

Numerical Methods I

MATH-GA 2010.001/CSCI-GA 2420.001

Benjamin Peherstorfer
Courant Institute, NYU

Based on slides by G. Stadler and A. Donev

Today

Last time

- ▶ Computing eigenvalues
- ▶ Singular value decomposition (SVD)

Today

- ▶ Computing the SVD
- ▶ Iterative methods for systems of linear equations

Announcements

- ▶ Homework 4 posted, is due Mon, Nov 7 before class
- ▶ There will be midterm grades, which are solely based on HW 1–2 (so not very informative!)

Recap: Singular value decomposition

Let $\underline{A} \in \mathbb{C}^{m \times n}$. A singular value decomposition of A is a factorization

$$A = \underline{U} \Sigma \underline{V}^H,$$

where

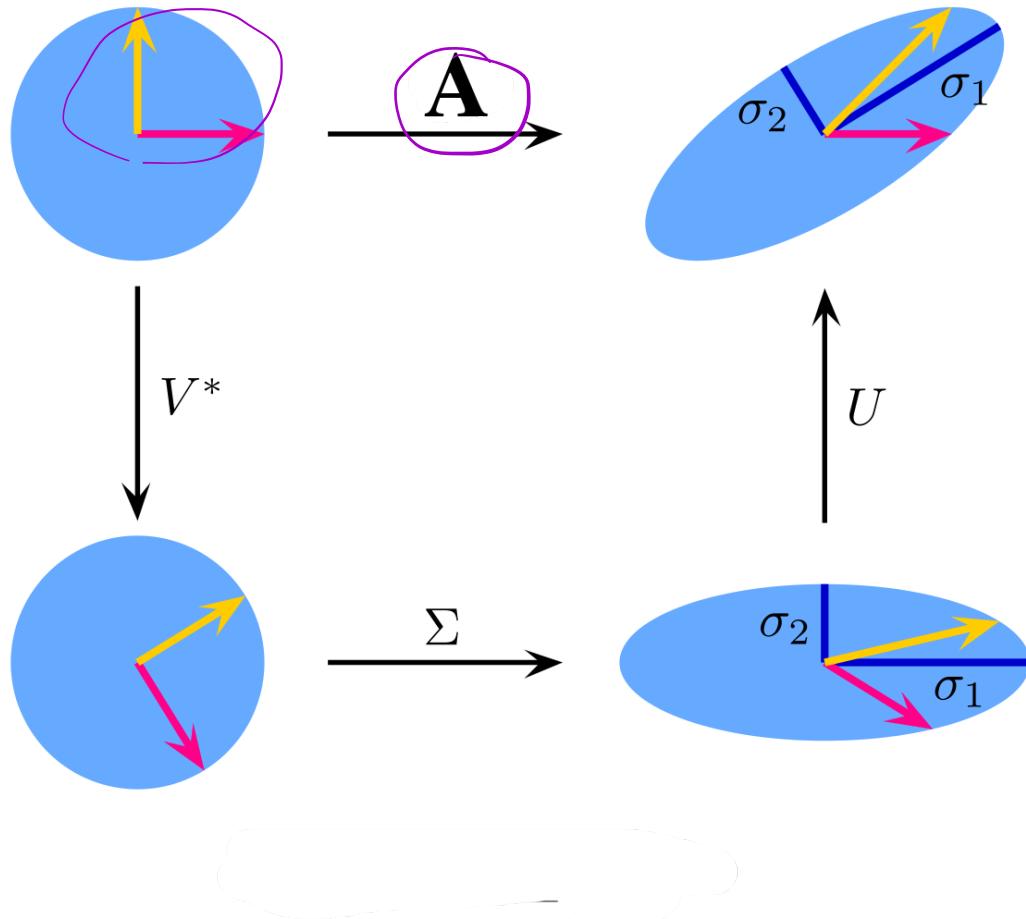
$$U \in \mathbb{C}^{m \times m} \text{ is unitary} \tag{8}$$

$$V \in \mathbb{C}^{n \times n} \text{ is unitary} \tag{9}$$

$$\Sigma \in \mathbb{R}^{m \times n} \text{ is diagonal.} \tag{10}$$

Additionally, the diagonal entries σ_j of Σ are non-negative and in non-decreasing order so that $\underline{\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0}$ where $p \in \min(m, n)$.

- ▶ The diagonal matrix Σ is real and has the same shape as A even when A is not square
- ▶ The matrices U and V are always square



The image of the unit sphere under a map A is a hyperellipse (in \mathbb{R}^m). Thus, with $A = U\Sigma V^H$, have

- ▶ The unitary map V^H preserves the sphere (rotating a sphere is a sphere)
- ▶ The diagonal matrix Σ stretches the sphere into a hyperellipse aligned with the canonical basis
- ▶ The unitary map U rotates or reflects the hyperellipse without changing shape

[Edited from figure by Georg-Johann, Wikipedia]

Recap: SVD vs. eigenvalue decomposition

SVD expresses a matrix in proper bases for column and row span to represent it as a diagonal matrix

$$A = U\Sigma V^H$$

We have seen something similar with eigenvectors: A non-defective square matrix A can be expressed as a diagonal matrix of eigenvalues Λ if the range and domain are presented in a basis of eigenvectors

$$A = X\Lambda X^{-1}$$

There are fundamental differences between the SVD and eigenvalue decomposition

- ▶ SVD uses two different bases (left and right singular vectors); eigenvalue decomposition uses just one (eigenvectors)
- ▶ SVD uses orthonormal bases, whereas eigenvalue basis generally is not orthogonal
- ▶ Not all matrices (even square ones) have an eigendecomposition, but all matrices (even rectangular ones) have a singular value decomposition

Typically, eigenvalues tell us something about the behavior of iterative processes that involve the matrix A such as A^k and e^{tA}

Singular values tend to tell us something about A itself

Recap: The SVD and matrix properties

In the following:

- ▶ The matrix A is of dimension $\underline{m \times n}$
- ▶ $\underline{p} = \min(m, n)$
- ▶ $\underline{r} \leq p$ is the number of non-zero singular values of A

We now list how the SVD is related to fundamental properties of the matrix A

- ▶ The rank of A is r , the number of non-zero singular values.
- ▶ The range (column span) of A is $\text{span}(u_1, \dots, u_r)$, the kernel is $\text{span}(v_{r+1}, \dots, v_n)$
- ▶ $\|A\|_2 = \sigma_1$ and $\|A\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_r^2}$
- ▶ For square A , $|\det(A)| = \prod_{i=1}^m \sigma_i$
- ▶ The non-zero singular values of A are the square roots of the non-zero eigenvalues of $A^H A$ and AA^H

Recap: Versions of SVD

Full SVD

$$\underbrace{\mathbf{A}}_{m \times n} = \underbrace{\mathbf{U}}_{m \times m} \underbrace{\Sigma}_{m \times n} \underbrace{\mathbf{V}^H}_{n \times n}$$

Reduced SVD

$$\underbrace{\mathbf{A}}_{m \times n} = \underbrace{\mathbf{U}}_{m \times n} \underbrace{\Sigma}_{n \times n} \underbrace{\mathbf{V}^H}_{n \times n}$$

Rank-revealing SVD of a rank r matrix with dimension $m \times n$

$$\underbrace{\mathbf{A}}_{m \times n} = \underbrace{\mathbf{U}}_{m \times r} \underbrace{\Sigma}_{r \times r} \underbrace{\mathbf{V}^H}_{r \times n}$$

Recap: SVD for low-rank approximation

Consider the SVD $A = \underline{U\Sigma V^H}$, then

$$A = \sum_{j=1}^r \sigma_j u_j v_j^H$$

Let us truncate the sum after $1 \leq \underline{q} \leq \underline{r}$ terms and define

$$A_q = \sum_{j=1}^{\underline{q}} \sigma_j u_j v_j^H.$$

Then, A_q is a best rank q approximation of A in the $\|\cdot\|_2$ norm

$$\|A - A_q\|_2 = \inf_{\substack{B \in \mathbb{C}^{m \times n} \\ \text{rank}(B) \leq q}} \|A - B\|_2 = \sigma_{q+1}$$

The error is the first left out singular value σ_{q+1} !

Recap: SVD and pseudo inverse

The (Moore-Penrose) pseudo inverse of an $m \times n$ matrix that is regular and square is $A^+ = A^{-1}$.

Otherwise, the pseudo inverse is given by

$$\underline{A}^+ = \underline{V} \underline{\Sigma}^+ \underline{U}^H,$$

where

$$\Sigma^+ = \text{diag}(\sigma_1^{-1}, \sigma_2^{-1}, \dots, \sigma_r^{-1}, 0, \dots, 0)$$

Thus, we can solve least-squares problems $\min_x \|b - Ax\|_2$ by taking the SVD of A , computing A^+ , and setting $x = A^+ b$ for the minimal norm solution w.r.t. $\|\cdot\|_2$

The approach we discussed via the QR decomposition is cheaper but the approach via the SVD is sometimes preferred because it allows to easily regularize the problem by truncating small singular values.

How to compute the SVD?

The SVD of an $m \times n$ ($m \geq n$) matrix A is related to the eigenvalue decomposition of $A^H A$

$$\underline{A^H A} = \underline{V \Sigma^H \Sigma V^H},$$

Since we know how to numerically compute the eigendecomposition, we could compute the SVD of A as follows

1. Form $\underline{A^H A}$
2. Compute the eigendecomposition $\underline{A^H A} = \underline{V \Lambda V^H}$ (Notice that $Z = A^H A$ is normal because $Z^H Z = Z Z^H$)
3. Let Σ be the $m \times n$ non-negative diagonal square root of Λ
4. Solve the system $U \Sigma = A V$ for unitary U

Is this a good idea?

[†]<https://nhigham.com/2022/10/11/seven-sins-of-numerical-linear-algebra/>

How **not** to compute the SVD?

The SVD of an $m \times n$ ($m \geq n$) matrix A is related to the eigenvalue decomposition of $A^H A$

$$A^H A = V \Sigma^H \Sigma V^H,$$

Since we know how to numerically compute the eigendecomposition, we could compute the SVD of A as follows

1. Form $A^H A$
2. Compute the eigendecomposition $A^H A = V \Lambda V^H$ (Notice that $Z = A^H A$ is normal because $Z^H Z = Z Z^H$)
3. Let Σ be the $m \times n$ non-negative diagonal square root of Λ
4. Solve the system $U \Sigma = A V$ for unitary U

Is this a good idea? \rightsquigarrow as we have seen before, it is typically dangerous[†] from a stability perspective to compute something of the matrix A via the matrix $A^H A$ (think of least-squares regression) \rightsquigarrow board

[†]<https://nhigham.com/2022/10/11/seven-sins-of-numerical-linear-algebra/>

• Perturbation of EV of $A^H A$

$$|\lambda_k(A^H A + \delta B) - \lambda_k(A^H A)| \leq \|\delta B\|_2 \quad (\text{for } k)$$

• singular values

$$|\sigma_k(A + \delta A) - \sigma_k(A)| \leq \|\delta A\|_2 \quad (\text{for } k)$$

Backward stable algo. for SV

$$\tilde{\sigma}_k = \sigma_k(A + \delta A), \quad \frac{\|\delta A\|}{\|A\|} \in O(\epsilon)$$

$$|\tilde{\sigma}_k - \sigma_k| \leq \|\delta A\| \quad \text{with} \quad \|\delta A\| \in O(\epsilon \|\underline{A}\|)$$

Now compute $\lambda_k(A^H A)$

$$\tilde{\lambda}_k = \lambda_k(A^H A + \delta B), \quad \frac{\|\delta B\|}{\|A^H A\|} \in O(\epsilon)$$

$$|\tilde{\lambda}_k - \lambda_k| \leq \|\delta B\| \in O(\epsilon \|A^H A\|) \\ = O(\epsilon \|A\|^2)$$

• Compute $\tilde{\sigma}_a = \sqrt{\tilde{\lambda}_a}$

$$|\tilde{\sigma}_a - \sigma_a| = |\sqrt{\tilde{\lambda}_a} - \sqrt{\lambda_a}| \leq \frac{|\tilde{\lambda}_a - \lambda_a|}{\underbrace{\sqrt{\tilde{\lambda}_a} + \sqrt{\lambda_a}}_{\geq 0} = \tilde{\sigma}_a}$$

$$|\tilde{\sigma}_a - \sigma_a| \text{ behaves as } O\left(\epsilon \frac{\|A\|^2}{\tilde{\sigma}_a}\right)$$

How we actually compute the SVD

For the sake of the argument, assume that A is a square $m \times m$ (following is applicable to rectangular matrices too) \rightsquigarrow board

$$H = \begin{bmatrix} 0 & A^H \\ A & 0 \end{bmatrix}$$

$$A = U \Sigma V^H$$

$$AV = U\Sigma$$

$$A^H U = V \Sigma^H U^H U = V \underbrace{\Sigma^H}_{\Sigma} = V \Sigma$$

$$\begin{aligned}
 & \underbrace{\begin{bmatrix} 0 & A^H \\ A & 0 \end{bmatrix}}_H \underbrace{\begin{bmatrix} V & V \\ U & -U \end{bmatrix}}_X = \begin{bmatrix} A^H U & -A^H U \\ A V & A V \end{bmatrix} \\
 &= \begin{bmatrix} V\Sigma & -V\Sigma \\ U\Sigma & U\Sigma \end{bmatrix} = \begin{bmatrix} V & V \\ U & -U \end{bmatrix} \begin{bmatrix} \Sigma & 0 \\ 0 & -\Sigma \end{bmatrix}
 \end{aligned}$$

1

How we actually compute the SVD

For the sake of the argument, assume that A is a square $m \times m$ (following is applicable to rectangular matrices too) \rightsquigarrow board

Consider the $2m \times 2m$ Hermitian matrix

$$H = \begin{bmatrix} 0 & A^H \\ A & 0 \end{bmatrix}$$

Since $A = U\Sigma V^H$, we have

$$AV = U\Sigma \text{ and } A^H U = V\Sigma^H = V\Sigma \quad (\text{recall that singular values are real})$$

which we write in matrix form as

$$\begin{bmatrix} 0 & A^H \\ A & 0 \end{bmatrix} \begin{bmatrix} V & V \\ U & -U \end{bmatrix} = \begin{bmatrix} V & V \\ U & -U \end{bmatrix} \begin{bmatrix} \Sigma & 0 \\ 0 & -\Sigma \end{bmatrix}$$

which is the eigendecomposition of $H \rightsquigarrow$ avoids using $A^H A$ and AA^H and is stable

Computing the SVD is typically a two-phase procedure:

First, reduce the matrix A to bidiagonal form, then diagonalize the bidiagonal matrix

$$\begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} \rightarrow \begin{bmatrix} * & * & & \\ * & * & * & \\ & * & * & * \\ & & * & * \\ & & & * \end{bmatrix} \rightarrow \begin{bmatrix} * & & & \\ & * & & \\ & & * & \\ & & & * \end{bmatrix}$$

- ▶ Phase 1 involves a finite number of operations that scale as $\mathcal{O}(mn^2)$
- ▶ Phase 2 (recall eigenvalue problems) is iterative but converges very quickly; in practice achieves convergence in $\mathcal{O}(n)$ to machine precision
- ▶ Thus, state-of-the-art computation of the SVD has costs that scale as $\mathcal{O}(mn^2)$

This is the celebrated Golub-Kahan approach from the 1960s

How can we bidiagonalize a matrix? Recall that we already know how to triangualize a matrix via (unitary) Householder reflection

$$\begin{array}{c} \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \xrightarrow{\hspace{1cm}} \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \xrightarrow{\hspace{1cm}} \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \xrightarrow{\hspace{1cm}} \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \xrightarrow{\hspace{1cm}} \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \\ A \\ \textcolor{purple}{\cancel{A}} \end{array}$$
$$\begin{array}{c} Q_1 A \\ \textcolor{purple}{\cancel{Q_1 A}} \end{array}$$
$$\begin{array}{c} Q_2 Q_1 A \\ \textcolor{purple}{\cancel{Q_2 Q_1 A}} \end{array}$$
$$\begin{array}{c} Q_3 Q_2 Q_1 A \\ \textcolor{purple}{\cancel{Q_3 Q_2 Q_1 A}} \end{array}$$
$$\begin{array}{c} Q_4 Q_3 Q_2 Q_1 A \\ \textcolor{purple}{\cancel{Q_4 Q_3 Q_2 Q_1 A}} \end{array}$$

How can we bidiagonalize a matrix? Recall that we already know how to triangualize a matrix via (unitary) Householder reflection

$$\begin{array}{c}
 \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \rightarrow \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \rightarrow \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \rightarrow \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \rightarrow \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \\
 A \qquad\qquad\qquad Q_1 A \qquad\qquad\qquad Q_2 Q_1 A \qquad\qquad\qquad Q_3 Q_2 Q_1 A \qquad\qquad\qquad Q_4 Q_3 Q_2 Q_1 A
 \end{array}$$

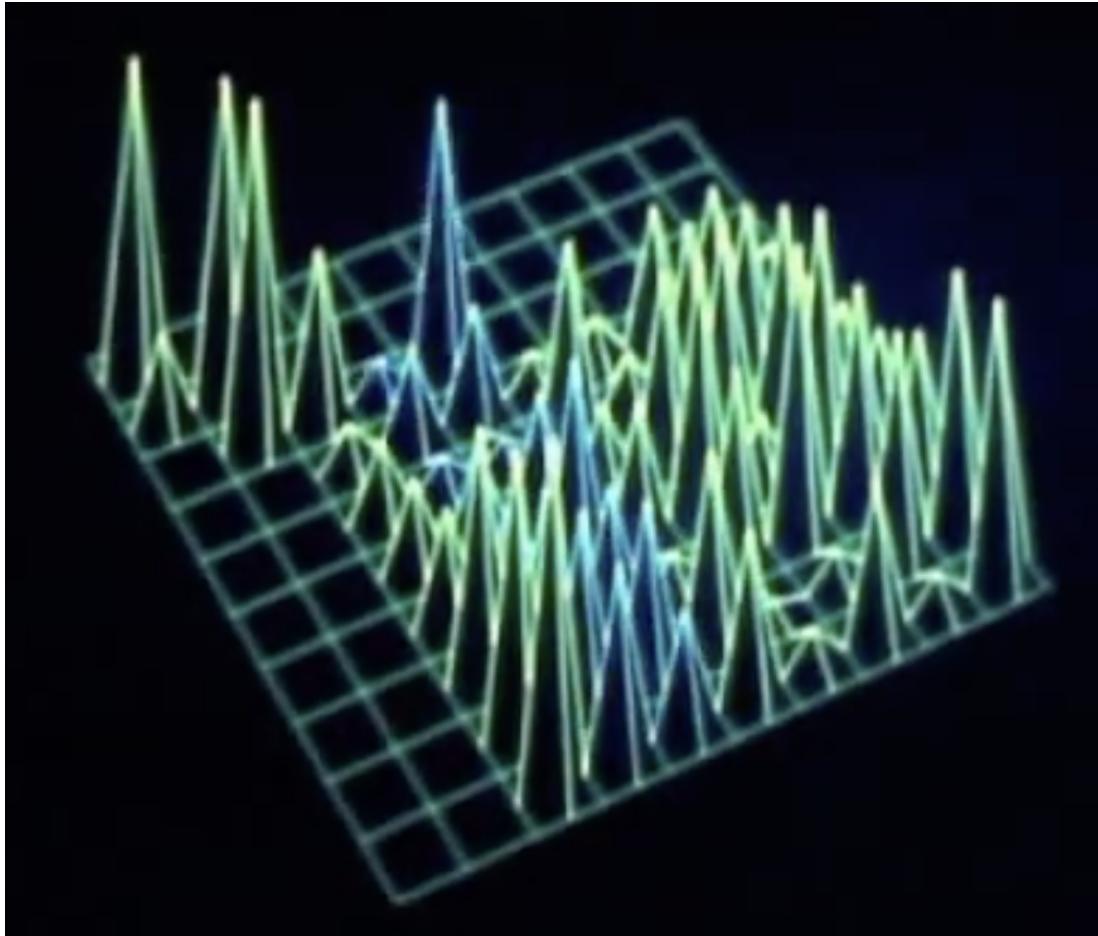
Now, we apply interleaved Householder reflection from left and right

$$\begin{array}{c}
 \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \rightarrow \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \rightarrow \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \rightarrow \left[\begin{array}{cccc} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{array} \right] \rightarrow \dots \\
 A \qquad\qquad\qquad U_1^H A \qquad\qquad\qquad U_1^H A V_1 \qquad\qquad\qquad U_2^H U_1^H A V_1 \qquad\qquad\qquad U_2^H U_1^H A V_1 V_2
 \end{array}$$

At the end, n reflectors are applied from the left and $n - 2$ from the right

A variant of the QR algorithm is then applied to the bidiagonal matrix (or other eigendecomposition algorithms) \rightsquigarrow details in Golub et al., Matrix Computations.

Visualizing SVD computation



Video: <http://youtu.be/R9UoFyqJca8>

SVD saves the universe

The first Star Trek movie came out in 1979. The producers had asked Los Alamos for computer graphics to run on the displays on the bridge of the Enterprise:

[Figure: Paramount Pictures]

Link: <https://blogs.mathworks.com/cleve/2012/12/10/1976-matrix-singular-value-decomposition-film/>

Matlab

```
1: >> X = randn(10, 4);
2: >> [U, S, V] = svd(X); % full SVD
3: >> [Ur, Sr, Vr] = svd(X, 0); % reduced (economic) SVD
4: >> size(S)
5: >> size(Sr)
6: ans =
7:      10      4
8: ans =
9:      4      4
```

In most cases, we want the reduced SVD and then we should explicitly compute it:

```
1: >> X = randn(10000, 50);
2: >> tic; [U, S, V] = svd(X); toc
3: Elapsed time is 6.746243 seconds.
4: >> tic; [U, S, V] = svd(X, 0); toc
5: Elapsed time is 0.064572 seconds.
```

For very large and sparse matrices, or if the matrix is unavailable and we only have a procedure that returns the matrix-vector product, we can iteratively compute the first few singular vectors via svds

```
1: >> X = gallery('poisson', 100);
2: >> whos X
3: Name          Size            Bytes  Class
   Attributes
4:
5:   X            10000x10000        1033608  double
   sparse
6:
7: >> tic; [U, S, V] = svds(X, 1); toc
8:
9: Elapsed time is 0.307015 seconds.
```

SVD is great because its computation is (very) stable \rightsquigarrow building block for computing many basic linear algebra quantities

type 'edit rank.m' in Matlab

```
1: function r = rank(A, tol)
2: %RANK Matrix rank.
3: % RANK(A) provides an estimate of the number of linearly
4: % independent rows or columns of a matrix A.
5: %
6: % RANK(A,TOL) is the number of singular values of A
7: % that are larger than TOL. By default , TOL = max( size(A)) * eps(norm(A))
8: %
9: % Class support for input A:
10: %     float: double, single
11:
12: % Copyright 1984–2015 The MathWorks, Inc.
13:
14: s = svd(A);
15: if nargin==1
16:     tol = max( size(A)) * eps(max(s));
17: end
18: r = sum(s > tol);
```

type 'edit pinv.m' in Matlab

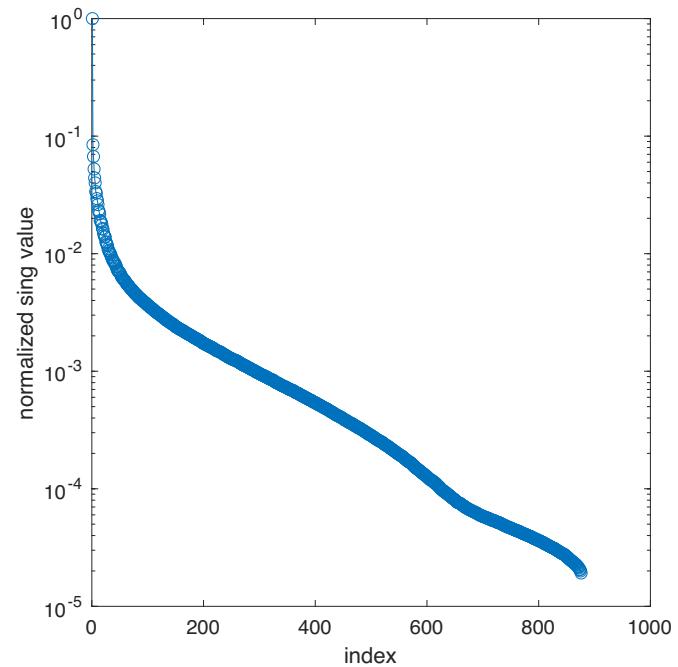
```
1: function X = pinv(A,tol)
2: %PINV    Pseudoinverse.
3: % X = PINV(A) produces a matrix X of the same dimensions
4: % as A' so that A*X*A = A, X*A*X = X and A*X and X*A
5: % are Hermitian. The computation is based on SVD(A) and any
6: % singular values less than a tolerance are treated as zero.
7: %
8: % PINV(A,TOL) treats all singular values of A that are less than TOL as
9: % zero. By default , TOL = max(size(A)) * eps(norm(A)).
10: %
11: % Class support for input A:
12: %     float: double, single
13: %
14: % See also RANK.
15:
16: % Copyright 1984–2015 The MathWorks, Inc.
17:
18: [U,S,V] = svd(A, 'econ');
19: s = diag(S);
20: if nargin < 2
21:     tol = max(size(A)) * eps(norm(s,inf));
22: end
23: r1 = sum(s > tol)+1;
24: V(:,r1:end) = [];
25: U(:,r1:end) = [];
26: s(r1:end) = [];
27: s = 1./s(:);
28: X = (V.*s.*')*U';
```

type 'edit orth.m' in Matlab

```
1: function Q = orth(A)
2: %ORTH    Orthogonalization.
3: %    Q = ORTH(A) is an orthonormal basis for the range of A.
4: %    That is, Q'*Q = I, the columns of Q span the same space as
5: %    the columns of A, and the number of columns of Q is the
6: %    rank of A.
7: %
8: %    Class support for input A:
9: %        float: double, single
10: %
11: %    See also SVD, RANK, NULL.
12: %
13: %    Copyright 1984–2015 The MathWorks, Inc.
14: %
15: [Q,S] = svd(A, 'econ'); %S is always square.
16: s = diag(S);
17: tol = max(size(A)) * eps(max(s));
18: r = sum(s > tol);
19: Q(:, r+1:end) = [];
```

Application of SVD for image compression

```
1: >> A = rgb2gray(imread('llama.jpg'));  
2: >> figure; imshow(A);  
3: >> [U, S, V] = svd(double(A));  
4: >> figure; semilogy(diag(S)/S(1, 1), '-o');  
5: >> xlabel('index'); ylabel('normalized sing value');  
6:  
7: >> r = 50;  
8: >> Aapprx = U(:, 1:r)*S(1:r, 1:r)*V(:, 1:r)';  
9: >> figure; imshow(uint8(Aapprx));
```

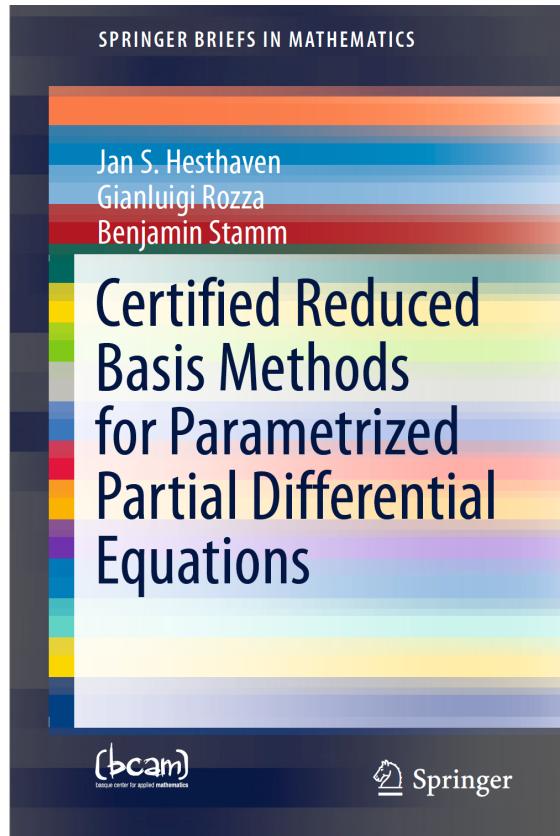


Original picture and singular values



Original picture (left) and reduced picture (right)

Outlook: Model reduction and latent dynamics



Arch Comput Methods Eng manuscript No.
(will be inserted by the editor)

G. Rozza · D.B.P. Huynh · A.T. Patera

Reduced basis approximation and *a posteriori* error estimation for affinely parametrized elliptic coercive partial differential equations

Application to transport and continuum mechanics

Received: August 2007 / Accepted: Date

Abstract In this paper we consider (hierarchical, Lagrange) reduced basis approximation and *a posteriori* error estimation for elliptic coercive output functions of affinely parametrized elliptic coercive partial differential equations. The essential ingredients are (primal-dual) Galerkin projection onto a low-dimensional space associated with a smooth “parametric manifold”—dimension reduction; efficient and effective greedy sampling methods for identification of optimal and numerically stable approximations; rapid convergence; *a posteriori* error bounds; procedures for rigorous a posteriori bounds for the linear-functional outputs of interest; and Offline-Online computational decomposition strategies—minimum *marginal cost* for high performance in the real-time/embedded (e.g., parameter-estimation, control) and many-query (e.g., design optimization, multi-model/scale) contexts. We present illustrative results for heat conduction and convection-diffusion, inviscid flow, and linear elasticity; outputs include transport rates, added mass, and stress intensity factors.

Keywords Partial differential equations, parameter variation, affine geometry description, Galerkin approx-

This work was supported by DARPA/APOSAR Grants FA9550-05-1-0114 and FA9550-07-1-0425, the Singapore-MIT Alliance, the Pappalardo MIT Mechanical Engineering Graduate Monograph Fund, and the Progetto Roberto Rocca Politecnico di Milano-MIT. We are grateful to many fruitful discussions with Professor Yvon Maday of University Paris.

G. Rozza
Massachusetts Institute of Technology, Mechanical Engineering Department, Room 3-264, 77 Mass Avenue, Cambridge, MA 02142-4307, USA. Tel.: +1 617-452-3285; E-mail: rozza@mit.edu

D.B.P. Huynh
National University of Singapore, Singapore-MIT Alliance, National University of Singapore, 4 Eng Drive, Singapore, 117576. Tel.: +65 9132 4000; E-mail: baphuynh@nus.edu.sg

A.T. Patera
Massachusetts Institute of Technology, Room 3-266, 77 Mass Avenue, Cambridge MA, 02142-4307, USA. Tel.: +1 617-253-8122; E-mail: patera@mit.edu

imation, *a posteriori* error estimation, reduced basis, reduced order model, coupling strategies, POD, greedy techniques, offline and online, marginals, marginal cost, coercivity lower bound, successive constraint method, real-time computation, many-query.

1 Introduction and Motivation

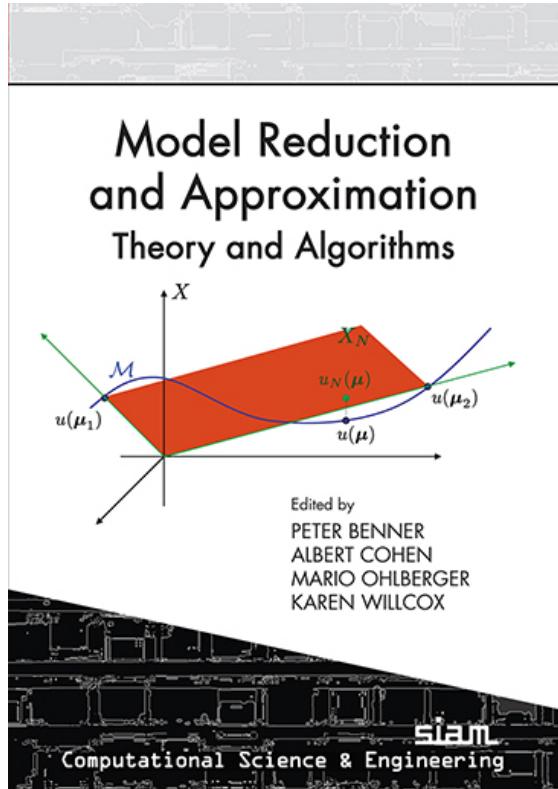
In this work we describe reduced basis (RB) approximation and *a posteriori* error estimation methods for rapid and reliable evaluation of *input-output relationships*, in which the *output* is expressed as a functional of a *field variable*—that is the solution of an *input-parametrized* partial differential equation (PDE). In this particular paper we shall focus on linear output functionals and affinely parametrized linear elliptic coercive PDEs; however the methodology is much more generally applicable, as we discuss below in §2.

We emphasize applications in transport and mechanics; unsteady and steady heat and mass transfer; acoustics;

and solid and fluid mechanics. (Of course we do not preclude other domains of inquiry within engineering (e.g., electromagnetics) or even more broadly within the quantitative disciplines (e.g., finance).)

The *input-parameter* vector typically characterizes the geometric configuration, the physical properties, and the boundary conditions and forces. The *output* of interest includes the maximum system temperature, an added mass coefficient, a crack stress intensity factor, an effective constitutive property, an acoustic waveguide transmission loss, or a channel flowrate or pressure drop. Finally, the *field variables* that connect the input parameters to the outputs can represent a distribution function, temperature or concentration, displacement, pressure, or velocity.

The methodology we describe in this paper is motivated by, optimized for, and applied within two particular contexts: the *real-time context* (e.g., parameter-estimation [54,96,154] or control [124]); and the *many-query context* (e.g., design optimization [107] or multi-model/scale simulation [26,49]). Both these contexts are



Iterative methods for solving systems of linear equations

Why iterative methods for linear systems?

- ▶ Costs of solving a linear system with direct method are $\mathcal{O}(n^3)$. This is too much! If n gets large, then n^3 is huge and costs become intractable.
- ▶ History of matrix computations over the years (according to Trefethen & Bau)
 - 1950: $m = 20$
 - 1965: $m = 200$
 - 1980: $m = 2000$
 - 1995: $m = 20000$

This is an increase of a factor 10^3 . However, computing power (FLOP/sec) increased by about 10^9 . Notice that $(10^3)^3 = 10^9$, which reflects the $\mathcal{O}(n^3)$ bottleneck

- ▶ (Side remark: There are direct methods that beat the complexity $\mathcal{O}(n^3)$ (think of Strassen's algorithm); however,

Why iterative methods for linear systems?

- ▶ Costs of solving a linear system with direct method are $\mathcal{O}(n^3)$. This is too much! If n gets large, then n^3 is huge and costs become intractable.
- ▶ History of matrix computations over the years (according to Trefethen & Bau)

1950: $m = 20$

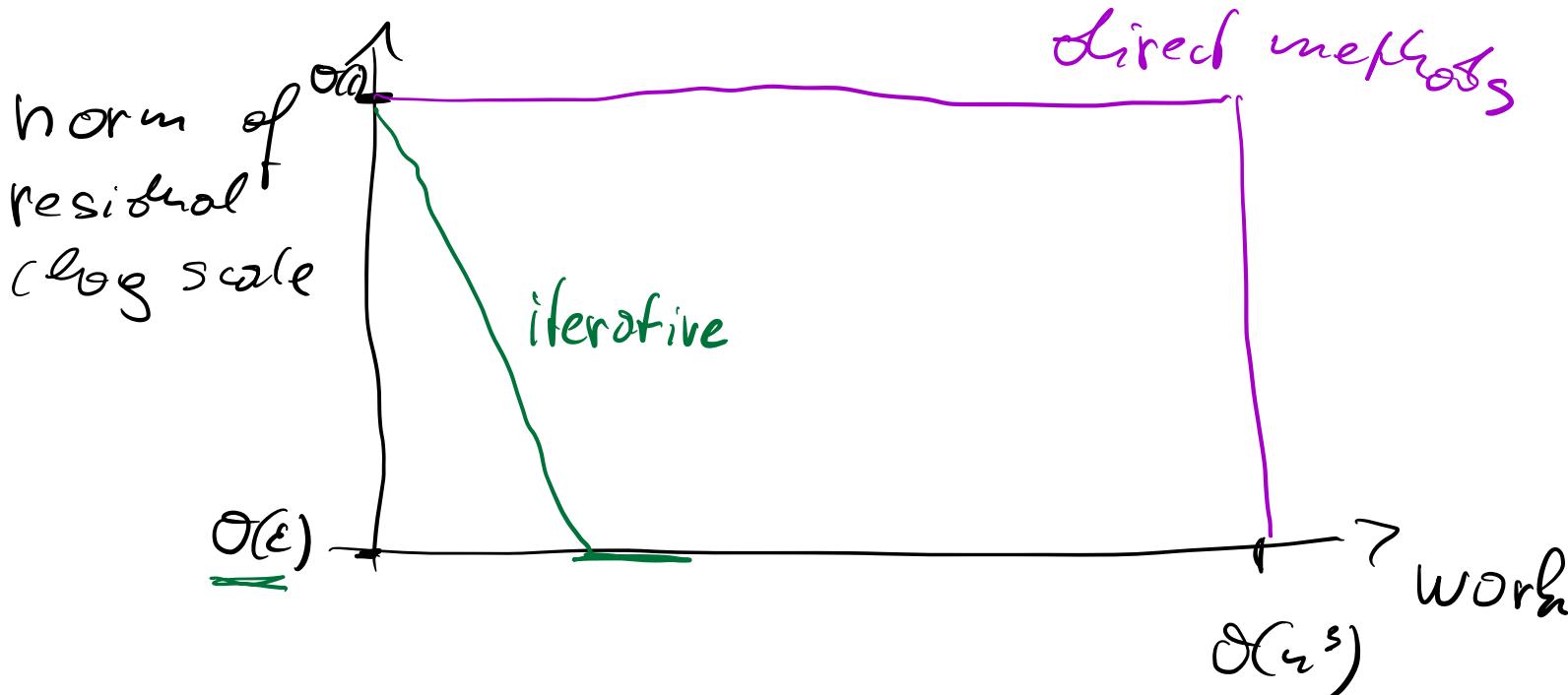
1965: $m = 200$

1980: $m = 2000$

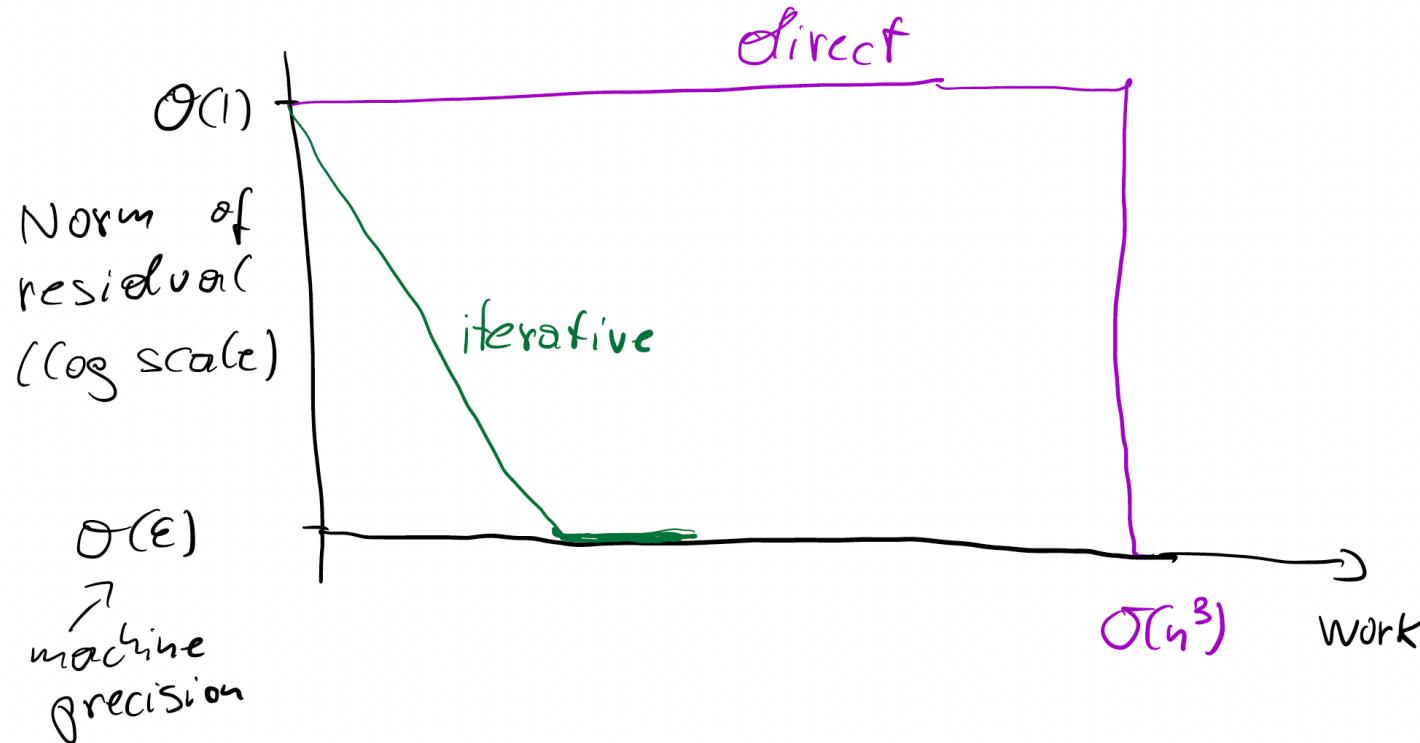
1995: $m = 20000$

This is an increase of a factor 10^3 . However, computing power (FLOP/sec) increased by about 10^9 . Notice that $(10^3)^3 = 10^9$, which reflects the $\mathcal{O}(n^3)$ bottleneck

- ▶ (Side remark: There are direct methods that beat the complexity $\mathcal{O}(n^3)$ (think of Strassen's algorithm); however, the numerical stability of these algorithms is not well understood and the constants hidden in the complexity results are huge so that we “never” see the improved rate in practice.)



- ▶ Iterative methods **can** converge geometrically until residual is below machine precision
- ▶ Direct methods make no progress at all until $\mathcal{O}(n^3)$ work is done, and then lead to residual on the order of machine precision



- ▶ Iterative methods **can** converge geometrically until residual is below machine precision
- ▶ Direct methods make no progress at all until $\mathcal{O}(n^3)$ work is done, and then lead to residual on the order of machine precision

Why iterative methods for linear systems? (cont'd)

- ▶ Classical direct methods (e.g., Gauss elimination) follow the pattern of taking $\mathcal{O}(n)$ steps and each step costs $\mathcal{O}(n^2)$ but don't exploit properties of the matrix
- ▶ Iterative methods reduce the number of steps and the costs of each step, depending on properties of the problem at hand (e.g., spectral properties of A in $Ax = b$)
- ▶ The ideal iterative method requires $\mathcal{O}(1)$ steps and $\mathcal{O}(n)$ costs per step to reach machine precision (think of multigrid and pre-conditioned CG)
- ▶ Iterative methods, even in the absence of rounding errors, do not deliver the exact answer (recall that LU, QR give us exact answer in finite number of steps if we are in exact arithmetic).
- ▶ However, we only can be as accurate as machine precision anyway if we do calculations on a computer. Furthermore, often matrices are stemming from discretizations of PDEs and then we need the solution only up to the discretization error.

Iterative solution of linear systems

Target problems: very large ($n = 10^5, 10^6, \dots$), A is usually sparse and has specific properties.

To solve

$$Ax = b$$

we construct a sequence

$$\underline{x_1, x_2, \dots}$$

of iterates that converges fast to the solution x , where x_{k+1} can be cheaply computed from $\{x_1, \dots, x_k\}$ (e.g., one matrix-vector multiplication).

Thought experiment: If we can compute one iteration with cost $\mathcal{O}(n)$ (e.g., one matrix-vector multiplication with a sparse matrix) and need a constant $\mathcal{O}(1)$ number of iterations to reach desired precision, then we solve $Ax = b$ with costs $\mathcal{O}(n)$. Intuitively, we cannot do better than that because we solve for n quantities and thus need to touch each at least once.

Prototype of iterative method

Let's start by trying to write $\underline{\mathbf{A}\mathbf{x} = \mathbf{b}}$ more generally as

$$\mathbf{x} = f(\mathbf{x})$$

For example set $f(\mathbf{x}) = \underline{(\mathbf{I} - \mathbf{A})\mathbf{x} + \mathbf{b}}$ to obtain

$$\mathbf{x} = f(\mathbf{x}) = (\mathbf{I} - \mathbf{A})\mathbf{x} + \mathbf{b}$$

We can now try to solve this via a fixed-point iteration

$$\mathbf{x}_{k+1} = \underline{f(\mathbf{x}_k)}$$

which is in our case

$$\mathbf{x}_{k+1} = \underline{(\mathbf{I} - \mathbf{A})\mathbf{x}_k + \mathbf{b}}$$

Instead of picking a specific f , let's look at the prototype

$$\underbrace{\mathbf{x}_{k+1}}_{\text{---}} = \underbrace{\mathbf{G}\mathbf{x}_k}_{\text{---}} + \underbrace{\mathbf{c}}_{\text{---}}$$

where $\underline{\mathbf{G}}$ is an iteration matrix somehow related to \mathbf{A} and $\underline{\mathbf{c}}$ is related to $\underline{\mathbf{b}}$.

We must have the actual solution $\underline{\mathbf{x}} = \underline{\mathbf{A}^{-1}\mathbf{b}}$ as a unique fixed point of the iteration:

$$\underbrace{\mathbf{x}}_{\text{---}} = \underbrace{\mathbf{G}\mathbf{x}}_{\text{---}} + \underbrace{\mathbf{c}}_{\text{---}}$$

And we need that the sequence $(\underline{\mathbf{x}}_k)$ converges

Theorem: The fixed point method $\underline{x_{k+1} = Gx_k + c}$ with an invertible G converges for each starting point $\underline{x_0}$ if and only if

$$\rho(G) < 1,$$

where $\rho(G) = \max_j |\lambda_j|$ is the largest magnitude of all eigenvalue of G (i.e., the spectral radius).

~~~ proof

$G$  is real, sym

$$Q G Q^T = D = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$$

Because  $|\lambda_i| \leq p(G) < 1$

$$\lim_{k \rightarrow \infty} \begin{bmatrix} \lambda_1^k & & \\ & \ddots & \\ & & \lambda_n^k \end{bmatrix} = 0$$
$$D^k$$

$$\lim_{k \rightarrow \infty} G^k = \lim_{k \rightarrow \infty} Q D^k Q^{-1} = 0$$

$m \geq 1$

$$\|x_{k+m} - x_k\| = \dots = \|G^k(x_m - x_0)\|$$

$$\xrightarrow[k \rightarrow \infty]{} 0$$

Often want  $\|G\| < 1$ , which implies

$$p(G) < 1$$

*Fix point*

$$\|x_k - x\| \leq \|G\|^k \|x_0 - x\|$$

$\Rightarrow$  smaller  $\|G\|$  means faster convergence

**Theorem:** The fixed point method  $\mathbf{x}_{k+1} = G\mathbf{x}_k + \mathbf{c}$  with an invertible  $G$  converges for each starting point  $\mathbf{x}_o$  if and only if

$$\rho(G) < 1,$$

where  $\rho(G) = \max_j |\lambda_j|$  is the largest magnitude of all eigenvalue of  $G$  (i.e., the spectral radius).

~~~ **proof**

In particular, for any induced matrix norm $\|\cdot\|$ we have $\rho(G) \leq \|G\|$ and thus $\|G\|$ is a sufficient criterion for convergence. We then obtain

$$\|\mathbf{x}_k - \mathbf{x}\| \leq \|G\|^k \|\mathbf{x}_0 - \mathbf{x}\|$$

Let Q be invertible, then

$$\begin{aligned} \mathbf{A}\mathbf{x} = \mathbf{b} &\Leftrightarrow Q^{-1}(\mathbf{b} - \mathbf{A}\mathbf{x}) = 0 \\ &\Leftrightarrow (\mathbf{I} - Q^{-1}\mathbf{A})\mathbf{x} + Q^{-1}\mathbf{b} = \mathbf{x} \\ &\Leftrightarrow \mathbf{G}\mathbf{x} + \mathbf{c} = \mathbf{x} \end{aligned}$$

Leads to fixed-point iteration

$$\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}$$

and $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ is stationary point

$$\mathbf{G}\mathbf{x} + \mathbf{c} = (\mathbf{I} - Q^{-1}\mathbf{A})\mathbf{x} + Q^{-1}\mathbf{b} = \mathbf{x} - Q^{-1}\mathbf{b} + Q^{-1}\mathbf{b} = \mathbf{x}$$

Extreme cases for selecting Q

What are two extreme cases for selecting Q (recall: Q needs to be invertible)?

Extreme cases for selecting \mathbf{Q}

What are two extreme cases for selecting \mathbf{Q} (recall: \mathbf{Q} needs to be invertible)?

Choose $\mathbf{Q} = \mathbf{A}$, then our iteration

$$\mathbf{x}_{k+1} = \mathbf{G}\mathbf{x}_k + \mathbf{c}, \quad \mathbf{G} = \mathbf{I} - \mathbf{Q}^{-1}\mathbf{A}, \quad \mathbf{c} = \mathbf{Q}^{-1}\mathbf{b}$$

becomes

$$(\mathbf{I} - \mathbf{A}^{-1}\mathbf{A})\mathbf{x}_k + \underbrace{\mathbf{A}^{-1}\mathbf{b}}_{\mathbf{x}} = \mathbf{x}_{k+1}$$

$$\mathbf{0} + \mathbf{x} = \mathbf{x}_{k+1}$$

and we are done in just a single step

$$\mathbf{x}_{k+1} = \mathbf{x}$$

Thus, if we “know the solution” (in form of having the inverse \mathbf{A}^{-1}) then no further work is needed here because we already did all the work when finding \mathbf{A}^{-1}

The other extreme is setting $Q = I$, this leads to the Richardson method

$$\mathbf{x}_{k+1} = (\mathbf{I} - \mathbf{A})\mathbf{x}_k + \mathbf{b}$$

We have invested zero costs in finding Q (and Q^{-1}) and so we can expect that $Q = I$ will require high costs in terms of number of iterations to converge in general, if it converges at all

When does Richardson method converge? \rightsquigarrow board

Let A spot, then

$$\rho(G) = \rho(I - A) = \text{not } \{ |1 - \lambda_{\max}(A)|, |1 - \lambda_{\min}(A)|\}$$

$$\lambda_{\max}(A) < 2$$

The Richardson method is consistent (solution is a stationary point) but it may not converge or converge very slowly

$$\mathbf{x}_{k+1} = (\mathbf{I} - \mathbf{A})\mathbf{x}_k + \mathbf{b}$$

What could we do instead of Richardson method?

The Richardson method is consistent (solution is a stationary point) but it may not converge or converge very slowly

$$\mathbf{x}_{k+1} = (\mathbf{I} - \mathbf{A})\mathbf{x}_k + \mathbf{b}$$

What could we do instead of Richardson method?

Let us think of \mathbf{Q} as a preconditioner of the Richardson method:

$$\mathbf{Q}^{-1} \approx \mathbf{A}^{-1}$$

and transform

$$\mathbf{Q}^{-1}\mathbf{A}\mathbf{x} = \mathbf{Q}^{-1}\mathbf{b}$$

Applying the Richardson iteration to this modified system gives

$$\mathbf{x}_{k+1} = (\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A})\mathbf{x}_k + \mathbf{Q}^{-1}\mathbf{b} = \mathbf{G}\mathbf{x}_k + \mathbf{c}$$

with an iteration matrix $\mathbf{I} - \mathbf{Q}^{-1}\mathbf{A} \approx 0$ and thus typically more rapid convergence

Common choices for Q

Requirements on a good Q are

Common choices for Q

Requirements on a good Q are

- ▶ Q^{-1} is representative of A^{-1}
- ▶ We can numerically solve $Qy = z$ much quicker than $Ax = b$ because in each step we solve

$$\boxed{Q\mathbf{x}_{k+1} = (Q - A)\mathbf{x}_k + b}$$

Common choices for Q

Requirements on a good Q are

- ▶ Q^{-1} is representative of A^{-1}
- ▶ We can numerically solve $Qy = z$ much quicker than $Ax = b$ because in each step we solve

$$Q\mathbf{x}_{k+1} = (\mathbf{Q} - \mathbf{A})\mathbf{x}_k + b$$

Split the matrix A as follows

$$A = L + D + U$$

Jacobi method

Theorem: Select now $Q = D \dots$ **Jacobi method**. The Jacobi method converges for any starting point \mathbf{x}_o to the solution of $A\mathbf{x} = \mathbf{b}$ if A is strictly diagonal dominant, i.e.,

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \quad \text{for } i = 1, \dots, n.$$

↔ board

$$A = L + D + R$$

$$x_{k+1} = (I - D^{-1}A)x_k + D^{-1}b$$

$$\begin{aligned} &= (I - D^{-1}(L + D + R))x_k + D^{-1}b \\ &= \underbrace{-D^{-1}(L + R)}_G x_k + D^{-1}b \end{aligned}$$

$$\rho(-D^{-1}(L + R)) = \rho(D^{-1}(L + R))$$

$$\leq \|D^{-1}(L + R)\|_\infty$$

$$= \max_i \underbrace{\left| \sum_{j \neq i} \frac{|a_{ij}|}{|a_{ii}|} \right|}_{\text{Diagonal dom.}} < 1$$

Jacobi method

Theorem: Select now $Q = D \dots$ **Jacobi method**. The Jacobi method converges for any starting point \mathbf{x}_o to the solution of $A\mathbf{x} = \mathbf{b}$ if A is strictly diagonal dominant, i.e.,

$$|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \quad \text{for } i = 1, \dots, n.$$

↔ board

Notice that $Q = D$ is a good choice in terms of computational costs because we can very quickly solve $Dy = z$ for the diagonal matrix D

Gauss-Seidel method

Theorem: Choose $Q = D + L \dots$ Gauss-Seidel method. The Gauss-Seidel method converges for any starting point x_0 if A is symmetric positive definite (spd).

~~ board

Gauss-Seidel method

Theorem: Choose $Q = D + L \dots$ Gauss-Seidel method. The Gauss-Seidel method converges for any starting point x_0 if A is symmetric positive definite (spd).

~ board

Notice that $Q = D + L$ is a good choice in terms of computational costs because we can very quickly solve $(D + L)y = z$ for the lower triangular matrix $D + L$ (forward substitution)