator preconditioning offers no hint how to improve the preconditioner. Hence, operator preconditioner may often achieve [...] the much-vaunted mesh independence of the preconditioner, but it may not perform satisfactorily on a given mesh."



Faber, Manteuffel and Parter, Adv. in Appl. Math. (1990):

"For a fixed *h*, using a preconditioning strategy based on an equivalent operator may not be superior to classical methods [...] Equivalence alone is not sufficient for a good preconditioning strategy. One must also choose an equivalent operator for which the bound is small.

There is no flaw in the analysis, only a flaw in the conclusions drawn from the analysis [...] asymptotic estimates ignore the constant multiplier. Methods with similar asymptotic work estimates may behave quite differently in practice."

Point.

A-priori approach is too rough. A-posteriori approach to algebraic errors is needed.



From Fermat, Descartes, Euler (principal axis theorem), Lagrange, ..., Cauchy, Jacobi ... Fredholm ... through Stieltjes to Hilbert, Schmidt, Lebesgue, Hellinger and Toeplitz, Wintner, Stone, von Neumann ...

L. A. Stein, Highlights in the history of spectral theory, Amer. Math. Monthly (1973)

In connection to Krylov subspace methods (CG), please recall

 $K_n(\tau \mathcal{A}, \tau r_0) \equiv \operatorname{span}\{\tau r_0, \tau \mathcal{A}(\tau r_0), (\tau \mathcal{A})^2(\tau r_0), \dots, (\tau \mathcal{A})^{n-1}(\tau r_0)\},\$

and the (Chebyshev-Markov-) Stieltjes moment problem (see Stieltjes (1894)), the basic (almost unknown) reference is

Y. V. Vorobyev, Method of Moments in Applied Mathematics, (1958,1965) (see also the approach in Liesen, S, Krylov Subspace Methods, Principles and Analysis. OUP, (2013).)



Using the orthogonal projection E_n onto K_n with respect to the inner product $(\cdot, \cdot)_V$, consider the orthogonally restricted operator

$$\tau \mathcal{A}_n : K_n \to K_n, \quad \tau \mathcal{A}_n \equiv E_n (\tau \mathcal{A}) E_n,$$

by formulating the following equalities

$$\tau \mathcal{A}_n (\tau r_0) \equiv \tau \mathcal{A} (\tau r_0),$$

$$(\tau \mathcal{A}_n)^2 \tau r_0 = \tau \mathcal{A}_n (\tau \mathcal{A} (\tau r_0)) \equiv (\tau \mathcal{A})^2 \tau r_0,$$

$$\vdots$$

$$(\tau \mathcal{A}_n)^{n-1} \tau r_0 = \tau \mathcal{A}_n ((\tau \mathcal{A})^{n-2} \tau r_0) \equiv (\tau \mathcal{A})^{n-1} \tau r_0,$$

$$(\tau \mathcal{A}_n)^n \tau r_0 = \tau \mathcal{A}_n ((\tau \mathcal{A})^{n-1} \tau r_0) \equiv E_n (\tau \mathcal{A})^n \tau r_0.$$



The *n*-dimensional approximation τA_n of τA matches the first 2n moments

$$((\tau \mathcal{A}_n)^{\ell} \tau r_0, \tau r_0)_V = ((\tau \mathcal{A})^{\ell} \tau r_0, \tau r_0)_V, \quad \ell = 0, 1, \dots, 2n - 1.$$

Denote symbolically $Q_n = (q_1, \ldots, q_n)$ a matrix composed of the columns q_1, \ldots, q_n forming an orthonormal basis of K_n determined by the Lanczos process

$$\tau \mathcal{A} Q_n = Q_n \mathbf{T}_n + \delta_{n+1} q_{n+1} \mathbf{e}_n^T$$

with $q_1 = \tau r_0 / \|\tau r_0\|_V$. We get $(\tau A_n)^{\ell} = Q_n \mathbf{T}_n^{\ell} Q_n^*$, $\ell = 0, 1, ...$ and the matching moments condition

$$\mathbf{e}_{1}^{*} \mathbf{T}_{n}^{\ell} \mathbf{e}_{1} = q_{1}^{*} (\tau \mathcal{A})^{\ell} q_{1}, \quad l = 0, 1, \dots, 2n - 1,$$



is the Jacobi matrix of the orthogonalization coefficients and the CG method (in Hilbert spaces) is formulated as

$$\mathbf{T}_n \mathbf{y}_n = \|\tau r_0\|_V \,\mathbf{e}_1, \qquad x_n = x_0 + Q_n \mathbf{y}_n, \quad x_n \in V.$$



Since τA is bounded and self-adjoint, its spectral decomposition is written using the Riemann-Stieltjes integral as

$$\tau \mathcal{A} = \int_{\lambda_L}^{\lambda_U} \lambda \, d\mathcal{E}_\lambda \,,$$

The spectral function \mathcal{E}_{λ} of $\tau \mathcal{A}$ represents a family of orthogonal projections which is

- non-decreasing, i.e., if $\mu > \nu$, then the subspace onto which \mathcal{E}_{μ} projects contains the subspace into which \mathcal{E}_{ν} projects;
- $\mathcal{E}_{\lambda_L} = 0$, $\mathcal{E}_{\lambda_U} = I$;
- \mathcal{E}_{λ} is right continuous, i.e. $\lim_{\lambda' \to \lambda_+} \mathcal{E}_{\lambda'} = \mathcal{E}_{\lambda}$.

The values of λ where \mathcal{E}_{λ} increases by jumps represent the eigenvalues of $\tau \mathcal{A}$, with the eigenvectors satisfying

$$\tau \mathcal{A} z = \lambda z, \quad z \in V.$$



4 Spectral moment problem

For the (finite) Jacobi matrix T_n we can analogously write

$$\mathbf{T}_n = \sum_{j=1}^n \theta_\ell^{(n)} \mathbf{s}_\ell^{(n)}, \quad \lambda_L < \theta_1^{(n)} < \theta_2^{(n)} < \dots < \theta_n^{(n)} < \lambda_U,$$

and the operator moment problem turns into

$$\int_{\lambda_L}^{\lambda_U} \lambda^{\ell} \, d\omega(\lambda) = \sum_{j=1}^n \{\theta_j^{(n)}\}^{\ell} \, \omega_j^{(n)}, \qquad \ell = 0, 1, \dots, 2n-1,$$

where $d\omega(\lambda) = q_1^* d\mathcal{E}_{\lambda} q_1$ represents the Riemann-Stieltjes distribution function associated with $\tau \mathcal{A}$ and q_1 . The distribution function $\omega^{(n)}(\lambda)$ approximates $\omega(\lambda)$ in the sense of the *n*-th Gauss-Christoffel quadrature; Gauss (1814), Jacobi (1826), Christoffel (1858).



4 CG \equiv Gauss-Christoffel quadrature



Condition number CG bounds should always be checked against this!



Up to now, the spectral moment problem linked directly the functional infinite dimensional CG formulation for Ax = b with the finite dimensional matrix representation

$$\mathbf{T}_n \mathbf{y}_n = \| \tau r_0 \|_V \mathbf{e}_1, \qquad x_n = x_0 + Q_n \mathbf{y}_n, \quad x_n \in V.$$

corresponding to the n-th iteration step.

In practice the infinite dimensional problem is first discretized, giving $A_h x_h = b_h$, and CG is then applied to the associated discretized matrix representation. Stability of the discretization scheme and of the moment problem representation concerns the relationship

$$\begin{array}{ccccc} \mathcal{A} & \longrightarrow & \mathcal{A}_{\{\boldsymbol{h},n\}} & \longrightarrow & \mathbf{A}_{\{\boldsymbol{h},n\}} & \longrightarrow & \mathbf{T}_{\{\boldsymbol{h},n\}} \,, \\ \mathcal{A} & \longrightarrow & \mathcal{A}_{\boldsymbol{h}} & \longrightarrow & \mathbf{A}_{\{\boldsymbol{h}\}} & \longrightarrow & & \mathbf{T}_{\{\boldsymbol{h},n\}} \,. \end{array}$$



For inverse problems polluted by noise the distribution function $\omega(\lambda)$ has a special shape.

The level of the noise can be determined simply by monitoring the first component of the eigenvector of \mathbf{T}_n associated with its smallest eigenvalue for n = 0, 1, 2, ... This can be done at almost no cost. Almost free lunches do exist.

Hnětynková, Plešinger, and S, The regularizing effect of the Golub-Kahan iterative bidiagonalization and revealing the noise level, BIT (2009); Michenková, MS Thesis, (2013).



4 Noise revealing in inverse problems



Distribution function (left) and the noise revealing indicator (right).



4 Dealing with blurred noisy elephants



The best approach to noisy ill-behaved elephants is not to apply the Golub-Kahan bidiagonalization, but to try to escape!



Theorem

1° The spectrum of A is $\{\lambda_1, \ldots, \lambda_N\}$ and GMRES(A, b) yields residuals with the prescribed nonincreasing sequence

 $||r_0|| \ge ||r_1|| \ge \cdots \ge ||r_{N-1}|| > ||r_N|| = 0.$

2° Matrix *A* is of the form $A = WRCR^{-1}W^*$ and b = Wh where *C* is the spectral companion matrix, *W* is unitary and *R* a nonsingular upper triangular matrix such that Rs = h.

Complete parametrization. Set of measure zero? Greenbaum, Ptak, Arioli and S (1994 - 98), Eirmann and Ernst (2001), Meurant (2012), Meurant and Tebbens (2012),



The bounds Const $F_n(sp(A), N)$ do not intersect the rectangle (1,0) - (1,N) - (0,N) - (0,0).







Becker, Johnson, and Rannacher (1995)

"Usually, ad hoc stopping criteria are used, e.g. requiring an initial (algebraic) residual to be reduced by a certain ad hoc factor, but these criteria have no clear connection to the actual error in the corresponding approximate solution, which is the quantity of interest. This leaves the user of iterative solutions methods in a serious dilemma: [...] one has either to continue the iterations until the discrete solution error is practically "zero", which increases the computational cost with possibly no gain in the overall precision, or take the risk of stopping the iterations prematurely. [...]

A solution to this problem can only be obtained by combining aspects of the underlying partial differential equations and the corresponding finite element discretization with aspects of the iterative discrete solution algorithm. A "pure" numerical linear algebra point of view, for instance based on the condition number of the stiffness matrix, does not appear to be able to lead to a balance of discretization and solution errors."



Babuška and Strouboulis (2001)

"In engineering practice it is not sufficient to estimate only the energy norm of the error because a small value of the global energy norm of the error does not necessarily imply that the error in the outputs of interest is also small (e.g. a 5% relative error in the global energy norm does not imply 5% relative error in the maximum stress in a region of interest). [...] An essential requirement is that the quantity of interest has to be well defined; for example, it is meaningless to ask for an estimate of the maximum error in the derivative, flux, or stress for a problem set in a polygonal non-convex domain, because the exact value does not exists (the derivative, flux, or stress in the neighborhood of a corner point is usually unbounded)."



Giles and Süli (2002)

"In many scientific and engineering applications [...] the objective is merely a rough, qualitative assessment of the details of the analytical solution over the computational domain, the quantitative concern being directed towards a few output functionals, derived quantities of particular engineering or scientific relevance."

In addition to discretization errors, algebraic errors can also affect the accuracy of the computed approximate solution. What is known on the spatial distribution of the algebraic errors over the computational domain?

Does the algebraic backward error approach resolve the matter? Unfortunately no.



Discrete (piecewise polynomial) FEM approximation $x_h = \Phi_h \mathbf{x}_n$.

- If \mathbf{x}_n is known exactly, then x_h is approximated over the given domain as the (exact) linear combination of the local basis functions.
- However, apart from trivial cases, \mathbf{x}_n that supply the global information is not known exactly. Then





5 Local discretisation





Theorem (x_h denotes the discrete Galerkin solution)

$$\begin{aligned} \|\nabla(x - x_h^{(n)})\|^2 &= \|\nabla(x - x_h)\|^2 + \|\nabla(x_h - x_h^{(n)})\|^2 \\ &= \|\nabla(x - x_h)\|^2 + \|\mathbf{x} - \mathbf{x}_n\|_{\mathbf{A}}^2 \end{aligned}$$

holds up to a small inaccuracy proportional to machine precision.

What is the distribution of the algebraic error in the functional space ?

5 L-shape domain, Papež, Liesen, S (2013)



Exact solution x (left) and the discretisation error $x - x_h$ (right) in the Poisson model problem, linear FEM.

5 L-shape domain, Papež, Liesen, S (2013)



Algebraic error $x_h - x_h^{(n)}$ (left) and the total error $x - x_h^{(n)}$ (right). Here $\|\nabla(x - x_h)\| > 0.1 \|\mathbf{x} - \mathbf{x}_n\|_{\mathbf{A}}$.



They should be based on a-posteriori error estimators which are fully computable and provide information on local distribution of the error (including the algebraic error) within the domain. Ideally, a-posteriori error estimators should satisfy the following additional properties:

- reliability (guaranteed upper bound);
- local efficiency;
- asymptotic exactness.

Verfürth (1996); Repin (1997; Ainsworth and Oden (2000); Babuška and Strouboulis (2001); Bangerth and Rannacher (2003); ... ; Deuflhard (1994); Bernd, Manteuffel, and McCormick (1996); ... ; Wohlmuth, Hoppe (1999); ... ; Arioli, Noulard, and Russo (2001); Arioli, Loghin, and Wathen (2005); Silvester and Simoncini (2011); ... ; Jiranek, S, and Vohralik (2011); Vohralik and Ern (2013); Arioli, Liesen, Miedlar, and S (2013); ...



Various approaches, influenced by relationship to moments and numerical quadrature (G. H. Golub). Lasting impact of the original paper on CG by Hestenes and Stiefel (1952).

- Dahlquist, Eisenstat, and Golub (1972); Dahlquist, Golub, and Nash (1978); ...; Fischer (1996); Golub and Meurant (1994, 1997, 2010); Golub and S (1994); ...; Brezinski (1999); ...; Calvetti, Morigi, Reichel, and Sgallari (2000, 2001); ...
- S and Tichý (2002); S and Tichý (2005); Meurant and S (2006); S and Tichý (2011); Meurant and Tichý (2013); ...

Point. Rounding error analysis is not an option, it is an imperative!



5 A-posteriori error norms estimates in CG

$$\int \lambda^{-1} d\omega(\lambda) = \sum_{i=1}^{n} \omega_i^{(n)} \left(\theta_i^{(n)}\right)^{-1} + R_n(f)$$

$$\frac{\|\mathbf{x} - \mathbf{x}_0\|_{\mathbf{A}}^2}{\|\mathbf{r}_0\|^2} = n \text{-th Gauss-Ch. quadrature} + \frac{\|\mathbf{x} - \mathbf{x}_n\|_{\mathbf{A}}^2}{\|\mathbf{r}_0\|^2}$$

$$\mathbf{r}_0^* \mathbf{A}^{-1} \mathbf{r}_0 = \sum_{j=0}^{n-1} \gamma_j \|\mathbf{r}_j\|^2 + \mathbf{r}_n^* \mathbf{A}^{-1} \mathbf{r}_n.$$

Formulas equivalent assuming exact arithmetic can (and do!) behave very differently in practical computations.

Hesteness and Stiefel (1952); Golub and S (1994); S and Tichý (2002); Meurant and S (2006); Golub and Meurant (2010); S and Tichý (2011); ...

Liesen, S, Krylov Subspace Methods. Principles and Analysis (2012)



Algebraic preconditioning can be viewed as the finite dimensional CG with setting $\mathbf{M} = \mathbf{I}$ (this corresponds in Galerkin discretization of the finite dimensional CG to taking discretization basis Φ orthonormal wrt $(.,.)_V$) applied to

$$\mathbf{B}\mathbf{w} = \mathbf{c}$$

with

$$\mathbf{B} = \mathbf{L}_h^{-1} \mathbf{A} \mathbf{L}_h^{-*}, \quad \mathbf{c} = \mathbf{L}_h^{-1} \mathbf{b}, \quad \mathbf{x} = \mathbf{L}_h^{-*} \mathbf{w}, \quad \mathbf{M}_h = \mathbf{L}_h \mathbf{L}_h^*.$$

Observation:

The associated Hilbert space formulation of CG in V_h corresponds to the transformation of the bases

$$\Phi_t = \Phi_h \mathbf{L}_h^{-*}, \quad \Phi_t^{\#} = \Phi_h^{\#} \mathbf{L}_h^*.$$



6 Preconditioning transforms the basis!

$$\mathbf{B} \equiv (\mathbf{B}_{ij}) = \left(\langle \mathcal{A}\phi_j^{(t)}, \phi_i^{(t)} \rangle \right)_{i,j=1,\dots,N} = (a(\phi_j^{(t)}, \phi_i^{(t)}))_{i,j=1,\dots,N},$$

where

$$\phi_{\ell}^{(t)} = \Phi_h \left(\mathbf{L}_h^{-*} \mathbf{e}_{\ell} \right), \quad \ell = 1, \dots, N$$

and the right hand side

$$\mathbf{c} = \Phi_h^\# \mathbf{L}_h^* \, \mathbf{b} \, .$$

Please recall, e.g., the hierarchical bases preconditioning Yserentant (1985, 1986); Axelsson; Vassilevski; ...; Gockenbach (2006)

Remark. Equivalently, with the orthonormal discretization basis Φ_t wrt $(.,.)_V$ we get $\mathbf{M} = \mathbf{I}$ and $\mathcal{A}_h \to \mathbf{B}$.

With the choice $(.,.)_V = (.,.)_a$ we get $\mathbf{B} = \mathbf{I}$.



Sparsity of the resulted matrices is always presented as the main advantage of FEM discretizations.

Sparsity means locality of information. In order to solve the problem, we need a global transfer of information. Therefore preconditioning! It is needed on the computational level in order to take care for the trouble caused by the (*computationally*) inconvenient approximation of the mathematical model when the *appropriate globally supported* basis functions are missing (cf. hierarchical bases preconditioning, DD with coarse space components, multilevel methods, ...).

Preconditioning can be interpreted as an intentional loss of sparsity (loss of locality of the supports of the basis functions).

Sparsity is important for efficiency, but perhaps in a different meaning; see, e.g., Schaeffer, Caflisch, Hauck and Osher (2013), .



7 Reaching an arbitrary accuracy?



It seems and it has been proved that an arbitrary prescribed accuracy can be reached using AFEM in a finite number of steps. Here linear FEM; see Morin, Nocheto, and Siebert (2002); Stevenson (2007). Something does not fit \longrightarrow maximal attainable accuracy in matrix computations.



- Assumptions (see, e.g., nonnormality and limitations of spectral bounds).
- Interpretations.
- Common views should always be given a critical second thought.

Analysis of finite dimensional algebraic problems can not be done on the PDE problem level using functional analysis in infinite dimensional Banach or Hilbert spaces. On the other hand, analysis of the finite dimensional algebraic problem must "do justice" to the original (non-algebraic) problem as much as possible.



8 Conclusions

Patrick J. Roache's book Validation and Verification in Computational Science, 1998, p. 387:

"With the often noted tremendous increases in computer speed and memory, and with the less often acknowledged but equally powerful increases in algorithmic accuracy and efficiency, a natural question suggest itself. What are we doing with the new computer power? with the new GUI and other set-up advances? with the new algorithms? What *should* we do? ... Get the right answer."

This requires to consider modelling, discretisation, analysis, and computation tightly coupled parts of a single solution process.



Modest steps in this direction:

- Operator and algebraic preconditioning is related to the discrerization basis.
- Krylov subspace methods viewed as the matching moments model reduction (infinite or finite dimensional setting).
- A-posteriori evaluation of the total error which is based on quantities of interest and includes the algebraic part. Algebraic <u>a-priori</u> reasoning is useful, but it addresses different questions.
- Adaptivity and stopping criteria for iterative solvers. Algebraic backward error theory not sufficient. We need backward error on the functional equation level.
- Numerical stability analysis of adaptive numerical schemes.



Recent references

- J. Liesen and Z.S., Krylov Subspace Methods, Principles and Analysis. Oxford University Press (2012)
- T. Gergelits and Z.S., Composite convergence bounds based on Chebyshev polynomials and finite precision conjugate gradient computations, Numerical Algorithms (2013) (DOI 10.1007/s11075-013-9713-z)
- J. Papez, J. Liesen and Z.S., On distribution of the discretization and algebraic error in numerical solution of partial differential equations, Preprint MORE/2012/03, (2013)
- M. Arioli, J. Liesen, A. Miedlar, and Z.S., Interplay between discretization and algebraic computation in adaptive numerical solution of elliptic PDE problems, GAMM Mitteilungen 36, 102-129 (2013)
- J. Málek and Z.S., From PDEs through functional analysis to iterative methods, or there and back again. In preparation.



Thank you very much for kind patience!

