# AMS526: Numerical Analysis I (Numerical Linear Algebra for Computational and Data Sciences)

## Lecture 19: Computing the SVD; Sparse Storage Formats

Xiangmin Jiao

Stony Brook University

# Outline

1. Computing the SVD (NLA§31)

2. Sparse Storage Format

# SVD of $A$ and Eigenvalues of $A^*A$

- Intuitive idea for computing SVD of $A \in \mathbb{C}^{m \times n}$:
  - Form $A^*A$ and compute its eigenvalue decomposition $A^*A = V \Lambda V^*$
  - Let $\Sigma = \sqrt{\Lambda}$, i.e., $\text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \ldots, \sqrt{\lambda_n})$
  - Solve system $U\Sigma = AV$ to obtain $U$
- This method is efficient if $m \gg n$.
- However, it may not be stable, especially for smaller singular values because of the squaring of the condition number
  - For SVD of $A$, $|\tilde{\sigma}_k - \sigma_k| = O(\epsilon_{\text{machine}} \|A\|)$, where $\tilde{\sigma}_k$ and $\sigma_k$ denote the computed and exact $k$th singular value
  - If computed from eigenvalue decomposition of $A^*A$, $|\tilde{\sigma}_k - \sigma_k| = O(\epsilon_{\text{machine}} \|A\|^2/\sigma_k)$, which is problematic if $\sigma_k \ll \|A\|$
- If one is interested in only relatively large singular values, then computing eigenvalues of $A^*A$ is not a problem. For general situations, a more stable algorithm is desired.

# A Different Reduction to Eigenvalue Problem

- Typical algorithm for computing SVD are similar to computation of eigenvalues

- Consider $A \in \mathbb{C}^{m \times n}$, then Hermitian matrix $H = \begin{bmatrix} 0 & A^* \\ A & 0 \end{bmatrix}$ has eigenvalue decomposition

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

where $A = U\Sigma V^*$ gives the SVD. This approach is stable.

- In practice, such a reduction is done implicitly without forming the large matrix

- Typically done in two phases

# Two-Phase Method

- In the first phase, reduce to bidiagonal form by applying different orthogonal transformations on left and right, which involves $O(mn^2)$ operations

- In the second phase, reduce to diagonal form using a variant of QR algorithm or divide-and-conquer algorithm, which involves $O(n^2)$ operations for fixed precision

$$
\begin{bmatrix}
\times & \times & \times & \times \\
\times & \times & \times & \times \\
\times & \times & \times & \times \\
\times & \times & \times & \times \\
\times & \times & \times & \times \\
\times & \times & \times & \times
\end{bmatrix}
\begin{array}{c} \text{Phase 1} \\ \longrightarrow \end{array}
\begin{bmatrix}
\times & \times & & \\
& \times & \times & \\
& & \times & \times \\
& & & \times \\
& & & \\
& & &
\end{bmatrix}
\begin{array}{c} \text{Phase 2} \\ \longrightarrow \end{array}
\begin{bmatrix}
\times & & & \\
& \times & & \\
& & \times & \\
& & & \times \\
& & & \\
& & &
\end{bmatrix}.
$$

$$\quad A \qquad\qquad\qquad B \qquad\qquad\qquad \Sigma$$

We hereafter focus on the first phase

# Golub-Kahan Bidiagonalization

- Apply Householder reflectors on both left and right sides

$$
\begin{bmatrix} \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times \\ \times & \times & \times & \times 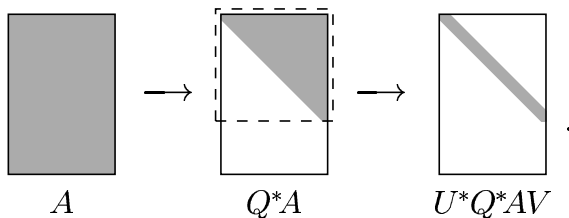\\ \times & \times & \times & \times \end{bmatrix} \quad \xrightarrow{U_1^* \cdot} \quad \begin{bmatrix} \mathbf{\times} & \mathbf{\times} & \mathbf{\times} & \mathbf{\times} \\ \mathbf{0} & \mathbf{\times} & \mathbf{\times} & \mathbf{\times} \\ \mathbf{0} & \mathbf{\times} & \mathbf{\times} & \mathbf{\times} \\ \mathbf{0} & \mathbf{\times} & \mathbf{\times} & \mathbf{\times} \\ \mathbf{0} & \mathbf{\times} & \mathbf{\times} & \mathbf{\times} \\ \mathbf{0} & \mathbf{\times} & \mathbf{\times} & \mathbf{\times} \end{bmatrix} \quad \xrightarrow{\cdot V_1} \quad \begin{bmatrix} \times & \mathbf{\times} & \mathbf{0} & \mathbf{0} \\ & \mathbf{\times} & \mathbf{\times} & \mathbf{\times} \\ & \mathbf{\times} & \mathbf{\times} & \mathbf{\times} \\ & \mathbf{\times} & \mathbf{\times} & \mathbf{\times} \\ & \mathbf{\times} & \mathbf{\times} & \mathbf{\times} \\ & \mathbf{\times} & \mathbf{\times} & \mathbf{\times} \end{bmatrix}
$$

$$
\qquad A \qquad\qquad\qquad U_1^* A \qquad\qquad\qquad U_1^* A V_1
$$

$$
\xrightarrow{U_2^* \cdot} \quad \begin{bmatrix} \times & \times & & \\ & \mathbf{\times} & \mathbf{\times} & \mathbf{\times} \\ & \mathbf{0} & \mathbf{\times} & \mathbf{\times} \\ & \mathbf{0} & \mathbf{\times} & \mathbf{\times} \\ & \mathbf{0} & \mathbf{\times} & \mathbf{\times} \\ & \mathbf{0} & \mathbf{\times} & \mathbf{\times} \end{bmatrix} \quad \xrightarrow{\cdot V_2} \quad \begin{bmatrix} \times & \times & & \\ & \times & \mathbf{\times} & \mathbf{0} \\ & & \mathbf{\times} & \mathbf{\times} \\ & & \mathbf{\times} & \mathbf{\times} \\ & & \mathbf{\times} & \mathbf{\times} \\ & & \mathbf{\times} & \mathbf{\times} \end{bmatrix}.
$$

$$
\qquad\qquad U_2^* U_1^* A V_1 \qquad\qquad\qquad U_2^* U_1^* A V_1 V_2
$$

- Work for Golub-Kahan bidiagonalization $\sim 4mn^2 - \frac{4}{3}n^3$ flops

# Lawson-Hanson-Chan Bidiagonalization

- Speed up by first performing QR factorization on $A$

Lawson–Hanson–Chan bidiagonalization
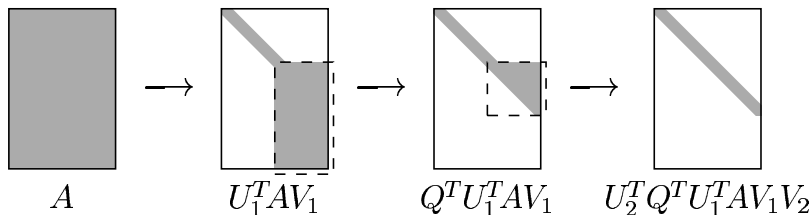


$A$       $Q^*A$       $U^*Q^*AV$

- Work for LHC bidiagonalization $\sim 2mn^2 + 2n^3$ flops, which is advantageous if $m \geq \frac{5}{3}n$

# Three-Step Bidiagonalization

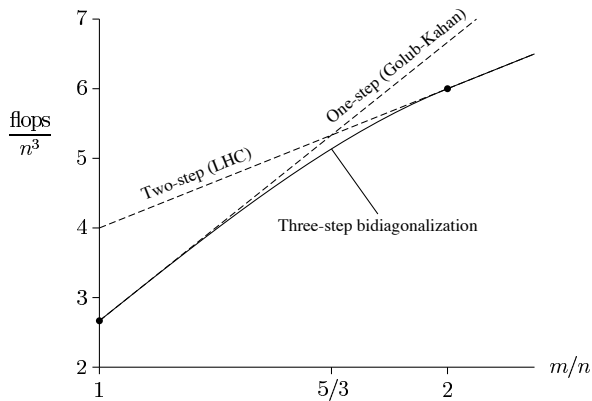- Hybrid approach: Apply QR at suitable time on submatrix with $2 : 1$ aspect ratio

## Three-step bidiagonalization



$$A \qquad U_1^T A V_1 \qquad Q^T U_1^T A V_1 \qquad U_2^T Q^T U_1^T A V_1 V_2$$

- Work for three-step bidiagonalization is $\sim 4mn^2 - \frac{4}{3}n^3 - \frac{2}{3}(m-n)^3$

# Comparison of Performance

# Outline

# Sparse Linear System

- Boundary value problems and implicit methods for time-dependent PDEs yield systems of linear algebraic equations to solve
- A matrix is *sparse* if it has relatively few nonzeros in its entries
- Sparsity can be exploited to use far less than $O(n^2)$ storage and $O(n^3)$ work required in standard approach to solving system with dense matrix, assuming matrix is $n \times n$

# Storage Format of Sparse Matrices

- Sparse-matrices are typically stored in special formats that store only nonzero entries, along with indices to identify their locations in matrix, such as
  - compressed-row storage (CRS)
  - compressed-column storage (CCS)
  - block compressed row storage (BCRS)
- Banded matrices have their own special storage formats (such as Compressed Diagonal Storage (CDS))
- See survey at http:/netlib.org/linalg/html_templates/node90.html
- Explicitly storing indices incurs additional storage overhead and makes arithmetic operations on nonzeros less efficient due to indirect addressing to access operands, so they are beneficial only for very sparse matrices
- Storage format can have big impact the effectiveness of different versions of same algorithm (with different ordering of loops)
- Besides direct methods, these storage formats are also important in implementing iterative and multigrid solvers

# Example of Compressed-Row Storage (CRS)

$$A = \begin{pmatrix} 10 & 0 & 0 & 0 & -2 & 0 \\ 3 & 9 & 0 & 0 & 0 & 3 \\ 0 & 7 & 8 & 7 & 0 & 0 \\ 3 & 0 & 8 & 7 & 5 & 0 \\ 0 & 8 & 0 & 9 & 9 & 13 \\ 0 & 4 & 0 & 0 & 2 & -1 \end{pmatrix} .$$

| val | 10 | -2 | 3 | 9 | 3 | 7 | 8 | 7 | 3 ⋯ 9 | 13 | 4 | 2 | -1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| col_ind | 1 | 5 | 1 | 2 | 6 | 2 | 3 | 4 | 1 ⋯ 5 | 6 | 2 | 5 | 6 |

| row_ptr | 1 | 3 | 6 | 9 | 13 | 17 | 20 |
|---|---|---|---|---|---|---|---|