#### Sparse Matrices in Numerical Mathematics

#### Miroslav Tůma

Faculty of Mathematics and Physics Charles University mirektuma@karlin.mff.cuni.cz

Praha, November 1, 2022

# Outline



#### Introduction

- Introductory notation and terminology
- Sparse matrices and data structures
- Graphs and sparse matrices
- 5 Factorizations
- 6 Reducibility and blocks
- Symbolic Cholesky factorization
- Cholesky factorization
- Sparse LU factorization
- Stability, ill-conditioning, indefiniteness
- Reorderings
- Algebraic preconditioning
- Incomplete factorizations
- Sparse approximate inverses

## Introductory notes

- Created as a material supporting online lectures of NMNV533.
- Assuming basic understanding of algebraic iterative (Krylov space) and direct (dense) solvers (elimination/factorization/solve) (A lot of these is repeated)
- The text deal prevailably with purely algebraic techniques. Such techniques often serve as building blocks for more complex approaches. In particular, some important techniques are mentioned at most. Like:
  - Multigrid/multilevel preconditioners,
  - Domain decomposition,
  - Row projection techniques.
- Only preconditioning of real systems considered here, although extension to complex field is typically straightforward.
- Orientation in variants of Cholesky and LU decompositions is assumed.

The main resource is:

Jennifer Scott and Miroslav Tůma: Algorithms for sparse linear systems, Birkhäuser- Springer, 2022, to appear.

- Printed parts of the resource will be provided to students until it will appear (expected open access then).
- Traditional material also the course text in Czech (nowadays outdated, not supported); see the web page of the course.

# Introductory notes: resources and history of the course

#### • A few other resources:

Davis, T. A. (2006). Direct Methods for Sparse Linear Systems. Fundamentals of Algorithms. SIAM, Philadelphia, PA.

Davis, T. A., Rajamanickam, S., & Sid-Lakhdar, W.M. (2016). A survey of direct methods for sparse linear systems. Acta Numer., 25, 383-566.

Duff, I. S., Erisman, A.M., & Reid, J. K. (2017). Direct Methods for Sparse Matrices (Second ed.). Oxford University Press, Oxford. George, A. & Liu, J. W. H. (1981). Computer Solution of Large Sparse Positive Definite Systems. Prentice Hall, Englewood Cliffs, NJ.

Saad, Y. (2003b). Iterative Methods for Sparse Linear Systems (Second ed.). SIAM, Philadelphia, PA.

#### Most of our activities around solving

$$Ax = b$$

- Direct methods
- Iterative methods
- Practical boundaries between them more and more fuzzy.
- Principially different.

#### Direct methods

- Direct methods: Transform *A* using a finite sequence of elementary transformations: An approach based on factorization (decomposition) and subsequent substitutions.
- The most simple case:  $A \rightarrow LL^T$  or  $LDL^T$  or LUIn principal = Gaussian elimination. Modern (decompositional) form based a lot on the work of Householder (end of 1950's)
  - ► Solving systems with triangular matrices like *L*, *U* is generally much cheaper and more straightforward that using *A*.
  - Factorizations are backbone of direct methods.
  - Occasionally other factorizations than LU or  $LL^T$  or  $LDL^T$
  - Most of the work is in the (Cholesky, indefinite, LU) decomposition.
  - But: also the computer model (sequential, concurrent processors, multicore, GPU) decides about relative complexity of the two steps.
- The algorithms can be made more efficient/stable using additional techniques before, after or during factorization.
- For example, the solution can be made more accurate by an auxiliary iterative method.

#### Iterative methods

Compute a sequence of approximations

 $x^{(0)}, x^{(1)}, x^{(2)}, \dots$ 

that (hopefully) converge to the solution x of the linear system.

- Iterative method are usually accompanied by a problem transformation based on a direct method called preconditioner.
- Usually have to be accompanied by a problem transformation based on a direct method called preconditioner.

#### Iterative methods

• Algebraic preconditioners are tools to convert the problem Ax = b into the one which is easier to solve. They are typically expressed in matrix form as a transformation like:

$$MAx = Mb$$

- *M* can be then used to apply approximation to *A*<sup>-1</sup> to vectors used in the iterative method.
- In practice, it can store approximation to *A* or *A*<sup>-1</sup> (approximate inverse).

#### Contrast: direct versus iterative methods

• Direct methods: designed to be robust, designed to solve

- Properly implemented, they can be used as block-box solvers for computing solutions with predictable accuracy.
- As we have seen, they can be expensive, requiring large amounts of memory, which increases with the size of A.
- Iterative methods: designed to approximate
  - The number of iterations depends on the initial guess  $x^{(0)}$ , A and b
  - Use the matrix A only indirectly, through matrix-vector products  $\rightarrow$  memory requirements are limited to a (small) number of vectors of length the size of A
  - A does not need to be available explicitly.
  - They can be terminated as soon as the required accuracy in the computed solution is achieved.
  - Typically must be preconditioned. Preconditioner computation is sometimes based on a relaxation of a direct method.

Where is the problem with direct methods?

 For example: sparse matrices and resulting factorizations may look like as follows:





Where is the problem with direct methods?

• For example: and they can look like as:



Figure: The locations of the nonzero entries in a symmetric permutation of the matrix from Figure **??** (left) and in  $\overline{L} + \overline{L}^T$  (right), where  $\overline{L}$  is the Cholesky factor of the permuted matrix.

Where is the problem with direct methods?

• For example: and they can look like as:



Figure: The locations of the nonzero entries in a symmetric permutation of the matrix from Figure **??** (left) and in  $\overline{L} + \overline{L}^T$  (right), where  $\overline{L}$  is the Cholesky factor of the permuted matrix.

#### Where is the problem with direct methods?

- We need exploit sparsity (mentioned later)
- See the figures above
- We need sparse (complete) factorizations  $A = LL^T$ , LU (up to the floating-point model)

Where is the problem with iterative methods?

- We must transform (precondition)
- We need sparse (incomplete) factorizations  $A = LL^T$ , LU (up to the floating-point model) like
  - incomplete decompositions ( $A \approx LL^T$ , LU etc.)
  - ▶ incomplete inverse decompositions ( $A^{-1} \approx ZZ^T$ ,  $WZ^T$  etc. )
- Or specific (PDE-based, model-based) approaches.

# Outline

Introductory notation and terminology

Interest in solving linear systems of equations

$$Ax = b, \tag{1}$$

• 
$$A \in \mathbb{R}^{n \times n}, 1 \le i \le n$$
, is nonsingular

- A is sparse
- $b \in \mathbb{R}^n$  (sparse or dense),  $x \in \mathbb{R}^n$  is the solutions
- Used throughout:

$$A = (a_{ij}), \ 1 \le i, j \le n.$$

• Matlab-like notation: nonzero (set a priori),  $A_{:,j}$ ,  $A_{i,:}$ ,  $A_{i:j,k:l}$ ,  $A_{*j}$ ,  $A_{i*}$ .

#### ٩

$$\binom{X}{2} = \{Y \subseteq X \mid |Y| = 2\}$$

- Vectors denoted by small letters as *v*, *u*, *x*, matrices by capital letters as *A*, *B*, . . .
- $A \in \mathbb{R}^{n \times n}, 1 \le i \le n$ , is nonsingular, the right-hand side vector  $b \in \mathbb{R}^n$  is given and  $x \in \mathbb{R}^n$  is the required solution vector. n is the order (or dimension) of A.

# **Basic Terminology**

- A is diagonal if for all  $i \neq j$ ,  $a_{ij} = 0$ ;
- A is lower triangular if for all i < j,  $a_{ij} = 0$
- A is upper triangular if for all i > j,  $a_{ij} = 0$ .
- A is unit triangular if it is triangular and all the entries on the diagonal are equal to one.
- A is structurally symmetric if for all *i* and *j* for which  $a_{ij}$  is nonzero the entry  $a_{ji}$  is also nonzero.
- A is symmetric if

$$a_{ij} = a_{ji}$$
, for all  $i, j$ .

Otherwise, *A* is nonsymmetric.

The symmetry index s(A) of A: the number of nonzeros a<sub>ij</sub>, i ≠ j, for which a<sub>ji</sub> is also nonzero divided by the total number of off-diagonal nonzeros. Small values of s(A): A is far from symmetric.

# Basic Terminology: special matrix classes

• *A* is symmetric positive definite (SPD) if it is symmetric and satisfies

 $v^T A v > 0$  for all nonzero  $v \in \mathbb{R}^n$ .

- Otherwise, *A* is symmetric indefinite.
- Symmetric and (typically) indefinite saddle point matrices have the form

$$A = \begin{pmatrix} G & R^T \\ R & B \end{pmatrix},$$

where  $G \in \mathbb{R}^{n_1 \times n_1}$ ,  $B \in \mathbb{R}^{n_2 \times n_2}$ ,  $R \in \mathbb{R}^{n_2 \times n_1}$  with  $n_1 + n_2 = n$ , G is a SPD matrix and B is a symmetric positive semidefinite matrix (that is  $v^T B v \ge 0$  for all nonzero  $v \in \mathbb{R}^{n_2}$ ). In some applications, B = 0.

• Symmetric block structure of A:

$$A = (A_{ib,jb}), \quad A_{ib,jb} \in \mathbb{R}^{n_i \times n_j}, \quad 1 \le ib, jb \le nb,$$

that is,

$$A = \begin{pmatrix} A_{1,1} & A_{1,2} & \cdots & A_{1,nb} \\ A_{2,1} & A_{2,2} & \cdots & A_{2,nb} \\ \vdots & \vdots & \ddots & \vdots \\ A_{nb,1} & A_{nb,2} & \cdots & A_{nb,nb} \end{pmatrix}$$

- Assuming the square blocks *A*<sub>*jb*, *jb*</sub> on the diagonal are nonsingular.
- Special cases: A is block diagonal if  $A_{ib, jb} = 0$  for all  $ib \neq jb$ , A is block lower triangular if  $A_{1:jb-1, jb} = 0$ ,  $2 \leq jb \leq nb$ , block upper triangular if  $A_{jb+1:nb, jb} = 0$ ,  $1 \leq jb \leq nb 1$ .

#### Definition

Matrix  $A \in \mathbf{R}^{n \times n}$  is **reducible**, if there is a permutation matrix P such that

$$P^{T}AP = \begin{pmatrix} A_{11} & 0\\ A_{21} & A_{22} \end{pmatrix},$$
(3)

where  $A_{11}$  and  $A_{22}$  are square nontrivial matrices (of dimension at least 1). If A is not reducible, it is called **irreducible**. Matrices of dimension 1 are always considered to be irreducible.

#### Remark

Symmetric reducible matrix is block diagonal.

# Basic Terminology: sparsity

- *A* is a sparse matrix if many of its entries are zero.
- Attempts to formalize matrix sparsity more precisely: matrix of order *n* may be said to be sparse if it has *O*(*n*) nonzeros.
- Our choice: *A* is sparse if it is advantageous to exploit its zero entries. Otherwise, *A* is dense.
- The sparsity pattern  $S{A}$  of A is the set of nonzeros, that is,

$$\mathcal{S}\{A\} = \{(i,j) \mid a_{ij} \neq 0, 1 \le i, j \le n\}.$$

- $S{A}$  is symmetric if for all *i* and *j*,  $a_{ij} \neq 0$  if and only if  $a_{ji} \neq 0$  (the values of the two entries need not be the same). If  $S{A}$  is symmetric then *A* is structurally symmetric.
- The number of nonzeros in A: denoted by nz(A) (or |S{A}|). A is structurally (or symbolically) singular if there are no values of the nz(A) entries of A whose row and column indices belong to S{A} for which A is nonsingular.

# Basic Terminology: sparsity

Sparsity: taking into account the structure of matrix nonzeros

#### Definition

Matrix  $A \in \mathbb{R}^{m \times n}$  is said to be sparse if it has  $O(\min\{m, n\})$  entries.



#### Definition

Matrix  $A \in \mathbb{R}^{m \times n}$  is said to be sparse if it has row counts bounded by  $r_{max} \ll n$  or column counts bounded by  $c_{max} \ll m$ .

#### Definition

Matrix  $A \in \mathbb{R}^{m \times n}$  is said to be sparse if its number of nonzero entries is  $O(n^{1+\gamma})$  for some  $\gamma < 1$ .

#### Definition

(pragmatic, application-based definition: J.H. Wilkinson) Matrix  $A \in \mathbb{R}^{m \times n}$  is said to be sparse if we can exploit the fact that a part of its entries is equal to zero.

#### An example showing importance of small exponent $\gamma$ for $n = 10^4$

$\gamma$	$n^{1+\gamma}$			
0.1	25119			
0.2	63096			
0.3	158489			
0.4	398107			
0.5	1000000			

# Rough comparison of dense and sparse (dimension, storage, time for decomposition)

Dense matrix

Sparse matrix

dim	space	dec time (s)	dim	space	dec time (s)
3000	4.5M	5.72	10000	40k	0.02
4000	8M	14.1	90000	0.36M	0.5
5000	12.5M	27.5	1M	4M	16.6
6000	18M	47.8	2M	8M	49.8

# Basic Terminology: sparsity

• Sparse vectors The sparsity pattern of  $v \in \mathbb{R}^n$  is given by

 $\mathcal{S}\{v\} = \{i \mid v_i \neq 0\},\$ 

and  $|\mathcal{S}\{v\}|$  is the length of v.

- Matrix A is factorizable (or strongly regular) if its principal leading minors (the determinants of its principal leading submatrices) are nonzero, that is, if its LU factorization without row/column interchanges does not break down.
- SPD matrices are factorizable.
- For more general *A*, in exact arithmetic the following standard result holds.

#### Theorem

If *A* is nonsingular then the rows of *A* can be permuted so that the permuted matrix is factorizable. The row permutations do not need to be known in advance. They can be constructed on-the-fly as the factorization proceeds.

# Basic Terminology: factorizations

- For symmetric positive definite A, the Cholesky factorization  $A = LL^T$ , where L is a lower triangular matrix with positive diagonal entries.
  - ▶ Rewritten as  $A = \hat{L}D\hat{L}^T$ , where  $\hat{L}$  is a unit lower triangular matrix and D is a diagonal matrix with positive diagonal entries: square root-free Cholesky factorization.
- For nonsymmetric A, the LU factorization A = LU, where L is a unit lower triangular matrix and U is an upper triangular matrix. Gaussian elimination is one process to put a matrix into LU form.
  - Rewritten as A = LDÛ, where Û is a unit upper triangular matrix and D is a diagonal matrix. This is called the LDU factorization.

# Basic Terminology: direct solver phases

- First look: The matrix *A* is factorized and then, given the right-hand side *b*, the factors used to compute the solution *x*.
- Second look:
- Most approaches further split the factorization into a symbolic phase (also called the analyse phase) and a numerical factorization phase that computes the factors.
- The symbolic phase: typically uses only  $S{A}$  to compute the nonzero structure of the factors of *A* without computing the numerical values of the nonzeros.
- The solve phase uses the factors to solve for a single *b* or for a block of multiple right-hand sides or for a sequence of right-hand sides one-by-one.
- Historically, the symbolic phase was much faster than the factorization phase. But parallelising the factorization → timings are much more closer.
- Series of problems in which the numerical values of the entries of *A* change but  $S{A}$  does not: symbolic phase just once.

 Basic sequential model: the von Neumann architecture:union of a central processing unit (CPU) and the memory, interconnected via input/output (I/O) mechanisms.



Figure: A simple uniprocessor von Neumann computer model.

- Nowadays: CPU → a mixture of powerful processors, co processors, cores, GPUs, and so on.
- Furthermore, performing arithmetic operations on the processing units is much faster than communication-based operations.
- Moreover, improvements in the speed of the processing units outpaces those in the memory-based hardware. Moore's Law is an example of an experimentally derived observation of this kind.

- Important milestones in processor development have been multiple functional units that compute identical numerical operations in parallel and data pipelining (also called vectorization) that enables the efficient processing of vectors and matrices.
- Vectorization often supported by additional tools like instruction pipelining, registers and by memory architectures with multiple layers, including small but fast memories called caches.
- Superscalar processors that enable the overlapping of identical (or different) arithmetic operations during run-time have been a standard component of computers since the 1990s.
- The ever-increasing heterogeneity of processing units and their hardware environment inside computers: expressing the code via units of scheduling and execution called threads.

#### Computer-based limitations:

- Compute throughput, that is, the number of arithmetic operations that can be performed per cycle.
- Memory throughput, that is, the number of operands than can be fetched from memory/cache and/or registers each cycle
- Latency, which is the time from initiating a compute instruction or memory request before it is completed and the result available for use in the next computation.
- Distinguishing: algorithms compute-bound, memory-bound or latency-bound.
- More ways to hide latency (blocks, prefetch, threads)
- Measuring computational intensity: the ratio of the number of operations to the number of operands read from memory.
- Most chips designed such that dense matrix-matrix multiply, which typically performs  $n^3$  operations on  $n^2$  data (with ratio k for a blocked algorithm with block size k), can run at full compute throughput, whilst matrix-vector multiply performs  $n^2$  operations on  $n^2$  data (ratio 1) and is limited by the memory throughput.

- The development of basic linear algebra subroutines (BLAS) for performing common linear algebra operations on dense matrices partially motivated by obtaining a high ratio. efficiently.
- Other important motivations behind using the BLAS (standardization, portability).
- Machine-specific optimized BLAS libraries available for a wide variety of computer architectures.

procedure	comm	ops	ratio
BLAS 1: AXPY: $y = y + \alpha x$	3n + 1	2n	2/3
BLAS 2: GEMV: $y = Ax$	$n^2 + 2n$	n(2n-1)	2
BLAS 3: GEMM: $C = AB$	$3n^2$	$n^2(2n-1)$	n/2

 Consequently, exploiting Level 3 BLAS when designing and implementing matrix algorithms (for both sparse and dense matrices) can improve performance compared to using Level 1 and Level 2 BLAS.

# Basic Terminology: finite precision arithmetic

- The IEEE standard (1985) expresses real numbers as  $a = \pm d_1. d_2...d_t \times 2^k$ , where k is an integer and  $d_i \in \{0,1\}, 1 \le i \le t$ , with  $d_1 = 1$  unless  $d_2 = d_3 = ... = d_t = 0$ .
- t = 24 (single precision), t = 53 (double precision), exponent k satisfies  $-126 \le k \le 127$  (single precision) and  $-1022 \le k \le 1023$  (double precision).
- Floating-point (FP) operations:

 $fl(a \ op \ b) = (a \ op \ b)(1+\delta), \qquad |\delta| \leq \epsilon,$ 

(*op* is a mathematical operation (such as  $=, +, -, \times, /, \sqrt{}$ ) and (*a* op b) is the exact result),  $\epsilon$  is the machine epsilon.

- 2 × ε is the smallest FP number which when added to the FP number 1.0 gives a result different from 1.0.
- $\epsilon$  is  $2^{-24} \approx 10^{-7}$  (single precision),  $\epsilon = 2^{-53} \approx 10^{-16}$  (double precision).
- rounding errors, truncation errors.
- catastrophic errors  $\rightarrow$  numerical instability

# Basic Terminology: bit compatibility

- Bit compatibility is essential for some users because of regulatory requirements (such as within the nuclear or financial industries) or to build trust in their software from non technical users.
- The critical issue is the way in which *N* numbers (or, more generally, matrices) are assembled:

$$sum = \sum_{j=1}^{N} S_j,$$

- where the  $S_j$  are computed using one or more processors. The assembly is commutative but, because of the potential rounding of the intermediate results, is not associative so that the result *sum* depends on the order in which the  $S_j$  are assembled.
- A straightforward approach to achieving bit compatibility is to enforce a defined order in such operations.
- This may adversely limit the scope for parallelism.
• The computational complexity of a numerical algorithm is typically based on estimating asymptotically the number of integer or floating-point operations or the memory usage.

#### Definition

A real function f(k) of a nonnegative real k satisfies f = O(g) if there exist positive constants  $c_u$  and  $k_0$  such that

$$f(k) \le c_u g(k) \text{ for all } k \ge k_0.$$
(4)

We say that  $f = \Theta(g)$  if, additionally, there exists a positive constant  $c_l$  such that

 $0 \le c_l g(k) \le f(k) \le c_u g(k)$  for all  $k \ge k_0$ .

- While O(g) bounds f asymptotically from above, Θ(g) represents an asymptotically tight bound.
- As a simple illustration, consider the quadratic function

$$f(k) = \alpha * k^2 + \beta * k - \gamma.$$

Provided  $\alpha \neq 0$ ,  $f(k) = \Theta(k^2)$  and the coefficient of the highest asymptotic term is  $\alpha$ . Computational complexity can estimate quantities related to the worst-case behaviour of an algorithm (worst-case complexity), or it can express average behaviour (average-case complexity).

• Unit costs, Sparse matrix algorithms that are  $\Theta(n^3)$  are considered to be computationally expensive.

### **Complexity here and in CS**

- Because of the development in computations, MFLOPs may be misleading
- Still terminology O(.) (bounding from above) or Θ(.) (bounding from both sides) sometimes relevant consists in replacing the bound (bounds) by constant × simpler function (etalon).
- Simpler functions are, e.g.,  $n^2, n^3, \log n, \ldots$
- Distinguish worst case and average case analysis
- Inverse Ackermann function will be introduced in exercises
- In CS: polynomial complexity versus superpolynomial complexity.
   Our case: even n<sup>3</sup> may be too much.
- Decision problems, polynomial reduction, class  $\mathcal{NP}$ , etc.

- Absolutely crucial for direct methods: complexity for generally dense matrices, sequential case:  $O(n^3)$  factorization,  $O(n^2)$  substitutions
- Useful for iterative methods as well: repeated multiplications and solve steps. But, expecting rather sparse matrices and (typically) dense vectors.
- Complexity in the sparse case depends on the decomposition model and computer architecture (implementation, completeness/incompleteness)

## Outline



### Sparse vector in a computer

### Example

Consider the sparse row vector  $v \in \mathbb{R}^8$ 

$$v = (1. -2. 0. -3. 0. 5. 3. 0.).$$
 (5)

The real array valV that stores the nonzero values and corresponding integer array of their indices indV are of length  $|S\{v\}| = 5$  and are as follows:

Subscripts	1	2	3	4	5
valV	1.	-2.	-3.	5.	3.
indV	1	2	4	6	7

#### Sparse vector in a computer

- Alternatively, a linked list can be used.
- linked list based format: stores matrix rows/columns as items connected by pointers
- Iinked lists can be cyclic, one-way, two-way
- A figure for demonstration, only values (not their indices) are shown



 rows/columns embedded into a larger array: emulated dynamic behavior

# Sparsity

#### Sparse vector in a computer

### • Linked list can be embedded into a large array.

### Example

Two possible ways of storing the sparse vector using linked lists.

Subscripts	1	2	З	3 4	5
Values	1.	-2.	_;	3. <b>5</b> .	3.
Indices	1	2	4	6	7
Links	2	3	4	- 5	0
Header	1				
Subscripts	1	2	3	4	5
Values	5.	3.	1.	-2.	-3.
Indices	6	7	1	2	4
Links	2	0	4	5	1
Header	3				

• Reasons for using linked lists: straightforward adds and removes.

### Example

On the left, an entry -4 has been added in position 5. On the right, an entry -2 in position 2 has been removed. \* indicates the entry is not accessed. The links that have changed are in bold.

Subscripts	1	2		3	4	5	6
Values	1.	-2.	-	-3.	5.	3.	-4.
Indices	1	2		4	6	7	5
Links	2	3		4	5	6	0
Header	1						
Subsc	ripts	1	2	3	4	ĻĘ	5
Values	5	1.	*	-3.	5	. 3	3.
Indice	s	1	*	4	6	5 7	7
Links		3	*	4	5	5 (	)
Heade	er	1					

#### Sparse matrix storage

- coordinate (or triplet format: the individual entries of A are held as triplets  $(i, j, a_{ij})$ , where *i* is the row index and *j* is the column index of the entry  $a_{ij} \neq 0$ . (dynamic storage format)
- CSR (Compressed Sparse Row) format. The column indices of the entries of A held by rows in an integer array (which we will call colindA) of length nz(A), with those in row 1 followed by those in row 2, and so on (with no space between rows). Sorted or unsorted. (static storage format)
- CSC (Compressed Sparse Columns): analogously by columns instead of rows.
- If *A* is symmetric, only the lower (or upper) triangular part is generally stored.
- Possible to store only  $\mathcal{S}{A}$ .

Sparse matrix in the coordinate format

• Example matrix  $A \in \mathbb{R}^{5 \times 5}$ 

### Example

Subscripts	1	2	3	4	5	6	7	8	9	10
rowindA	3	2	3	4	1	1	2	5	3	5
colindA	3	2	1	4	4	1	5	5	5	2
valA	3.	1.	-1.	1.	-2.	3.	4.	6.	1.	7.

Sparse matrix stored using linked lists

• Easy adding and deleting entries is possible if t linked lists are used: the matrix held as a collection of columns, each in a linked list. colA\_head holds header pointers.

Exam	ple										
-	Subscripts	1	2	3	4	5	6	7	8	9	10
	rowindA	3	2	3	4	1	1	2	5	3	5
	valA	3.	1.	-1.	1.	-2.	3.	4.	6.	1.	7.
	link	0	10	0	0	4	3	9	0	8	0
	colA_head	6	2	1	5	7					

If we consider column 4, then colA\_head(4) = 5, rowindA(5) = 1 and valA(5) = -2., so the first entry in column 4 is  $a_{1,4} = -2$ . Next, link(5) = 4, rowindA(4) = 4 and valA(4) = 1., so the next entry in column 4 is  $a_{4,4} = 1$ ..

#### Sparse matrix in the CSR format

• CSR format represents *A* as follows. Here the entries within each row are in order of increasing column index.

Example	Э											
Su	ubscripts	1	2	3	4	5	6	7	8	9	10	
ro	wptrA	1	3	5	8	9	11					
со	lindA	1	4	2	5	1	3	5	4	2	5	
va	lA	3.	-2.	1.	4.	-1.	3.	1.	1.	7.	6.	

In our codes we often use: ia: rowptrA, ja: colindA, aa: valindA

Sparse matrix: static versuis dynamic formats

- o dynamic data structures:
  - more flexible but this flexibility might not be needed
  - difficult to vectorize
  - difficult to keep spatial locality
  - used preferably for storing vectors
- static data structures:
  - ad-hoc insertions/deletions should be avoided (better algorithms)
  - much simpler to vectorize
  - efficient access to rows/columns

### Simulating dynamic storage format

- A disadvantage of linked list storage: prohibits the fast access to rows (or columns) of the matrix. And this is needed!
- Simulated dynamism of storage schemes: storage format with some additional elbow space for new non zero entries of *A* is needed.
- Often the case in approximate factorizations where new non zero entries can be added and/or removed and it is hard to predict the necessary space in advance.
- In this case, the elbow space can embed new non zeros.
- The format is called the DS format.

#### Sparse matrix: DS formats

• Consider again the sparse matrix  $A \in \mathbb{R}^{5 \times 5}$  (6). The DS format represents A as follows.

Example	e
---------	---

Subscripts	1	2	3	4	5	6	7	8	9	10	11	12	13	14
rowptrA	1	5	8	12	14									
colindA	1	4			2	5		1	3	5		4		2
valAR	3.	-2.			1.	4.		-1.	3.		1.		1.	7.
rowlength	2	2	3	1	2									
colptrA	1	4	6	9	12									
rowindA	1	3		2	5	3			1	4		2	3	5
valAC	3.	-1.		1.	7.	3.			-2.	1.		4.	1.	6.
collength	2	2	1	2	3									

#### Sparse matrix: DS formats

- It can happen that the free space between row and/or column segments disappears throughout a computational algorithm. Then the DS format must be reorganized.
- In particular, a row segment can be moved to the end of the arrays valAR and colindA implying also a corresponding update in rowptrA. The space where the row *i* originally resided is then denoted as free.
- If there is no free space at the end of the arrays valAR and colindA, a compression of the row segments or full reallocation should be done.
- While the DS format seems to be complicated, it can be extremely useful in some cases. Surprisingly efficient if the amount of changes is limited as it often is in approximate factorizations.

#### **Block formats**

- Blocked formats may be used to accelerate multiplication between a sparse matrix and a dense vector.
- The Variable Block Row (VBR) format groups together similar adjacent rows and columns.
- The data structure of the VBR format uses six arrays. Integer arrays rptr and cptr hold the index of the first row in each block row and the index of the first column in each block column. respectively. In many cases, the block row and column partitionings are conformal and only one of these arrays is needed. The real array vala contains the entries of the matrix block-by-block in column-major order. The integer array indx holds pointers to the beginning of each block entry within valA. The index array bindx holds the block column indices of the block entries of the matrix and, finally, the integer array bptr holds pointers to the start of each row block in bindx.

### Sparse matrix: DS formats

### Example

Consider the sparse matrix  $A \in \mathbb{R}^{8 \times 8}$ 

	1	2	3	4	5	6	7	8
1	/ 1.	2.				3.		
2	4.	5.				6.		
3			7.	8.	9.	10.		
4	11.	12.					15.	16.
5		13.					17.	
6	14.							18.
7			19.		20.			
8			21.	22.				)

Here the row blocks comprise rows 1:2, 3, 4:6 and 7:8. The column blocks comprise columns 1:2, 3:5, 6, 7:8. The VBR format stores A as follows.

#### Sparse matrix: DS formats

Example																	
Subscripts	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
rptr	1	3	4	7	9												
cptr	1	3	6	7	9												
valA	1.	4.	2.	5.	3.	6.	7.	8.	9.	10.	11.	14.	12.	13.	15.	17.	16.
indx	1	5	7	10	11	15	19										
bindx	1	3	2	3	1	4	2										
bptr	1	3	5	7													

Matmats in CSR/CSC

1) CSR - CSC

$$C = AB, A = \begin{pmatrix} a_1 \\ \vdots \\ a_m \end{pmatrix}, B = (b_1, \dots, b_n), C = (c_{ij})$$
(7)

- Each entry  $c_{ij}$  computed as a product of a compressed row of A and compressed column of B
- Not clear whether the result  $c_{ij}$  is nonzero
- Consequently:  $O(n^3)$  operations, not useful for sparse matrices.

Matmats in CSR/CSC

2) CSR - CSR

$$C = AB, A = \begin{pmatrix} a_1 \\ \vdots \\ a_m \end{pmatrix}, B = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}, C = (c_{ij})$$

В

	a <sub>i</sub>	
* *	* *	*

(8)

Matmats in CSR/CSC

#### 2) CSR - CSR



Matmats in CSR/CSC

3) CSC - CSR

$$C = AB, A = \begin{pmatrix} a_1, \dots, a_m \end{pmatrix}, B = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}, C = (c_{ij})$$
(10)

- How one can store A by CSC and pass it by rows?
- Pointers to first entries in columns: (array *first*)
- First test: nonzero in the first row → move one step down, add next nonzero into the list value(next)
- Complexity: O(nonzeros) + O(n)

### Matmats in CSR/CSC

### 3) CSC - CSR



### Based on forming virtual rows in A

## Outline

Graphs and sparse matrices

### Directed and undirected graphs

- Many sparse matrix algorithms exploit the relationship matrices ↔ graphs.
- A graph G = (V, E) is a finite set V of vertices (or nodes), and a set E of edges defined as pairs of distinct vertices.
- No distinction between the pairs (u, v) and (v, u): the edges are represented by unordered pairs, the graph is <u>undirected</u>.

### Definition

Simple undirected graph *G* is an ordered pair of sets (V, E), where  $V = \{v_1, \ldots, v_n\}$  is called the set of vertices of *G*,  $E = \{e_1, \ldots, e_m\}$  is called the set of edges satisfying (no multiple edges, no loops)

$$E \subseteq \binom{V}{2}.$$

• If the pairs are ordered: directed graph (a digraph).

## Graphs and sparse matrices

#### Directed and undirected graphs



Figure: An example of an undirected graph.



Figure: An example of a directed graph (digraph). The arrows indicate the direction of an edge. There is an edge  $(4 \rightarrow 5)$  and an edge  $(5 \rightarrow 4)$ .

### **Bipartite graphs**

### Definition

A simple bipartite graph is an ordered pair of sets (R, C, E) such that  $E = \{\{i, j\} | i \in R, j \in C\}$ . *R* is called the row vertex set, *C* is called the column vertex set and *E* is called the edge set.



## Graphs and their matrices

### Graph terminology

- A labelling (or ordering) of a graph G = (V, E) with n vertices is a bijection of {1, 2, ..., n} onto V. The integer i (1 ≤ i ≤ n) assigned to a vertex in V is called the label (or simply the number).
- Our standard choice of vertices will be  $\mathcal{V} = \{1, \dots, n\}$  so that the vertices are directly identified by their labels.
- Another example of a labelled undirected graph:



Simple undirected graph  $G = (\{1, 2, 3, 4, 5, 6\}, \{\{1, 2\}, \{2, 3\}, \{1, 4\}, \{3, 4\}, \{3, 6\}, \{3, 5\}, \{5, 6\}\})$ 

- $\mathcal{G}_s = (\mathcal{V}_s, \mathcal{E}_s)$  is a subgraph of  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  if and only if  $\mathcal{V}_s \subseteq \mathcal{V}$  and  $\mathcal{E}_s \subseteq \mathcal{E}$  and  $(u_s, v_s) \in \mathcal{E}_s$  implies  $u_s, v_s \in \mathcal{V}_s$ .
- The subgraph is an induced subgraph if  $\mathcal{E}_s$  contains all the edges in  $\mathcal{E}$  that have both u and v in  $\mathcal{V}_s$ .
- Two graphs  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and  $\mathcal{G}_s = (\mathcal{V}_s, \mathcal{E}_s)$  are isomorphic if there is a bijection  $g : \mathcal{V} \to \mathcal{V}_s$  that preserves adjacency, that is  $(u, v) \in \mathcal{E}$  if and only if  $(g(u), g(v)) \in \mathcal{E}_s$ .
- Undirected graph: two vertices u and v in  $\mathcal{V}$  are said to be adjacent (or neighbours) if  $e = (u, v) \in \mathcal{E}$ ; the edge e is incident to the vertex u and to the vertex v. u and v are the endpoints of e.
- We also use the notation  $(u \leftrightarrow v)$  for an edge (or  $(u \leftrightarrow v)$  to emphasise the edge belongs to the graph  $\mathcal{G}$ ).

- The degree deg<sub>G</sub>(u) of u ∈ V is the number of vertices in V that are adjacent to u, and the adjacency set adj<sub>G</sub>{u} is the set of these adjacent vertices (thus |adj<sub>G</sub>{u}| = deg<sub>G</sub>(u)).
- If V<sub>s</sub> is a subset of the vertices, then the adjacency set adj<sub>G</sub>{V<sub>s</sub>} is the set of vertices in V \ V<sub>s</sub> that are adjacent to at least one vertex in V<sub>s</sub>.
- A subgraph is a clique when every pair of vertices is adjacent. In the example in Figure 4, deg<sub>G</sub>(2) = 4 and adj<sub>G</sub>{2} = {1,3,4,6}. The induced subgraph with vertices V<sub>s</sub> = {2,4,6} is a clique.

## Graphs and their matrices

### Graph terminology

- Notation  $(u \rightarrow v)$  for a directed edge.
- Emphasising the graph to which the edge belongs:  $(u \xrightarrow{\mathcal{G}} v)$ .
- In a digraph there can be an edge  $(u \rightarrow v)$  but no edge  $(v \rightarrow u)$ .
- The adjacency set of u can be split into two parts

$$adj_{\mathcal{G}}^{+}\{u\} = \{v \mid (u \to v) \in \mathcal{E}\} \quad \text{and} \quad adj_{\mathcal{G}}^{-}\{u\} = \{v \mid (v \to u) \in \mathcal{E}\}.$$

In the example given in Figure:  $adj_{\mathcal{G}}^+\{2\} = \{3, 4\}, adj_{\mathcal{G}}^-\{2\} = 1.$ 

• A sequence of k edges in an undirected graph G: a walk.

$$u_0 \longleftrightarrow u_1 \longleftrightarrow \ldots \longleftrightarrow u_{k-1} \longleftrightarrow u_k$$

• If G is a digraph then the sequence: a directed walk.

$$u_0 \longrightarrow u_1 \longrightarrow \ldots \longrightarrow u_{k-1} \longrightarrow u_k$$

• k is the length of the walk. A walk of zero length: k = 0.

- The vertices  $u_0$  and  $u_k$  are connected by the walk and for k > 0,  $u_k$  is said to be reachable from  $u_0$ ;
- The set of vertices that are reachable from  $u_0$  is denoted by  $\mathcal{R}each(u_0)$ . The walk is closed if  $u_0 = u_k$ ; a closed walk is called a cycle.
- Graphs that do not contain cycles are acyclic.
- A (directed) trail is a (directed) walk in which all the edges are distinct and a (directed) path is a (directed) trail in which all the vertices (and therefore also all the edges) are distinct.
- The distance between two vertices is the number of edges in the shortest path connecting them (this is also called the length of the path).

 If the vertices V are labelled 1, 2, ..., n then in the undirected graph G = (V, E) a path between a pair of its vertices with labels i and j is denoted by

$$i \xleftarrow{\mathcal{G}} j$$

or, if it is clear which graph the path is in, by

$$i \Longleftrightarrow j.$$

If all intermediate vertices on the path are less than  $\min\{i, j\}$  then the path is called a fill-path and is denoted by

$$i \xleftarrow{\mathcal{G}} j \quad \text{or} \quad i \xleftarrow{\min} j.$$

 If all intermediate vertices on the path belong to a subset V<sub>s</sub> then the path is denoted by

$$i \xleftarrow{\mathcal{G}} j \quad \text{or} \quad i \xleftarrow{\mathcal{V}_s} j$$

If  $\mathcal{G}$  is a digraph, the double-sided arrow symbols are replaced by one-sided ones  $\implies$  in the direction of the edges. For example,

$$i \xrightarrow{\mathcal{G}} j, \quad i \Longrightarrow j, \quad i \xrightarrow{\min} j \quad \text{and} \quad i \xrightarrow{\mathcal{V}_s} j.$$
### Graph terminology

- Digraph with no cycles (directed acyclic graph): DAG.
- In a DAG, if there is a path u ⇒ v then u is called an ancestor of v and v is said to be a descendant of u.
- DAG with two different orderings. Left: vertices 2, 3, 5 and 6 are descendants of 1. Only vertices 5 and 6 are descendants of vertex 4.



Figure: An example of a DAG with two different topological orderings.

### Graph terminology: reachability

- Given a graph and a subset V<sub>s</sub> of vertices, if u and v are two distinct vertices that do not belong to V<sub>s</sub>, then v is reachable from u through V<sub>s</sub> if u and v are connected by a path that is either of length 1 or is composed entirely of vertices that belong to V<sub>s</sub> (except for the endpoints u and v).
- Given V<sub>s</sub> and u ∉ V<sub>s</sub>, the reachable set Reach(u, V<sub>s</sub>) of u through V<sub>s</sub> is the set of all vertices that are reachable from u through V<sub>s</sub>. Note that if V<sub>s</sub> is empty or u does not belong to adj<sub>G</sub>(V<sub>s</sub>) then Reach(u, V<sub>s</sub>) = adj<sub>G</sub>(u).
- A simple example is given in Figure.

### Graph terminology: reachability



Figure: An example of an undirected graph to illustrate reachability. If  $\mathcal{V}_s = \{4, 5\}$  then  $\mathcal{R}each(2, \mathcal{V}_s) = \{1, 3, 6\}$  and  $\mathcal{R}each(6, \mathcal{V}_s) = \{2, 3, 7\}$ .

## Graph terminology: connectedness, trees

- An undirected graph is connected if every pair of vertices is connected by a path.
- A connected acyclic graph is called a tree, that is, a tree is an undirected graph in which any two vertices are connected by exactly one path.
- Every tree has at least two vertices of degree 1. Such vertices are called leaf vertices.
- A graph is a forest if it consists of a disjoint union of trees. This is illustrated in Figure. Connectivity is an equivalence relation and consequently, it provides a partition of V into disjoint equivalence classes.

#### Graph terminology: connectedness, trees



Figure: An example of an undirected graph with 12 vertices that is a forest (it consists of two disjoint trees). Vertices 1, 2, 3, 6, 7, 8 and 11 are leaf vertices.

### Graph terminology: connected, strongly connected

- If G is connected then a spanning tree of G is a subgraph of G that is a tree containing every vertex of G.
- A directed graph G = (V, E) is strongly connected if for every pair of vertices u, v ∈ V there is a path from u to v and a path from v to u.
- Strong connectivity is an equivalence relation on V. It induces a partitioning V = V<sub>1</sub> ∪ ... ∪ V<sub>s</sub> such that each V<sub>i</sub> (1 ≤ i ≤ s) is strongly connected and is maximal with this property: no additional vertices from G can be included in V<sub>i</sub> without breaking its strong connectivity.
- The  $V_i$  are called strongly connected components (or sometimes just strong components) of G.

### Graph terminology: rooted trees

- Any undirected tree \$\mathcal{T} = (\mathcal{V}, \mathcal{E})\$ can be converted to a directed rooted tree \$\mathcal{T}' = (\mathcal{V}, \mathcal{E}')\$ by specifying a root vertex \$r\$.
- r can be chosen arbitrarily: any choice gives a directed rooted tree. An edge with endpoints u and v in E becomes a directed edge (u → v) in E' if there is a path from u to r such that the first edge of this path is from u to v.
- A rooted tree is a special case of a DAG.
- v is called the parent of u if the directed edge (u → v) ∈ E'; u is said to be a child of v (two or more child vertices are referred to as children). Two vertices in a rooted tree are siblings if they have the same parent. Leaf vertices have no children.
- Given *r*, this directed path is unique.

### Graph terminology: rooted trees



Figure: An example of a an undirected tree  $\mathcal{T}$  (left) and the rooted tree  $\mathcal{T}'$  (right) obtained from  $\mathcal{T}$  by choosing the root r = 4. The arrows indicate the direction of the edges.

## Another example of a rooted tree



- The root of this tree is 12
- Then, for example, 10 is an ancestor of vertices 2, 8 a 9. These vertices are descendants of 10. Set of ancestors of 10 is  $anc(10) = \{10, 11, 12\}$ . parent(i) for vertices  $1 \dots 11$  is 8, 10, 7, 7, 6, 9, 9, 10, 10, 11, 12, null.

- Adjacency graphs provide a link between sparse matrices and graphs. If A is a sparse matrix of order n then an adjacency graph G(A) (often written simply as G) with n vertices V(A) = {1,...,n} can be associated with it.
- If A is structurally symmetric, then the edge set is

 $\mathcal{E}(A) = \{(i,j) \mid a_{ij} \neq 0, \ i \neq j\}.$ 

A digraph can be associated with a nonsymmetric A by setting

$$\mathcal{E}(A) = \{ (i \to j) \mid a_{ij} \neq 0, \ i \neq j \}.$$

Each diagonal nonzero  $a_{ii}$  corresponds to a loop or self edge. Loops are generally omitted from  $\mathcal{G}$  since many algorithms that use A assume that the diagonal entries of A are present.

 To capture not only the sparsity pattern of A but also the values of the entries, G can be transformed into a weighted graph using a mapping E(A) → ℝ and/or V(A) → ℝ.

## Adjacency matrix of an undirected graph

### Definition

For a simple undirected graph G = (V, E) with  $V = \{1, ..., n\}$  the adjacency matrix is the (0, 1) matrix  $A_G = (a_{ij})$   $(i, j \in V)$ , where  $a_{ij}$  is 0 if there is an edge  $\{i, j\}$  in E where  $i, j \in V$  and zero otherwise.

Adjacency matrix for the undirected graph above.

### Incidence matrix of an undirected graph

# Definition

For a simple undirected graph G = (V, E) with  $V = \{1, ..., n\}$  and  $E = \{1, ..., m\}$  the incidence matrix is the (0, 1) matrix  $A_G = (a_{ij})$   $(i \in \{1, ..., m\}, j \in \{1, ..., n\})$ , where  $a_{ij} = 1$ , if j is the vertex of the edge i and  $a_{ij} = 0$  otherwise.

- Incidence matrix is generally rectangular.
- Several slightly differing definitions of adjacency and incidence matrices.
- Various generalizations.

Transfer between the classes of undirected and directed graphs

- Symmetrization: directed  $\rightarrow$  undirected
  - ▶ Just considering edges from  $V \times V$  as from  $\begin{pmatrix} V \\ 2 \end{pmatrix}$
- Orientation: undirected  $\rightarrow$  directed
  - ▶ Not unique. Instead from an edge from  $\binom{V}{2}$  we can have one or two edges from  $V \times V$ .
- In any case, part of terminology is shared between the classes of undirected and directed graphs



Figure: An example of a structurally symmetric sparse matrix and its undirected graph (left) and a nonsymmetric sparse matrix and its digraph (right). Arrows indicate the direction of the edges in the digraph.

### Graphs of triangular matrices

 A special case is the directed graph associated with a triangular matrix. If L is a lower triangular matrix and U is an upper triangular matrix then the directed graphs G(L) and G(U) have edge sets

$$\mathcal{E}(L) = \{ (i \to j) | l_{ij} \neq 0, \ i > j \}$$
  
$$\mathcal{E}(U) = \{ (i \to j) | u_{ij} \neq 0, \ i < j \}$$

• It is sometimes convenient to use  $\mathcal{G}(L^T)$  in which the direction of the edges is reversed

$$\mathcal{E}(L^T) = \{ (j \to i) \, | \, l_{ij} \neq 0, \ i > j \}.$$
(11)

It is straightforward to see that  $\mathcal{G}(L)$ ,  $\mathcal{G}(L^T)$  and  $\mathcal{G}(U)$  are DAGs.

## Permutation matrices

- A permutation matrix is a square matrix that has exactly one entry equal to 1 in each row and column, and all remaining entries are zeros (that is, it is a permutation of the identity matrix).
- Premultiplying a matrix by *P* reorders the rows and postmultiplying by *P* reorders the columns. *P* can be represented by an integer-valued permutation vector *p*, where *p<sub>i</sub>* is the column index of the 1 within the *i*-th row of *P*. For example,

$$P = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix} \text{ and } p = \begin{pmatrix} 2 \\ 3 \\ 1 \end{pmatrix}$$

• The graph of a matrix A is unchanged if a symmetric permutation  $A' = PAP^T$  is performed, only the labelling of the vertices changes and thus relabelling  $\mathcal{G}(A)$  can be used to permute A.

Graphs of arrowhead matrices

Consider the arrowhead matrix A and its graph  $\mathcal{G}(A)$ . The symmetrically permuted matrix A' and  $\mathcal{G}(A')$  are also shown, with P such that the first row and column of A are the last row and column of A'.



89/609

## Topological orderings of digraphs

- The digraph of a general matrix A is not invariant under nonsymmetric permutations PAQ, with Q ≠ P<sup>T</sup>. A topological ordering of a directed graph is a labelling of its vertices such that for every edge (i → j), vertex i precedes vertex j (i.e., i < j).</li>
- A topological ordering is possible if and only if the graph has no directed cycles, that is, if it is a DAG. Any DAG has at least one topological ordering.
- Non-unique (see above).

### Lists

# • A list is an ordered sequence of arbitrary elements

$$(u_0, u_1, \dots, u_{k-1}, u_k)$$
 (12)

u<sub>0</sub> is the head of the list and u<sub>k</sub> is its tail. An empty list: ().
A stack is a list in which elements can only be added to or removed from the head. A pointer locates the head of the stack. Let S = (u<sub>0</sub>, u<sub>1</sub>, ..., u<sub>k-1</sub>, u<sub>k</sub>) be a stack. push(S, v) denotes adding v onto the stack by incrementing the pointer by one, giving (v, u<sub>0</sub>, ... u<sub>k</sub>). pop(S, u<sub>0</sub>) denotes the stack (u<sub>1</sub>, ... u<sub>k</sub>) that results from decreasing the pointer by one (removing u<sub>0</sub> from the head).

 A queue is a list in which elements can be added to the tail (appended) or removed (popped) from the head. Two pointers locate the head and the tail. Consider the queue
 Q = (u\_0, u\_1, ..., u\_{k-1}, u\_k). The append operation append(Q, u\_{k+1}) results in the queue (u\_0, ..., u\_k, u\_{k+1}) and the pop operation pop(Q, u\_0) results in the queue (u\_1, ..., u\_k).

### Stack and queue once more

Definition

List is called queue, if it enables efficient

- access head of a list,
- deletion of the list head (pop) and
- adding an element behind the current tail.

List is called stack, if it enables efficient

- access head of a list,
- deletion of the list head (pop) and
- adding and element before the current head (push).

### Stack and queue

Queue and stack are schematically depicted below. The arrows represent efficient (easily implementable) operations.



# Searching the adjacency graph $\mathcal{G}(A) = (\mathcal{V}, \mathcal{E})$

- The sequence in which the vertices are visited can be used to reorder the graph and hence permute the matrix.
- Given a start vertex, a graph search (graph traversal) performs exploration of the vertices and edges of  $\mathcal{G}(A)$
- It generates sets of visited vertices and explored edges.
- $V_v$ : the set of visited vertices,  $V_n$ : the set of vertices not yet visited.
- The search: selects an unexplored edge in  $\mathcal{E}$  with vertices in  $\mathcal{V}_v$ . If the other vertex belongs to  $\mathcal{V}_n$  then this vertex is moved into  $\mathcal{V}_v$  and the edge is flagged as explored.
- The explored edge may be directed or undirected; in an undirected graph, the edge (v, w) formally corresponds to the pair of edges (w → v) and (v → w).

- Starting from a chosen start vertex s, a breadth-first search (BFS) explores all the vertices adjacent to s.
- Then all the vertices whose shortest path from s is of length 2, and then length 3, and so on
- A queue is used in its implementation.
- The search terminates when there are no unexplored edges (v, w) with v ∈ V<sub>v</sub> and w ∈ V<sub>n</sub> that are reachable from s.



Figure: An illustration of a BFS of a connected undirected graph, with the labels indicating the order in which the vertices are visited.

## DFS

- A depth-first search (DFS) visits child vertices before visiting sibling vertices
- Starting from a chosen vertex s, the set of vertices that are visited are those vertices u for which a directed path from s to u exists in G.
- Different results depending on *s* and how ties are broken.
- Like the BFS, all vertices in  $\mathcal{R}each(s)$  are visited.
- Traversed edges form a DFS spanning tree. Visiting all the edges of a graph results in a DFS forest that consists of exactly one DFS spanning tree for each connected component.



Figure: An illustration of a DFS of a connected directed graph. The labels indicate the order in which the vertices are visited. The edges of the DFS spanning tree are in bold.

## DFS

- Number of ways to construct the output vertex order for a DFS.
- Given a start vertex s, in a preorder list, the vertices are returned in the order in which they are added into V<sub>v</sub>
- In a postorder list, the vertices are in the order in which they are last visited during the DFS algorithm For the example in Figure, the vertices are added into  $V_v$  in the order 1, 2, 3, 4, 5, 6, 7 and this is the preorder list.
- The sequence in which the DFS visits the vertices is 1, 2, 3, 2, 4, 2, 1, 5, 6, 5, 1, 7, 1. In this sequence, vertex 3 is the first vertex to appear for the last time so the postordering starts with vertex 3. The next vertex to appear for the last time is vertex 4, followed by vertex 2, and so on, resulting in the postorder list 3, 4, 2, 6, 5, 7, 1.

# Algorithm (Find preorder and postorder lists using a DFS)

1:  $\mathcal{V}_{v} = \emptyset$ , preorder = () and postorder = ()

2: for all  $v \in \mathcal{V}$  do

*3:* if  $v \notin \mathcal{V}_v$  then

- **4**: <math>push(preorder, v)
- 5:  $\mathcal{V}_v = \mathcal{V}_v \cup \{v\}$
- 6: dfs\_step(v)
- 7: end if
- 8: end for
- 9: recursive function (dfs\_step(v))
- 10: for all  $(v 
  ightarrow w) \in \mathcal{E}$  do
- 11: **if**  $w \notin \mathcal{V}_v$  then
- $12: \qquad push(preorder, w)$
- 13:  $\mathcal{V}_v = \mathcal{V}_v \cup \{w\}$
- 14: dfs\_step(w)
- 15: end if
- 16: end for
- $17: \quad push(postorder, v)$
- 18: end recursive function

Add v onto the preorder stack
 Add v to the set of visited vertices

Add w onto the preorder stack
 Add w to the set of visited vertices
 recursive search

 $\triangleright$  Add v onto the postorder stack

99/609

# Outline

**Factorizations** 

### Introduction to factorizations

- Entries of factors outside  $S{A}$ : filled entries. Together the filled entries are called the fill-in.
- Adding the fill-in to  $\mathcal{G}(A)$ : the filled graph.
- Numerical cancellations in LU factorizations rarely happen. Also difficult to predict.
- non-cancellation assumption: assuming that the result of adding, subtracting or multiplying two nonzeros is nonzero. It follows that if A = LU and E<sub>L</sub> denotes the set of (directed) edges of the digraph G(L) then for i > j

```
a_{ij} \neq 0 implies (i \rightarrow j) \in \mathcal{E}_L.
```

## Observation

The sparsity structures of the LU factors of A satisfy

 $\mathcal{S}\{A\} \subseteq \mathcal{S}\{L+U\}.$ 

### Introduction to factorizations

- The traditional way of describing Gaussian elimination is based on the systematic annihilation of the entries in the lower triangular part of *A* column-by-column.
- Assuming *A* is factorizable, this can be written formally as sequential multiplications by column elimination matrices
- Getting the elimination sequence:

$$A = A^{(1)}, A^{(2)}, \dots, A^{(n)}$$
(13)

of partially eliminated matrices as follows:

$$A^{(1)} \to A^{(2)} = C_1 A^{(1)} \to A^{(3)} = C_2 C_1 A^{(1)} \to \dots \to A^{(n)} = C_{n-1} \dots C_1 A^{(1)}$$

The unit lower triangular matrices C<sub>i</sub> (1 ≤ i ≤ n − 1) are the column elimination matrices.

# **Factorizations**

Elementwise, assuming  $a_{11} = a_{11}^{(1)} \neq 0$ , the first step  $C_1 A^{(1)} = A^{(2)}$  is

$$\begin{pmatrix} 1 & & & \\ -a_{21}^{(1)}/a_{11}^{(1)} & 1 & & \\ -a_{31}^{(1)}/a_{11}^{(1)} & 1 & & \\ \vdots & & & 1 & \\ -a_{n1}^{(1)}/a_{11}^{(1)} & & & 1 \end{pmatrix} \begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ a_{21}^{(1)} & a_{22}^{(1)} & \dots & a_{2n}^{(1)} \\ a_{31}^{(1)} & a_{32}^{(1)} & \dots & a_{3n}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1}^{(1)} & a_{n2}^{(1)} & \dots & a_{nn}^{(1)} \end{pmatrix} = \begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ 0 & a_{32}^{(2)} & \dots & a_{3n}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{n2}^{(2)} & \dots & a_{nn}^{(2)} \end{pmatrix}$$

Provided  $a_{22}^{(2)} \neq 0$ , the second step  $C_2 A^{(2)} = A^{(3)}$  is

$$\begin{pmatrix} 1 & & & \\ & 1 & & \\ & -a_{32}^{(2)}/a_{22}^{(2)} & 1 & \\ & \vdots & & 1 & \\ & -a_{n2}^{(2)}/a_{22}^{(2)} & & & 1 \end{pmatrix} \begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ 0 & a_{32}^{(2)} & \dots & a_{3n}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{n2}^{(2)} & \dots & a_{nn}^{(2)} \end{pmatrix} = \begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ 0 & 0 & a_{33}^{(3)} & \dots & a_{3n}^{(3)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & a_{n2}^{(2)} & \dots & a_{nn}^{(2)} \end{pmatrix} = \begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & a_{1n}^{(1)} \\ 0 & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ 0 & 0 & a_{33}^{(3)} & \dots & a_{3n}^{(3)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & a_{n3}^{(3)} & \dots & a_{nn}^{(3)} \end{pmatrix}$$

## Introduction to factorizations

- The *k*-th partially eliminated matrix is  $A^{(k)}$ .
- The active entries in  $A^{(k)}$ :  $a_{ij}^{(k)}$ ,  $1 \le k \le i, j \le n$ . The submatrix of  $A^{(k)}$  with the active entries: active submatrix.
- G(A<sup>(k)</sup>) is the k-th elimination graph and is denoted by G<sup>k</sup>. If S{A} is nonsymmetric then G<sup>k</sup> is a digraph.
- The inverse of each  $C_k$  is the unit lower triangular matrix obtained by changing the sign of all the off-diagonal entries.
- The product of unit lower triangular matrices (beware the order) is a unit lower triangular matrix: provided a<sup>(k)</sup><sub>kk</sub> ≠ 0 (1 ≤ k < n)</li>

$$A = A^{(1)} = C_1^{-1} C_2^{-1} \dots C_{n-1}^{-1} A^{(n)} = LU,$$

• The subdiagonal entries of *L* are the negative of the subdiagonal entries of the matrix  $C_1 + C_2 + \ldots + C_{n-1}$ .

## Introduction to factorizations: SPD

• If A is a symmetric positive definite (SPD) matrix then setting  $U = DL^T$ , the LU factorization can be written as

$$A = LDL^T,$$

- This is the square root-free Cholesky factorization.
- Alternatively, it can be expressed as the Cholesky factorization

$$A = (LD^{1/2})(LD^{1/2})^T,$$

where the lower triangular matrix  $LD^{1/2}$  has positive diagonal entries.

# **Factorizations**

## Introduction to factorizations: generic scheme

LU factorization can be stated in the generic form.  $l_{ik}$ : multipliers,  $a_{kk}^{(k)}$ : pivots. A is factorizable  $\Rightarrow a_{kk}^{(k)} \neq 0$  for all k.

Algorithm (Generic LU factorization)

Input: Factorizable matrix A. Output: LU factorization A = LU.



#### Introduction to factorizations: generic scheme

- Three nested loops: six ways of assigning the indices to the loops.
- The performance differs based on sparsity, computer architecture.
- In exact arithmetic: the same L and U.
- To identify the variants, names that derive from the order in which the indices are assigned to the loops can be used.
  - ► *kij* and *kji*: submatrix LU factorizations,
  - *jik* and *jki*: column factorizations.
  - ► The remaining ones: row factorizations (column LU factorization applied to *A*<sup>T</sup>.)

## Submatrix LU

Each outermost step of the submatrix LU variants computes one row of U and one column of L. The first step (k = 1) can be described using a combination of matrix and vector notation as follows:

$$C_1 A = \begin{pmatrix} 1 \\ -v/a_{11} & I \end{pmatrix} \begin{pmatrix} a_{11} & u^T \\ v & A_{2:n,2:n} \end{pmatrix} = \begin{pmatrix} a_{11} & u^T \\ A_{2:n,2:n} - vu^T/a_{11} \end{pmatrix},$$

where

$$v = (a_{21}, \dots, a_{n1})^T$$
,  $(l_{21}, \dots, l_{n1})^T = v/a_{11}$ ,  $u^T = (a_{12}, \dots, a_{1n})$ .

The  $(n-1) \times (n-1)$  active submatrix

$$S = A_{2:n,2:n} - vu^T / a_{11}$$

is the Schur complement of A with respect to  $a_{11}$ . If A is factorizable then so too is S and the process can be repeated.
## Submatrix LU

- The operations performed at each step *k* correspond to a sequence of rank-one updates.
- After k-1 steps  $(1 < k \le n)$ , the  $(n-k+1) \times (n-k+1)$  Schur complement of A with respect to its  $(k-1) \times (k-1)$  principal leading submatrix is the active submatrix of the partially eliminated matrix  $A^{(k)}$  given by

$$S^{(k)} = \begin{pmatrix} a_{kk} & \dots & a_{kn} \\ \vdots & \ddots & \vdots \\ a_{nk} & \dots & a_{nn} \end{pmatrix} - \sum_{j=1}^{k-1} \begin{pmatrix} l_{kj} \\ \vdots \\ l_{nj} \end{pmatrix} (u_{jk} & \dots & u_{jn})$$
$$= \begin{pmatrix} a_{kk}^{(k)} & \dots & a_{kn}^{(k)} \\ \vdots & \ddots & \vdots \\ a_{nk}^{(k)} & \dots & a_{nn}^{(k)} \end{pmatrix} = A_{k:n,k:n}^{(k)}.$$
(14)

 If A is SPD then the Cholesky and LDLT factorizations are termed right-looking (fan-out) factorizations.

$$A = \begin{pmatrix} 4 & -1 & -5 & -11 & 6 \\ -1 & 4 & -2 & -1 & 3 \\ -1 & 0 & 4 & -1 & 5 \\ 6 & 7 & 8 & 10 & 2 \\ 8 & -1 & -1 & 4 & 7 \end{pmatrix}$$

$$A = \begin{pmatrix} 4 & -1 & -5 & -11 & 6 \\ -0.25 & 3.75 & -3.25 & -3.75 & 4.5 \\ -0.25 & -0.25 & 2.75 & -3.75 & 6.5 \\ 1.5 & 8.5 & 15.5 & 26.5 & -7 \\ 2 & 1 & 9 & 26 & -5 \end{pmatrix}$$

$$A = \begin{pmatrix} 4 & -1 & -5 & -11 & 6 \\ -0.25 & 3.75 & -3.25 & -3.75 & 4.5 \\ -0.25 & -0.07 & 2.53 & -4 & 6.8 \\ 1.5 & 2.27 & 22.87 & 35 & -17.2 \\ 2 & 0.27 & 9.87 & 27 & -6.2 \end{pmatrix}$$

$$A = \begin{pmatrix} 4 & -1 & -5 & -11 & 6 \\ -0.25 & 3.75 & -3.25 & -3.75 & 4.5 \\ -0.25 & -0.07 & 2.53 & -4 & 6.8 \\ 1.5 & 2.27 & 9.03 & 71.1 & -78.6 \\ 2 & 0.27 & 3.89 & 42.6 & -32.7 \end{pmatrix}$$

$$A = \begin{pmatrix} 4 & -1 & -5 & -11 & 6 \\ -0.25 & 3.75 & -3.25 & -3.75 & 4.5 \\ -0.25 & -0.07 & 2.53 & -4 & 6.8 \\ 1.5 & 2.27 & 9.03 & 71.1 & -78.6 \\ 2 & 0.27 & 3.89 & 0.6 & 14.4 \end{pmatrix}$$

## Submatrix LU: kij algorithm

# Algorithm

kij lu decomposition (row oriented submatrix dense algorithm)  $l = I_n$   $u = O_n$ for k=1:n-1 for i=k+1:n  $l_{ik} = a_{ik}/a_{kk}$ for j=k+1:n  $a_{ij} = a_{ij} - l_{ik} * a_{kj}$ end  $u_{kk:n} = a_{kk:n}$ end  $u_{nn} = a_{nn}$ 

## Submatrix LU: kji algorithm

# Algorithm

*kji* lu decomposition (column oriented submatrix dense algorithm)  $l = I_n, u = O_n$ for k=1:n-1for s=k+1:n  $l_{sk} = a_{s,k}/a_{k,k}$ end for j=k+1:nfor i=k+1:n  $a_{ij} = a_{ij} - l_{ik} * a_{kj}$ end  $u_{kk:n} = a_{kk:n}$ end  $u_{nn} = a_{nn}$ 

## Submatrix LU: depiction



# Column LU

- In the column LU factorization, the outermost index is *j*.
- *l*<sub>11</sub> = 1, the off-diagonal entries in column 1 of *L* are obtained by dividing the corresponding entries in column 1 of *A* by *u*<sub>11</sub> = *a*<sub>11</sub>.
- Assume j 1 columns  $(1 < j \le n)$  of L and U have been computed. Then

$$\begin{pmatrix} L_{1:j-1,1:j-1} \\ L_{j:n,1:j-1} \end{pmatrix} U_{1:j-1,1:j-1} = \begin{pmatrix} A_{1:j-1,1:j-1} \\ A_{j:n,1:j-1} \end{pmatrix}$$

Column j of U and then column j of L are computed as

$$U_{1:j-1,j} = L_{1:j-1,1:j-1}^{-1} A_{1:j-1,j}, \quad u_{jj} = a_{jj} - L_{j,1:j-1} U_{1:j-1,j},$$

 $l_{jj} = 1, L_{j+1:n,j} = (A_{j+1:n,j} - L_{j+1:n,1:j-1}U_{1:j-1,j})/u_{jj}.$ 

• The strictly upper triangular part of U:j is determined from

$$L_{1:j-1,1:j-1}U_{1:j-1,j} = A_{1:j-1,j},$$

• The strictly lower triangular part of column *j* of *L* computed as a linear combination of column *A*<sub>*j*+1:*n*,*j*</sub> of *A* and previously computed columns of *L*.

## Column LU: A symmetric

### Observation

 $d_{jj}$  ( $1 \le j \le n$ ) of the LDLT factorization of the symmetric A is

$$d_{jj} = u_{jj} = a_{jj} - \sum_{k=1}^{j-1} d_{kk} l_{jk}^2.$$

The L factor is the same as is computed by the column LU factorization and

$$d_{jj}L_{j+1:n,j} = A_{j+1:n,j} - \sum_{k=1}^{j-1} L_{j+1:n,k} d_{kk} l_{jk}.$$

The U factor is equal to  $DL^T$ . Computing L and D in this way is called the left-looking (fan-in) factorization.

Algorithm (Column LU factorization with interchanging rows (partial pivoting))

**Input:** Nonsingular nonsymmetric matrix A. **Output:** LU factorization PA = LU, where P is a row permutation matrix.

- 1: Interchange rows of *A* so that  $|a_{11}| = \max\{|a_{i1}| | 1 \le i \le n\}$
- 2:  $l_{11} = 1$ ,  $u_{11} = a_{11}$ ,  $L_{2:n,1} = A_{2:n,1}/a_{11}$
- 3: for j = 2 : n do

4: Solve 
$$L_{1:j-1,1:j-1}U_{1:j-1,j} = A_{1:j-1,j}$$

5: 
$$z_{1:n-j+1} = A_{j:n,j} - L_{j:n,1:j-1}U_{1:j-1,j}$$

6: Apply row interchanges to z, A and L so that

$$|z_1| = \max\{|z_i| \mid 1 \le i \le n - j + 1\}.$$

7: 
$$l_{jj} = 1, u_{jj} = z_1$$
 and  $L_{j+1:n,j} = z_{2:n-j+1}/z_1$   
8: end for

$$A = \begin{pmatrix} 4 & -1 & -5 & -11 & 6 \\ -1 & 4 & -2 & -1 & 3 \\ -1 & 0 & 4 & -1 & 5 \\ 6 & 7 & 8 & 10 & 2 \\ 8 & -1 & -1 & 4 & 7 \end{pmatrix}$$

$$A = \begin{pmatrix} 4 & -1 & -5 & -11 & 6 \\ -0.25 & 4 & -2 & -1 & 3 \\ -0.25 & 0 & 4 & -1 & 5 \\ 1.5 & 7 & 8 & 10 & 2 \\ 2 & -1 & -1 & 4 & 7 \end{pmatrix}$$

$$A = \begin{pmatrix} 4 & -1 & -5 & -11 & 6 \\ -0.25 & 3.75 & -2 & -1 & 3 \\ -0.25 & -0.07 & 4 & -1 & 5 \\ 1.5 & 2.27 & 8 & 10 & 2 \\ 2 & 0.27 & -1 & 4 & 7 \end{pmatrix}$$

• Two completely different phases of the column construction.

- applying  $L_{i-1}^{-1}$  (as a forward substitution)
- computing subdiagonal part of L as a linear combination of a column of A with previously computed columns of L
- They are separated in the Cholesky factorization
- Row algorithm is column algorithm for  $A^T$

## Column algorithm: jki

## Algorithm

```
jki lu decomposition (delayed column dense algorithm)

l = I_{n_1} u = O_{n_1} u_{11} = a_{11}

for j=2:n

for s=j:n

l_{sj-1} = a_{sj-1}/a_{j-1j-1}

end

for k=1:j-1

for i=k+1:n

a_{ij} = a_{ij} - l_{ik} * a_{kj}

end

u_{1:jj} = a_{1:jj}
```

## Column algorithm: jik

## Algorithm

```
jik lu decomposition (dot product - based column dense algorithm) l=I_n,\,u_{11}=a_{11} for j=2:n
                 for s=j:n
                 l_{sj-1} = a_{sj-1}/a_{j-1j-1} end
                 for i=2:j
for k=1:i-1
                a_{ij} = a_{ij} - l_{ik} * a_{kj}
                 for i=j+1:n
for k=1:j-1
                a_{ij} = a_{ij} - l_{ik} * a_{kj}
end
           \mathop{\textit{end}}\limits^{u_{1:jj}} = a_{1:jj}
```

# Column algorithm: depiction



# Row algorithm: ikj

## Algorithm

## Row algorithm: ijk

## Algorithm

## Row algorithm: depiction



#### Outside the generic scheme

- An alternative is factorization by bordering.
- Set all diagonal entries of *L* to 1 and assume the first k-1 rows of *L* and first k-1 columns of *U* ( $1 < k \le n$ ) have been computed (that is,  $L_{1:k-1,1:k-1}$  and  $U_{1:k-1,1:k-1}$ ). At step k,  $A_{1:k,1:k}$  satisfies

$$\begin{pmatrix} A_{1:k-1,1:k-1} & A_{1:k-1,k} \\ A_{k,1:k-1} & a_{kk} \end{pmatrix} = \begin{pmatrix} L_{1:k-1,1:k-1} & 0 \\ L_{k,1:k-1} & 1 \end{pmatrix} \begin{pmatrix} U_{1:k-1,1:k-1} & U_{1:k-1,k} \\ 0 & u_{kk} \end{pmatrix}$$

• That is, the lower triangular part of row *k* of *L* and the upper triangular part of column *k* of *U* are obtained by solving

$$L_{k,1:k-1}U_{1:k-1,1:k-1} = A_{k,1:k-1},$$
  

$$L_{1:k-1,1:k-1}U_{1:k-1,k} = A_{1:k-1,k}.$$

The diagonal entry  $u_{kk}$  is then given by

$$u_{kk} = a_{kk} - L_{k,1:k-1}U_{1:k-1,k}$$
 (with  $u_{11} = a_{11}$ ).

### $\boldsymbol{A}$ not factorizable

- What if A is not factorizable?
- Then there exists a row permutation matrix *P* such that *PA* is factorizable.
- Consider the simple  $2 \times 2$  matrix A and its LU factorization

$$A = \begin{pmatrix} \delta & 1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 \\ \delta^{-1} & 1 \end{pmatrix} \begin{pmatrix} \delta & 1 \\ 1 - \delta^{-1} \end{pmatrix}.$$

If  $\delta = 0$  this factorization does not exist and if  $\delta$  is very small then the entries in the factors involving  $\delta^{-1}$  are very large.

• Interchanging the rows of A we have

$$PA = \begin{pmatrix} 1 & 1 \\ \delta & 1 \end{pmatrix} = \begin{pmatrix} 1 & \\ \delta & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ & 1 - \delta \end{pmatrix},$$

which is valid for all  $\delta \neq 1$ .

#### Lemma

Consider one step of the submatrix factorization of an SPD A. Schur complement of A with respect to (positive)  $a_{1,1}$  is positive definite.

Proof.

For 
$$(\alpha \quad z^T)^T$$
 we have  $x^T A x =$   

$$(\alpha \quad z^T) \begin{pmatrix} a_{1,1} & a_{1,2:n} \\ a_{2:n,1} & A_{2:n,2:n} \end{pmatrix} \begin{pmatrix} \alpha \\ z \end{pmatrix} =$$
(15)  

$$\alpha^2 a_{1,1} + \alpha a_{1,2:n} z + \alpha z^T a_{2:n,1} + z^T A_{2:n,2:n} z =$$
(16)  

$$(\alpha + a_{1,1}^{-1} a_{1,2:n} z)^T a_{1,1} (\alpha + a_{1,1}^{-1} a_{1,2:n} z) + z^T (A_{2:n,2:n} - a_{2:n,1} a_{1,1}^{-1} a_{1,2:n}) z$$
(17)

Choosing  $z \neq 0$  and setting  $\alpha = -a_{1,1}^{-1}a_{1,2:n}z$  we get

$$x^{T}Ax = z^{T}Sz$$
 where  $S = A_{2:n,2:n} - a_{2:n,1}a_{1,1}^{-1}a_{1,2:n}$ .

### Schemes for Cholesky

- Left-looking schemes (second part of the column LU)
- Right-looking schemes (just the submatrix scheme)
- The row scheme corresponds to the first part of the Cholesky algorithm.





## Submatrix (right-looking) Cholesky

$$\begin{split} A &= \begin{pmatrix} a_{1,1} & a_{1,2} & a_{1,3;n} \\ a_{2,1} & a_{2,2} & a_{2,3;n} \\ a_{3:n,1} & a_{3:n,2} & A_{3:n,3;n} \end{pmatrix} \\ &= \begin{pmatrix} \sqrt{a_{1,1}} & 0 \\ \frac{a_{2,1}}{\sqrt{a_{1,1}}} & \sqrt{a_{2,2}^{(1)}} \\ \frac{a_{3:n,1}}{\sqrt{a_{1,1}}} & \frac{a_{3:n,2}^{(1)}}{\sqrt{a_{2,2}^{(1)}}} & I_{n-2} \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & A_{3:n,3:n}^{(2)} - \frac{a_{3:n,1}a_{1,3:n}}{a_{1,1}} - \frac{a_{3:n,2}^{(1)}a_{2,3:n}^{(1)}}{a_{2,2}^{(1)}} \end{pmatrix} \\ &= \begin{pmatrix} \sqrt{a_{1,1}} & \frac{a_{2,1}}{\sqrt{a_{1,1}}} & \frac{a_{1,3:n}}{\sqrt{a_{1,1}}} \\ 0 & \sqrt{a_{2,2}^{(1)}} & \frac{a_{2,3:n}}{\sqrt{a_{2,2}^{(1)}}} \\ & & I_{n-2} \end{pmatrix} \\ &= \begin{pmatrix} l_{1,1} & 0 & 0 \\ l_{2,1} & l_{2,2} & 0 \\ l_{3:n,1} & l_{3:n,2} & I_{n-2} \end{pmatrix} \begin{pmatrix} l_{1,1} & l_{2,1} & l_{1,3:n} \\ 0 & 0 & I_{n-2} \end{pmatrix} \end{split}$$

### Column (left-looking) Cholesky

# Algorithm

Column Cholesky factorization:  $A \rightarrow$  square-root factor  $L = (l_{ij})$ 1. for j = 1 : n do 2. Compute an auxiliary vector  $t_{j:n}$ 

$$\begin{pmatrix} t_j \\ \vdots \\ t_n \end{pmatrix} = \begin{pmatrix} a_{jj} \\ \vdots \\ a_{nj} \end{pmatrix} - \sum_{\{k \mid l_{jk} \neq 0\}} l_{jk} \begin{pmatrix} l_{jk} \\ \vdots \\ l_{nk} \end{pmatrix}$$
(18)

3. Get a column of L by scaling  $t_{j:n}$ 

$$\begin{pmatrix} l_{jj} \\ \vdots \\ l_{nj} \end{pmatrix} = \frac{1}{\sqrt{t_j}} \begin{pmatrix} t_j \\ \vdots \\ t_n \end{pmatrix}$$
(19)

#### 4. end *j*

# Row (up-looking) Cholesky

## Algorithm

Row Cholesky factorization: :  $A \rightarrow$  square-root factor  $L = (l_{ij})$ . 1. for i = 1 : n do

2. Solve the triangular system

$$L_{1:i-1,1:i-1} \begin{pmatrix} l_{i1} \\ \vdots \\ l_{i,i-1} \end{pmatrix} = \begin{pmatrix} a_{i1} \\ \vdots \\ a_{i,i-1} \end{pmatrix}$$
(20)

3. Compute the diagonal entry  $l_{ii} = \sqrt{\left(a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2\right)}$ 4. end *i* 

### Elimination versus decomposition (factorization)

- Householder (end of 1950's, beginning of 1960's): expressing Gaussian elimination as a decomposition
- Various reformulations of the same decomposition: different properties in
  - sparse implementations
  - vector processing
  - parallel implementations
- Six algorithms of the generic scheme for Cholesky as for LU, but there are also others (bordering, Dongarra-Eisenstat)

## Remind the fill-in

• Arrow matrix - original matrices: example showing how bad the fill-in problem can be



## Remind the fill-in

### • Arrow matrix - structure after elimination

## Remind the fill-in

### Arrow matrix - structure after elimination



Fill-in description and ways to avoid it must capture it dynamically!

### Simple fill-in results

- Assume S{A} is symmetric and consider the elimination graph G<sup>k</sup> at step k.
- Its vertices are the n k + 1 uneliminated vertices. Its edge set contains the edges in G(A) connecting these vertices and additional edges corresponding to filled entries produced during the first k 1 elimination steps.
- The sequence of graphs  $\mathcal{G}^1 \equiv \mathcal{G}(A), \mathcal{G}^2, \dots$  is generated recursively using Parter's rule:

To obtain the elimination graph  $\mathcal{G}^{k+1}$  from  $\mathcal{G}^k$ , delete vertex k and add all possible edges between vertices that are adjacent to vertex k in  $\mathcal{G}^k$ .

- Denoting  $\mathcal{G}^k = (\mathcal{V}^k, \mathcal{E}^k), \ \mathcal{G}^{k+1} = (\mathcal{V}^{k+1}, \mathcal{E}^{k+1})$ , the Parter's rule:  $\mathcal{V}^{k+1} = \mathcal{V}^k \setminus \{k\}, \ \mathcal{E}^{k+1} = \mathcal{E}^k \cup \{(i,j) \mid i, j \in adj_{\mathcal{G}^k}\{k\}\} \setminus \{i \mid i \in adj_{\mathcal{G}^k}\{k\}\}.$
- If S{A} is symmetric then Parter's rule says that the adjacency set of vertex k becomes a clique when k is eliminated: Gaussian elimination systematically generates cliques.
- As the elimination process progresses, cliques grow or more than one clique joins to form larger cliques: clique amalgamation.
- A clique with m vertices has m(m-1)/2 edges. It can be represented by storing a list of its vertices, without any reference to edges.
- If *S*{*A*} is nonsymmetric then the elimination graphs are digraphs and Parter's rule generalizes as follows:
- To obtain the elimination graph  $\mathcal{G}^{k+1}$  from  $\mathcal{G}^k$ , delete vertex k and add all edges  $(i \xrightarrow{\mathcal{G}^{k+1}} j)$  such that  $(i \xrightarrow{\mathcal{G}^k} k)$  and  $(k \xrightarrow{\mathcal{G}^k} j)$ .

• Simple example (symmetric):



**Figure:** Illustration of Parter's rule. The original undirected graph  $\mathcal{G} = \mathcal{G}^1$  and the elimination graph  $\mathcal{G}^2$  that results from eliminating vertex 1 are shown on the left and right, respectively. The red dashed lines denote fill edges. The vertices  $\{2, 3, 4\}$  become a clique.
## **Factorizations**

• Simple example (nonsymmetric):



**Figure:** Illustration of Parter's rule for a nonsymmetric  $S\{A\}$ . The original digraph  $\mathcal{G} = \mathcal{G}^1$  and the directed elimination graph  $\mathcal{G}^2$  that results from eliminating vertex 1 are shown on the left and right, respectively. The red dashed lines denote fill edges.

#### Parter's rule for factors

• The repeated application of Parter's rule specifies all the edges in  $\mathcal{G}(L+L^T)$ :

(i, j) is an edge of  $\mathcal{G}(L + L^T)$  if and only if (i, j) is an edge of  $\mathcal{G}(A)$  or (i, k) and (k, j) are edges of  $\mathcal{G}(L + L^T)$  for some k < i, j.

• This generalizes to a nonsymmetric matrix *A* and its LU factorization:

 $(i \rightarrow j)$  is an edge of the digraph  $\mathcal{G}(L + U)$  if and only if  $(i \rightarrow j)$  is an edge of the digraph  $\mathcal{G}(A)$  or  $(i \rightarrow k)$  and  $(k \rightarrow j)$  are edges of  $\mathcal{G}(L + U)$  for some k < i, j.

## **Factorizations**

 Parter's rule is a local rule that uses the dependency on nonzeros obtained in previous steps of the factorization. The following result fully characterizes the nonzero entries in the factors using only paths in *G*(*A*).

#### Theorem

(a) Let S{A} be symmetric and A = LL<sup>T</sup>. Then (L+L<sup>T</sup>)<sub>ij</sub> ≠ 0 if and only if there is a fill-path i ∉ G(A) min j.
(b) Let S{A} be nonsymmetric and A = LU. Then (L+U)<sub>ij</sub> ≠ 0 if and only if there is a fill-path i ∉ G(A) min j.
The fill-paths may not be unique.

## **Factorizations**

Symmetric  $S{A}$ : a filled entry in position (8,6) of *L* because of a fill-path 8  $\xleftarrow{\mathcal{G}(A)}{min}$  6: 8  $\longleftrightarrow$  2  $\longleftrightarrow$  5  $\longleftrightarrow$  1  $\longleftrightarrow$  6.



## Corollary

(Edges of  $\mathcal{G}^k$  in terms of reachable sets) Assume  $\mathcal{S}\{A\}$  is symmetric. Let  $\mathcal{V}^k$  be the set of k-1 vertices of  $\mathcal{G}(A)$  that have already been eliminated and let v be a vertex in the elimination graph  $\mathcal{G}^k$ . Then the set of vertices adjacent to v in  $\mathcal{G}^k$  is the set  $\mathcal{R}each(v, \mathcal{V}^k)$  of vertices reachable from v through  $\mathcal{V}^k$  in  $\mathcal{G}(A)$ .

#### Proof.

The proof is by induction on k. The result holds trivially for k = 1 because the  $\mathcal{R}each(v, \mathcal{V}^1) = adj_{\mathcal{G}(A)}\{v\}$ . Assume the result holds for  $\mathcal{G}^1, \ldots, \mathcal{G}^k$  with  $k \ge 1$  and let v be a vertex in the graph  $\mathcal{G}^{k+1}$  that is obtained after eliminating  $v_k$  from  $\mathcal{G}^k$ . If v is not adjacent to  $v_k$  in  $\mathcal{G}^k$  then  $\mathcal{R}each(v, \mathcal{V}^{k+1}) = \mathcal{R}each(v, \mathcal{V}^k)$ . Otherwise, if v is adjacent to  $v_k$  in  $\mathcal{G}^k$  then  $adj_{\mathcal{G}^{k+1}}\{v\} = \mathcal{R}each(v, \mathcal{V}^k) \cup \mathcal{R}each(v_k, \mathcal{V}^k)$ . In both cases Parter's rule implies that the new adjacency set is exactly equal to the vertices that are reachable from v through  $\mathcal{V}^{k+1}$ , that is,  $\mathcal{R}each(v, \mathcal{V}^{k+1})$ .

## **Factorizations**

Figure depicts a graph  $\mathcal{G}(A)$ . The adjacency sets of the vertices in  $\mathcal{G}^4$  that result from eliminating vertices  $\mathcal{V}^4 = \{1, 2, 3\}$  are  $adj_{\mathcal{G}^4}\{4\} = \mathcal{R}each(4, \mathcal{V}^4) = \{5\}, adj_{\mathcal{G}^4}\{5\} = \mathcal{R}each(5, \mathcal{V}^4) = \{4, 6, 7\}, adj_{\mathcal{G}^4}\{6\} = \mathcal{R}each(6, \mathcal{V}^4) = \{5, 7\}, adj_{\mathcal{G}^4}\{7\} = \mathcal{R}each(7, \mathcal{V}^4) = \{5, 6, 8\}, adj_{\mathcal{G}^4}\{8\} = \mathcal{R}each(8, \mathcal{V}^4) = \{7\}.$ 



**Figure:** An example to illustrate reachable sets in  $\mathcal{G}(A)$ . The grey vertices 1, 2, and 3 are eliminated in the first three elimination steps ( $\mathcal{V}^4 = \{1, 2, 3\}$ ).

### So far, only implicit results

- Neither the local characterization of filled entries using Parter's rule nor Theorem on paths provide a direct answer as to whether a certain edge belongs to  $\mathcal{G}(L + L^T)$  (or  $\mathcal{G}(L + U)$ ); without performing the eliminations
- Results presented so far do not tell us whether a given entry of a factor of *A* is nonzero.
- Such questions are addressed by deeper theoretical and algorithmic results that are presented later.

## **Factorizations**

#### From factorization to the solution

• Once an LU factorization has been computed, the solution x of the linear system Ax = b is computed by solving the lower triangular system

$$Ly = b, (21)$$

followed by the upper triangular system

$$Ux = y. \tag{22}$$

- Triangular solves with a dense right-hand side vectorstraightforward.
- First forward substitution: the component  $y_1$  is determined from the first equation. Substitution into the second equation to obtain  $y_2$ , and so on.
- Once y is available, the solution can be obtained by back substitution in which the last equation is used to obtain  $x_n$ , which is then substituted into equation n - 1 to obtain  $x_{n-1}$ , and so on.

# **Factorizations**

Algorithm (Forward substitution: simple lower triangular solve Ly = b with b dense)

**Input:** Lower triangular matrix L with nonzero diagonal entries and dense right-hand side b.

**Output:** The dense solution vector y.

- 1: Initialise y = b
- 2: for j = 1 : n do
- 3:  $y_j = y_j / l_{jj}$
- 4: **for** i = j + 1 : n **do**
- 5: if  $l_{ij} \neq 0$  then
- $b: y_i = y_i l_{ij}y_j$
- 7: end if
- 8: end for
- 9: end for

### Sparse right-hand side

- When *b* is sparse, the solution *y* is also sparse. In particular, if in Algorithm  $y_k = 0$ , then the outer loop with j = k can be skipped.
- Furthermore, if b<sub>1</sub> = b<sub>2</sub> = ... = b<sub>k</sub> = 0 and b<sub>k+1</sub> ≠ 0, then y<sub>1</sub> = y<sub>2</sub> = ... = y<sub>k</sub> = 0. Scanning y to check for zeros adds O(n) to the complexity.
- But if the set of indices  $\mathcal{J} = \{j \mid y_j \neq 0\}$  is known beforehand then we can use the following Algorithm,
- A possible way to determine  $\mathcal{J}$  is discussed later.

# **Factorizations**

Algorithm (Forward substitution: lower triangular solve Ly = b with *b* sparse)

**Input:** Lower triangular matrix L with nonzero diagonal entries, sparse vector b and the set  $\mathcal{J}$ .

**Output:** The sparse solution vector y.

1: Initialise y = b2: for  $j \in \mathcal{J}$  do  $\triangleright$  Take indices from  $\mathcal{J}$  in increasing order 3:  $y_i = y_i / l_{ii}$ for i = j + 1 : n do 4: if  $l_{ij} \neq 0$  then 5:  $y_i = y_i - l_{ij}y_j$ 6: end if 7: end for 8: 9: end for

# Outline

Reducibility and blocks

Enhancements due to permuting into a block form

- Permuting to block form is closely connected to matrix reducibility.
- *A* is said to be reducible if there is a permutation matrix *P* such that

$$PAP^{T} = \begin{pmatrix} A_{p_{1},p_{1}} & A_{p_{1},p_{2}} \\ 0 & A_{p_{2},p_{2}} \end{pmatrix},$$

where  $A_{p_1,p_1}$  and  $A_{p_2,p_2}$  are non trivial square matrices (that is, they are of order at least 1).

- If *A* is not reducible, it is irreducible. Matrices of order 1 are irreducible.
- If  $S{A}$  is symmetric then  $A_{p_1,p_2} = 0$  and  $PAP^T$  is block diagonal.
- A one-sided permutation can transform an irreducible matrix *A* into a reducible matrix *AQ*.

$$A = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & \\ 1 & - \end{pmatrix}, \quad Q = \begin{pmatrix} 1 & 1 & \\ 1 & - \\ 1 & - \end{pmatrix}, \quad AQ = \begin{pmatrix} 1 & 1 & 1 \\ - & 1 & 1 \\ - & 1 \end{pmatrix}.$$

### Hall properties

- A matrix A is said to be a Hall matrix (or to have the Hall property) if every set of k columns has nonzeros in at least k rows  $(1 \le k \le n)$ .
- A is a strong Hall matrix (or to have the strong Hall property) if every set of k columns  $(1 \le k < n)$  has nonzeros in at least k + 1 rows.
- If A is square then A has the strong Hall property if and only if the directed graph G(A) is strongly connected.

#### Theorem

Given a nonsingular nonsymmetric matrix A there exists a permutation matrix P such that

$$PAP^{T} = \begin{pmatrix} A_{1,1} & A_{1,2} & \cdots & A_{1,nb} \\ 0 & A_{2,2} & \cdots & A_{2,nb} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & A_{nb,nb} \end{pmatrix},$$
(23)

where the square matrices  $A_{ib,ib}$  on the diagonal are irreducible. The set  $\{A_{ib,ib} | 1 \le ib \le nb\}$  is uniquely determined (but the blocks may appear on the diagonal in a different order). The order of the rows and columns within each  $A_{ib,ib}$  may not be unique.

• This upper block triangular form (BTF) is also known as the Frobenius normal form. It is said to be non trivial if nb > 1, and this is the case if A does not have the strong Hall property.

#### Permutation to BTF

An example of a matrix that can be symmetrically permuted to block triangular form with nb = 2 is given in Figure 16.

	1	2	3	4	5	6		6	3	5	4	1	2
1	(*	*		*	*		6	(*	*				* \
2	[	*			*		3	*	*	*			*
3		*	*		*	*	5			*		*	
4		*		*			4				*		*
5	*				*		1			*	*	*	*
6		*	*			*/	2			*			*/

**Figure:** The sparsity patterns of *A* (left) and the upper block triangular form  $PAP^T$  with two blocks  $A_{ib,ib}$ , i = 1, 2, of order 2 and 4 (right).

## Using BTF

## Algorithm (Solve a sparse linear system in upper BTF)

*Input:* Upper block triangular matrix (23) and a conformally partitioned right-hand side vector c. **Output:** The conformally partitioned solution vector y.

1: for ib = 1 : nb do 2: Compute  $P_{ib}A_{ib,ib} = L_{ib}U_{ib}$ 3: end for 4: Solve  $L_{nb}U_{nb} y_{nb} = P_{nb}c_{nb}$ 5: for ib = nb - 1 : 1 do 6: for jb = ib + 1 : nb do 7:  $c_{ib} = c_{ib} - A_{ib,jb}y_{jb}$ 8: end for 9: Solve  $L_{ib}U_{ib} y_{ib} = P_{ib}c_{ib}$ 10: end for

> Perform forward and back substitutions

Sparse matrix-vector operation

> Perform substitutions

#### Transversal

- The transversal of a matrix A is the set of its nonzero diagonal elements.
- A has a full or maximum transversal if all its diagonal entries are nonzero. There exist permutation matrices P and Q such that PAQ has a full transversal if and only if A has the Hall property.
- Moreover, if *A* is nonsingular then it can be nonsymmetrically permuted to have a full transversal. However, the converse is clearly not true (for example, a matrix with all its entries equal to one has a full transversal but it is singular).
- If A has a full transversal then there exists a permutation P<sub>s</sub> such that P<sub>s</sub>AP<sup>T</sup><sub>s</sub> has the BTF form. In the other words, once A has a full transversal, a symmetric permutation is sufficient to obtain the BTF form.

### Permutation to BTF

- Finding P<sub>s</sub> is identical to finding the strongly connected components (SCCs) of the digraph G(A) = (V, E).
- To find the SCCs, V is partitioned into non-empty subsets V<sub>i</sub> with each vertex belonging to exactly one subset. Each vertex i in the quotient graph corresponds to a subset V<sub>i</sub> and there is an edge in the quotient graph with endpoints i and j if E contains at least one edge with one endpoint in V<sub>i</sub> and the other in V<sub>j</sub>.
- The condensation (or component graph) of a digraph is a quotient graph in which the SCCs form the subsets of the partition, that is, each SCC is contracted to a single vertex. This reduction provides a simplified view of the connectivity between components.

#### Permutation to BTF

An example of SCCs. Here five SCCs:  $\{p, q, r\}$ ,  $\{s, t, u\}$ ,  $\{v\}$ ,  $\{w\}$ , and  $\{x\}$ .



**Figure:** An illustration of the strong components of a digraph. On the left, the five SCCs are denoted using different colours and on the right is the condensation DAG  $\mathcal{G}_C$  formed by the  $_{162/609}$ 

### Relationship between SCCs and DAGs

#### Theorem

The condensation  $\mathcal{G}_C$  of a digraph is a DAG (directed acyclic graph).

- Any DAG can be topologically ordered.
- Consequently,  $\mathcal{G}_C = (\mathcal{V}_C, \mathcal{E}_C)$  can be topologically ordered and if  $\mathcal{V}_i$  and  $\mathcal{V}_j$  are contracted to  $s_i$  and  $s_j$  and  $(s_i \longrightarrow s_j) \in \mathcal{E}_C$  then  $s_i < s_j$ .
- It follows that to permute *A* to block triangular form it is sufficient to find the SCCs of *G*(*A*).
- That is, topologically ordering the vertices of the condensation  $\mathcal{G}_C$  induced by the SCCs is the quotient graph that implies the block triangular form.

## Finding SCCs

- Tarjan's algorithm: the key idea here is that vertices of a SCC form a subtree in the DFS spanning tree of the graph.
- The algorithm performs depth-first searches, keeping track of two properties for each vertex v: when v was first encountered (held in *invorder*(v)) and the lowest numbered vertex that is reachable from v (called the low-link value and held in *lowlink*(v)).
- It pushes vertices onto a stack as it goes and outputs a SCC when it finds a vertex for which *invorder*(v) and *lowlink*(v) are the same. The value *lowlink*(v) is computed during the DFS from v, as this finds the vertices that are reachable from v.
- The outermost loop of the algorithm visits each vertex that has not yet been visited, ensuring vertices that are not reachable from the starting vertex are eventually visited.
- The complexity bound is  $O(|\mathcal{V}| + |\mathcal{E}|)$ .

## Algorithm (Tarjan's algorithm to find the SCCs)

```
1: \mathcal{V}_v = \emptyset, S = (), index = 0,
                                                                                                         Each vertex is initially unvisited
2: for each v \in V do

3: if v \notin V_v then
4: so
5: end
6: end for
           scomp step(v)
       end if
7: recursive function (scomp_step(v))
8:
       \mathcal{V}_v = \mathcal{V}_v \cup \{v\}
                                                                                                      Add v to the set of visited vertices
9:
        index = index + 1
                                                                                           \triangleright Set the index for v to smallest unused index
10:
          invorder(v) = index, lowlink(v) = index
11:
          push(S, v)
                                                                                                                       \triangleright Put v on the stack
12:
          Set v = head(S)
                                                                                                               \triangleright v is the current head of S
13:
          for each (v \rightarrow w) \in \mathcal{E} do
                                                                                                           ▷ Look in the adjacency list of v
14:
             if w \notin \mathcal{V}_v then
                                                                                                    > w not yet been visited; recurse on it
15:
                 scomp step(w)
16:
17:
                 lowlink(v) = min(lowlink(v), lowlink(w))
             else if w \in S then
                                                                                            \triangleright w is in the stack and hence in current SCC
18:
                 lowlink(v) = min(lowlink(v), invorder(w))
19:
20:
21:
22:
23:
             end if
          end for
          if low link(v) = invorder(v) then
             pop all vertices down to v from S to obtain a new SCC
          end if
24: end recursive function
```

Finding blocks: indistinguishability

• Finding sets of columns of *A* frequently have identical sparsity patterns. For example, for *A* from a finite element discretisation.

### Definition

We say that two vertices u and v of an undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  are indistinguishable if they have the same neighbours, that is,  $adj_{\mathcal{G}}\{u\} \cup \{u\} = adj_{\mathcal{G}}\{v\} \cup \{v\}$ . A set of mutually indistinguishable vertices is called an indistinguishable vertex set. If  $\mathcal{U} \subseteq \mathcal{V}$  is an indistinguishable vertex set then  $\mathcal{U}$  is maximal if  $\mathcal{U} \cup \{w\}$  is not indistinguishable for any  $w \in \mathcal{V} \setminus \mathcal{U}$ .

- $\bullet\,$  Indistinguishability is an equivalence relation on  ${\cal V}$
- Maximal indistinguishable vertex sets represent its classes → a partitioning of V into nsup ≥ 1 non-empty disjoint subsets

$$\mathcal{V} = \mathcal{V}_1 \cup \mathcal{V}_2 \cup \ldots \cup \mathcal{V}_{nsup}. \tag{24}$$

 An indistinguishable vertex set can be represented by a single vertex, called a supervariable.

### Finding blocks: indistinguishability

- S{A} is symmetric, G(A) = (V, E). Let V be partitioned into indistinguishable vertex sets and reorder the vertices such that those belonging to each subset V<sub>1</sub>,..., V<sub>nsup</sub> are numbered consecutively, with those in V<sub>i</sub> preceding those in V<sub>i+1</sub> (1 ≤ i < nsup).</li>
- If *P* is the permutation matrix corresponding to this ordering then the permuted matrix *PAP*<sup>T</sup> has a block structure in which the blocks are dense.

### Finding supervariables

- Initially place all the vertices in a single vertex set (that is, into a single supervariable).
- Split into two supervariables by taking the first vertex j = 1 and moving those vertices that are in the adjacency set of j into a new vertex set (a new supervariable).
- Each vertex *j* is considered in turn and each vertex set V<sub>sv</sub> that contains a vertex in adj<sub>G</sub>{j} ∪ j is split into two by moving the vertices in adj<sub>G</sub>{j} ∪ j that belong to V<sub>sv</sub> into a new vertex set.
- As a result of the splitting and moving of vertices a vertex set can become empty, in which case it is discarded.
- Once the supervariables have been determined, the permuted matrix *PAP<sup>T</sup>* can be condensed to a matrix of order equal to *nsup*.

### Finding supervariables

### Algorithm (Find the supervariables of an undirected graph)

Input: Graph G of a symmetrically structured matrix. Output: Partitioning of V into indistinguishable vertex sets.

1:  $\mathcal{V}_1 = \{1, 2, \dots, n\}$ *2:* for j = 1 : n do 3: for  $i \in adj_{\mathcal{G}}\{j\} \cup j$  do 4: Find sv such that  $i \in \mathcal{V}_{sv}$ 5: if this is the first occurrence of sv for the current index j then 6: Establish a new set  $\mathcal{V}_{nsv}$  and move *i* from  $\mathcal{V}_{sv}$  to  $\mathcal{V}_{nsv}$ 7: else 8: Move i from  $\mathcal{V}_{sv}$  to  $\mathcal{V}_{nsv}$ 9: end if 10: Discard  $\mathcal{V}_{sv}$  if it is empty 11: end for 12: end for

#### Finding supervariables: illustration

- Initially, 1, 2, 3, 4, 5 are put into a single vertex set  $V_1$ .
- Vertices i = 1, 2 and 5 belong to adj<sub>G</sub>{1} ∪ {1}; they are moved from V₁ into a new vertex set.

### • There is no further splitting of the vertex sets for j = 2.

- adj<sub>G</sub>{3} ∪ {3} = {3,4,5}. i = 3 and 4 are moved from V<sub>1</sub> into a new vertex set. V<sub>1</sub> is now empty and can be discarded. Vertex i = 5 is moved from the vertex set that holds vertices 1 and 2 into a new vertex set. For j = 4 and 5 no additional splitting is performed.
- Thus, three supervariables are found, namely  $\{1, 2\}$ ,  $\{3, 4\}$  and  $\{5\}$ .

### Approach by Ashcraft

- Number graph nodes/vertices (use numbers as their labels)
- Ocmpute vertex checksums:

$$chksum = \sum_{\{u,v\}\in E} w$$

- Sort vertices by their checksums: in O(|E|+|V|log(|V|)) time
- Different checksum means different block
- Signature First tie-breaking rule: if chksum(u) = chksum(v): compare |adj(u)| and |adj(v)|
- Second tie-breaking rule: compare adjacency sets of *u* and *v* (the most time consuming)

### Approximate indistinguishability: Saad

- Using symbolic dot products between the rows of the matrix.
- Here we assume that  $S{A}$  is symmetric but modifications exist.
- Rewrite A as row vectors

$$A = \left(a_1^T, \dots, a_n^T\right)^T,$$

and consider  $\mathcal{G}(A) = (\mathcal{V}, \mathcal{E})$ .

- A partition V = V<sub>1</sub> ∪ ... ∪ V<sub>nb</sub> is constructed using row products a<sup>T</sup><sub>i</sub>a<sub>k</sub> between different rows of A that express the level of orthogonality between the rows;
- if  $a_i^T a_k$  is large then *i* and *k* are assigned to different vertex sets.

### Approximate indistinguishability

- Algorithm treats all entries of A as unity (symbolically)
- The symbolic row products can be considered as a generalization of the angles between rows expressed by their cosines, hence
- On output, if  $adjmap(i_1) = adjmap(i_2)$  then vertices  $i_1$  and  $i_2$  belong to the same vertex set.
- Symmetry of S{A} simplifies the computation of the symbolic row products because for row *i* only *k* > *i* is considered, that is, only the symbolic row products that correspond to one triangle part of A<sup>T</sup>A are checked.
- The procedure outlined in Algorithm 6.4 is controlled by a threshold parameter  $\tau \in (0, 1]$ .

### Algorithm

```
Find approximately indistinguishable vertex sets
1: nb = 0, adjmap(1:n) = 0, cosine(1:n) = 0
2: for i = 1 : n do
3:
       if adjmap(i) = 0 then
4:
           nb = nb + 1
                                                                                      Start a new set
5:
           adjmap(i) = ib
6:
           for (i, j) \in \mathcal{E} do
                                                                  \triangleright Corresponds to an entry in A_{i,1:n}
7:
               for (k, j) \in \mathcal{E} with k > i do
                                                      \triangleright Both rows i and k have an entry in column j
8:
                  if adjmap(k) = 0 then \triangleright k has not been yet added to some partitioning set
9:
                      cosine(k) = cosine(k) + 1
                                                                        Increase partial dot product
10:
                   end if
11:
                end for
12:
                for k with cosine(k) \neq 0 do
13:
                   if cosine(k)^2 \ge \tau^2 * nz_i * nz_k then
                                                                      > Test similarity of row patterns
14:
                       adjmap(k) = nb
15:
                   end if
16:
                   cosine(k) = 0
17:
                end for
18:
            end for
19:
        end if
                                                                                                   174/609
```



**Figure:** A example to illustrate Algorithm. The original matrix is given (left) together with the permuted matrix with indistinguishable vertex sets  $\mathcal{V} = \{1,3\} \cup \{2,6\} \cup \{4\} \cup \{5\}$  obtained using  $\tau = 1$  (centre) and the permuted matrix with approximately indistinguishable vertex sets  $\mathcal{V} = \{1,3,5\} \cup \{2,6\} \cup \{4\}$  obtained using  $\tau = 0.5$  (right). The threshold  $\tau = 0.5$  results in putting row 5 into the same set as row 1, making the vertex sets only approximately indistinguishable. The permuted matrix has an approximate block form.

# Outline

Symbolic Cholesky factorization

### Symbolic Cholesky

- Symbolic Cholesky: only  $S{A}$  used to determine the nonzero structure of the Cholesky factor *L* without computing the numerical values of the nonzeros.
- Implicitly assumed that all the diagonal entries of A are included in S{A} (even if they are zero this cannot happen if A is SPD).
- A fundamental difference between dense and sparse Cholesky factorizations is that, in the latter, each column of *L* depends on only a subset of the previous columns.
- Operations to update a computed column should be also sparse.
- We will see that the symbolic Cholesky can be clearly described using a tree structure.

## Symbolic Cholesky

### Column replication

Let us start from the patterns of subsequent Schur complements.

$$S^{(k)} = A_{k:n,k:n} - \sum_{j=1}^{k-1} \begin{pmatrix} l_{kj} \\ \vdots \\ l_{nj} \end{pmatrix} (l_{kj} \dots l_{nj}).$$
(25)

• Consider column j of L  $(1 \le j \le k - 1)$  and let  $l_{ij} \ne 0$  for some i > j.

Observation

For any  $i > j \ge 1$  such that  $l_{ij} \ne 0$ 

$$\mathcal{S}\{L_{i:n,j}\} \subseteq \mathcal{S}\{L_{i:n,i}\}.$$
(26)

This is called the column replication principle because the pattern of column j of L (rows i to n) is replicated in the pattern of column i of L.
## Column replication

 Denote the row index of the first subdiagonal nonzero entry in column j of L by parent(j), that is,

 $parent(j) = \min\{i \mid i > j \text{ and } l_{ij} \neq 0\}.$ (27)

- If there is no such entry, set parent(j) = 0. The row index parent(parent(j)) is denoted by  $parent^2(j)$ , and so on.
- Applying column replication recursively implies the sparsity pattern of column *j* of *L* is replicated in that of column *parent*(*j*), which in turn is replicated in the pattern of column *parent*<sup>2</sup>(*j*), and so on.

## **Column replication**

Example: Consider j = 1. Because the first subdiagonal entry in column 1 is in row 3, parent(1) = 3, parent(3) = parent<sup>2</sup>(1) = 5.



**Figure:** An illustration of column replication. On the left are the entries in *L* before step 1 of a Cholesky factorization (that is, the entries in the lower triangular part of *A*); in the centre we show the replication of the nonzeros from column 1 in the pattern of column parent(1) = 3 (red entries *f*); on the right, we show the subsequent replication in column  $parent^2(1) = 5$ .

#### Fill-in as replication of column structures



## Existence of nonzeros in columns

The following result shows that, provided *A* is irreducible, the mapping parent(j) has nonzero values given by (27) for all j < n.

#### Theorem

If *A* is SPD and irreducible then in each column *j* ( $1 \le j < n$ ) of its Cholesky factor *L* there exists an entry  $l_{ij} \ne 0$  with i > j.

## Proof.

From Parter's rule, each step of the Cholesky factorization corresponds to adding new edges into the graph of the corresponding Schur complement. If *A* is irreducible then the graphs corresponding to the Schur complements are connected. Consequently, for any vertex j ( $1 \le j < n$ ) in any of these graphs there is at least one vertex i with i > j to which j is connected. This corresponds to the nonzero entry in column j of *L*.

## Row replication

With the convention  $parent^1(j) = parent(j)$ , the next theorem shows that if entry  $l_{ij}$  of L is nonzero then  $parent^t(j) = i$  for some  $t \ge 1$  and there is row replication in the sequence  $i = parent^1(i) = parent^2(i)$ 

 $j, parent^1(j), parent^2(j), \ldots, parent^t(j).$ 

## Theorem

Let A be SPD and let L be its Cholesky factor. If  $l_{ij} \neq 0$  for some  $j < i \le n$  then there exists  $t \ge 1$  such that  $parent^t(j) = i$  and  $l_{ik} \ne 0$  for k = j,  $parent^1(j)$ ,  $parent^2(j)$ , ...,  $parent^t(j)$ .

## Proof.

If  $i = parent^{1}(j)$ , the result is immediate. Otherwise, there exists an index k, j < k < i of a subdiagonal entry in column j of L such that  $k = parent^{1}(j)$ . Column replication implies  $l_{ik} \neq 0$ . Applying an inductive argument to  $l_{ik}$ , the result follows after a finite number of steps.  $\Box$ 

## Row replication: example



 Replication of column structures → replication of nonzeros in a row. Necessary and sufficient condition (NSC) for a fill-in entry

- If there is sequence of nonzeros in a row of *L*, it is natural to ask where the sequence begins:
- If there is no k ≥ 1 such that a<sub>ik</sub> ≠ 0, no replication of nonzeros can start in row i.

## Theorem

Let *A* be SPD and let *L* be its Cholesky factor. If  $a_{ij} = 0$  for some  $1 \le j < i \le n$  then there is a filled entry  $l_{ij} \ne 0$  if and only if there exists k < j and  $t \ge 1$  such that  $a_{ik} \ne 0$  and  $parent^t(k) = j$ .

## Elimination tree

- The discussion of column replication is significantly simplified using elimination trees.
- The elimination tree (or etree)  $\mathcal{T}(A)$  (or simply  $\mathcal{T}$ ) of a SPD matrix has vertices  $1, 2, \ldots, n$  and an edge between each pair (j, parent(j)), where parent(j) has been defined above;
- j is a root vertex of the tree if parent(j) = 0.
- The edges of  $\mathcal{T}$  are considered to be directed from a child to its parent, that is,

$$\mathcal{E}(\mathcal{T}) = \{ (j \longrightarrow i) \mid i = parent(j) \}.$$

If  $\mathcal{T}$  has a single component then the root vertex is n.

 Despite the terminology, the elimination tree need not be connected and in general is a forest. For simplicity, in our discussions, we assume T has a single component and we say that T is described by the vector *parent*.

#### Elimination tree



**Figure:** An illustration of a sparse matrix *A* with a symmetric sparsity pattern and its elimination tree T(A). The root vertex is 8. The filled entries in  $S\{L + L^T\}$  are denoted by *f*.

## Elimination tree: terminology

- Concepts such as child, leaf, ancestor and descendant vertices for directed rooted trees can be applied to  $\mathcal{T}$ .
- Additionally, anc<sub>T</sub>{j} and desc<sub>T</sub>{j} are defined to be the sets of ancestors and descendants of vertex j in T.
- *T*(*j*): the subtree of *T* induced by *j* and *desc*<sub>*T*</sub>{*j*); *j* is the root vertex of *T*(*j*).
- The size  $|\mathcal{T}(j)|$  is the number of vertices in the subtree.
- A pruned subtree of T(j) is the connected subgraph induced by j and a subset of desc<sub>T</sub>{j). That is, for any vertex i in a pruned subtree of T(j), all the ancestors of i belong to T(j).
- A pruned subtree of  $\mathcal{T}$  shares the mapping *parent* with  $\mathcal{T}$ .

## Elimination tree: simple properties

The following observation is straightforward.

Observation If  $i \in anc_{\mathcal{T}}\{j\}$  for some  $j \neq i$  then i > j.

• The connection between the mapping *parent* and the sets of ancestors and descendants is emphasized by the next observation.

## Observation

If *i* and *j* are vertices of the elimination tree  $\mathcal{T}$  with  $j < i \leq n$  then  $i \in anc_{\mathcal{T}}\{j\}$  if and only if  $j \in desc_{\mathcal{T}}\{i\}$  if and only if  $parent^t(j) = i$  for some  $t \geq 1$ .

## Elimination tree: simple properties

- Replications can be expressed using rooted trees.
- For example, instead of stating that there exists  $t \ge 1$  such that  $parent^t(j) = i$ , we can write that  $i \in anc_T\{j\} \setminus \{j\}$ .
- Rewriting the necessary and sufficient condition in Theorem above as the following corollary we get a clear characterization of the sparsity patterns of the rows of *L*.

## Corollary

Consider the elimination tree  $\mathcal{T}$  and the Cholesky factor L of A. If i and j are vertices of  $\mathcal{T}$  with  $j < i \leq n$  and  $a_{ij} = 0$  then  $l_{ij} \neq 0$  if and only if there exists k < j such that  $j \in anc_{\mathcal{T}}(k)$  and  $a_{ik} \neq 0$ .

#### Elimination tree: row subtrees

- The subtree of T with vertices that correspond to the nonzeros of row i of L is called the i-th row subtree and is denoted by Tr(i).
- Formally, row subtree is a pruned subtree of  $\mathcal{T}$  induced by the union of the vertex set

$$\{i\} \cup \{k \mid a_{ik} \neq 0 \text{ and } k < i\}$$

with all vertices on the directed paths in T from k to i, that is, with all their ancestors from  $T_r(i)$ .

- The root vertex is *i* and the leaf vertices are a subset of the column indices in the *i*-th row of the lower triangular part of *A*.
- The row subtrees are connected subgraphs of *T*, even if *T* is not connected. If *T* can be found without determining the pattern of *L*, then *T<sub>r</sub>(i)* can be used to derive the sparsity pattern of row *i* of *L*, without having to store each entry explicitly.

Elimination tree: row subtrees



**Figure:** The row subtree  $T_r(5)$  of the elimination tree T from above (left). Vertex 3 has been pruned because  $a_{35} = 0$ . The row subtree  $T_r(8)$  (right) differs from T = T(A) because vertex 1 has been pruned ( $a_{18} = 0$ ).

#### Elimination tree: row subtrees: another example





## Row subtrees



## Theorem

If *i* and *j* are vertices in the elimination tree T with  $j < i \le n$  then  $i \in anc_T\{j\}$  if and only if there exists a path

$$j \xleftarrow{\mathcal{G}(A)}{\{1,\dots,i\}} i$$

## Proof.

Assume  $i \in anc_{\mathcal{T}}\{j\}$ . Then there is a path  $j \stackrel{\mathcal{T}}{\Longrightarrow} i$  of length  $l \ge 1$ . Each edge of this path belongs to  $\mathcal{G}(L)$  and corresponds either to an edge in  $\mathcal{G}(A)$  or to a fill-path in  $\mathcal{G}(A)$ . Connecting these paths together gives (28). Conversely, if the path (28) exists, induction on its length can be used to prove the result. If the path is of length 1 then the result holds because i and j are connected in  $\mathcal{G}(A)$  by an edge. Consequently, i is an ancestor of j. Now assume that the result is true for all paths of length less than l (l > 1), and consider a path of length l. Let m be the largest vertex on this path. If m < j

then (28) is a fill-path connecting i and j and, therefore,  $i \in anc_T\{j\}$ . Otherwise, for  $m \ge j$  the

assumption implies  $i \in anc_{\mathcal{T}}\{m\}$  and  $m \in anc_{\mathcal{T}}\{j\}$ , that is,  $i \in anc_{\mathcal{T}}\{j\}$ .

196/609

(28)

## Locality in characterization of ancestors

- Given a vertex *j* in *T*, the following corollary indicates how to find parent(j) (if it exists).
- If the set of ancestors of *j* is non empty then the lowest numbered one is its parent.

# Corollary

Vertex *i* is the parent of vertex *j* in T if and only if *i* is the lowest numbered vertex satisfying  $j < i \le n$  for which there is a path (28).

 The existence of (28) is equivalent to requiring *i* and *j* belong to the same component of the graph G(A<sub>1:i,1:i</sub>) corresponding to the *i* × *i* principal leading submatrix A<sub>1:i,1:i</sub> of A.

## Locality in characterization of ancestors

• Figure depicts  $\mathcal{G}(A)$  for the matrix A given in Figure above.

• Consider vertex 4. Its set of ancestors for which the paths mentioned above exist comprises vertices 5, 6 and 8. Vertex 7 is not an ancestor of 4 because there is no path from 7 to 4 in the graph  $\mathcal{G}(A_{1:7,1:7})$ . Among the ancestors of 4, vertex 5 fulfils the condition from Corollary above and is thus the parent of 4.



Figure: The graph  $\mathcal{G}(A)$  of the matrix from Figure 20 illustrating Theorem and Corollary above.

- Locality: T = T(A) can be constructed by stepwise extensions of the elimination trees of the principal leading submatrices of A.
- Assume we have  $\mathcal{T}(A_{1:i-1,1:i-1})$  and we want  $\mathcal{T}(A_{1:i,1:i})$ .
- Initialise *T*(*A*<sub>1:*i*,1:*i*</sub>) = *T*(*A*<sub>1:*i*-1,1:*i*-1</sub>). If there are no entries in row *i* of *A* to the left of the diagonal then there is nothing to do, only an isolated vertex *i* is added. Otherwise, *i* is the root of the row subtree *T<sub>r</sub>*(*i*) and an ancestor of some vertex in *T*.
- For each such vertex, say *j*, its ancestors *k* with *k* < *i* are vertices of *T*(*A*<sub>1:*i*-1,1:*i*-1</sub>). Consider the directed path in *T*(*A*<sub>1:*i*,1:*i*</sub>) *t* ≥ 1 such that parent<sup>t</sup>(*j*) = *i*, we can write that *i* ∈ anc<sub>*T*</sub>{*j*} \ {*j*}. are connected subgraphs of *T*.
- If  $parent^t(j) \neq i$  then the new root of the subtree of  $\mathcal{T}(A_{1:i,1:i})$  that contains j is added by setting  $i = parent^{t+1}(j)$ . Otherwise, i has already been added.

## Constructing elimination tree

## Algorithm (Construction of an elimination tree)

**Input:** A with a symmetric sparsity pattern and its undirected graph G. **Output:** Elimination tree T described by the vector *parent*.





Constructing elimination tree



Constructing elimination tree  $* \ / \ * \ * \ \rangle$ 



1





















#### Constructing elimination tree: comments

- The most expensive part of Algorithm to find the elimination tree is the while loop that searches for subtree roots. This search is based on tracing the directed path from *j* to its root *parent*<sup>t</sup>(*j*).
- Because this path is unique for a given *j*, shortcuts can be incorporated; this is called path compression.
- Having found a directed path from j to k, subsequent searches can be made more efficient by introducing a vector ancestor and setting ancestor(j) = k.
- The modified algorithm is outlined below. It maintains two structures side-by-side using the current values of *parent* and *ancestor*. The tree described by *ancestor* is termed the virtual tree.

# Algorithm (Construction of an elimination tree using path compression)

**Input:** A with a symmetric sparsity pattern and its undirected graph G. **Output:** Elimination tree T described by the vector *parent*.

```
1: for i = 1 : n do
2:
       parent(i) = 0, ancestor(i) = 0
3:
       for j \in adj_{\mathcal{G}}\{i\} and j < i do
4:
          jroot = j
5:
           while ancestor(jroot) = 0 and ancestor(jroot) = i do
6:
              l = ancestor(jroot)
7:
              ancestor(jroot) = i
                                                         To accelerate future searches
8:
              jroot = l
<u>9</u>:
          end while
10:
           if ancestor(jroot) = 0 then
11:
               ancestor(jroot) = i and parent(jroot) = i
12.
           end if
13:
        end for
14: end for
                                                                                     203/609
```

Constructing elimination tree with compression

• Figure below shows a matrix for which path compression makes building T significantly more efficient.



Figure: A sparse matrix for which computing the elimination tree using Algorithm 7.2 is muc more efficient than using Algorithm 7.1.

• For this example,  $\mathcal{T}$  is determined by the mapping parent(6) = 0; parent(i) = i + 1 for i = 1, ..., 5.

## Constructing elimination tree with compression

- The complexity of the first Algorithm is  $O(n^2)$ .
- For this example the complexity of the Algorithm with compression is *O*(*n*).
- Formally, the complexity of the new Algorithm is  $O(nz(A) \log_2(n))$ , where nz(A) is the number of nonzeros of A but the logarithmic factor is rarely reached.
- Additional modifications can reduce the theoretical complexity to O(nz(A) g(nz(A), n)), where g(nz(A), n) is a very slowly increasing function called the functional inverse of Ackermann's function.

# Constructing elimination tree with compression Independence of subtrees

- The following simple theorem states that there is no edge in  $\mathcal{G}(L + L^T)$  between vertices belonging to subtrees of  $\mathcal{T}$  with different vertex sets.
- Importance for parallel computations.

#### Theorem

Consider the elimination tree  $\mathcal{T}$  and the Cholesky factor L of A. Let  $\mathcal{T}(i)$  and  $\mathcal{T}(j)$  be two vertex-disjoint subtrees of  $\mathcal{T}$ . Then for all  $s \in \mathcal{T}(i)$  and  $t \in \mathcal{T}(j)$  the entry  $l_{st}$  of L is zero.

## Sparsity pattern of *L*

- The explicit structure of *L* is not always required, only the numbers of nonzeros in each row and column of *L* are needed.
- For example, to allocate the storage
- Let row<sub>L</sub>{i} denote the sparsity pattern of the off-diagonal part of row i of L, that is,

$$row_L\{i\} = \mathcal{S}\{L_{i,1:i-1}\} = \{j \mid j < i, \ l_{ij} \neq 0\}, \quad 1 \le i \le n.$$

• The number of entries in L is

$$nz(L) = \sum_{i=1}^{n} |row_L\{i\}| + n.$$

- $row_L\{i\}$  is given by the vertices of the row subtree  $\mathcal{T}_r(i)$ .
- The complexity of the algorithm is O(nz(L)).
# Algorithm (Computation of the row sparsity patterns of the Cholesky factor L)

**Input:** A with a symmetric sparsity pattern, its undirected graph  $\mathcal{G}$  and elimination tree  $\mathcal{T}$  described by the vector *parent*.

**Output:** Row sparsity patterns  $row_L\{i\}$  of the Cholesky factor L of A  $(1 \le i \le n)$ .

1: for i = 1 : n do  $\triangleright$  Loop over the rows of A  $row_L\{i\} = \emptyset$ 2: Initialisation 3: mark(i) = i4: for  $k \in adj_{\mathcal{G}}\{i\}$  and k < i do  $\triangleright$  Loop over below diagonal entries in column i of A 5: j = k6: while  $mark(j) \neq i$  do  $\triangleright$  Column j not yet encountered in row i 7: mark(j) = i $\triangleright$  Flag *j* as encountered in row *i* 8:  $row_L\{i\} = row_L\{i\} \cup \{j\}$  $\triangleright$  Add *j* to the sparsity pattern of row *i* Move up the elimination tree 9: j = parent(j)10: end while 11: end for 12: end for 208/609

### Cholesky factorization skeleton

- Efficiency can be improved by employing the skeleton graph *G*(*A*<sup>-</sup>) that is obtained from *G*(*A*) by removing every edge (*i*, *j*) for which *j* < *i* and *j* is not a leaf vertex of *T<sub>r</sub>*(*i*).
- $\mathcal{G}(A^-)$  is the smallest subgraph of  $\mathcal{G}(A)$  with the same filled graph as  $\mathcal{G}(A)$ . The corresponding matrix is the skeleton matrix.
- The complexity of constructing the elimination tree using the skeleton matrix and its graph G(A<sup>-</sup>) is O(nz(A<sup>-</sup>) g(nz(A<sup>-</sup>), n)), where nz(A<sup>-</sup>) is the number of entries in the skeleton matrix. Because nz(A<sup>-</sup>) is often significantly smaller than nz(A), an implementation that processes G(A<sup>-</sup>) rather than G(A) can be substantially faster.

Cholesky factorization skeleton



**Figure:** An illustration of the sparsity pattern of A and its graph  $\mathcal{G}(A)$  (left) and the sparsity pattern of the corresponding skeleton matrix  $A^-$  and graph  $\mathcal{G}(A^-)$  (right). The entries in A **200d**609

Row counts: more spohisticated approach Row counts: more sophisticated algorithm: the idea



- Needed: fast algorithm to determine junctions of branches in the elimination tree,
- and fast algorithm to find leaves of the elimination tree.
- Can be done by traversing the postordered elimination tree.
- The complexity can be then nearly linear in m.

## Sparse Cholesky factorization - components

It would be nice to know column structures of L as well



• Analogously to the row sparsity patterns, let  $col_L\{j\}$  denote the sparsity pattern of the off-diagonal part of column *j* of *L*, that is,

$$col_L\{j\} = \mathcal{S}(L_{j+1:n,j}) = \{i \mid i > j, \ l_{ij} \neq 0\}, \quad 1 \le j \le n.$$

• The column replication principle can be written as

 $col_L\{j\} \subseteq col_L\{parent(j)\} \cup parent(j).$ 

• Theorem describes  $col_L\{j\}$  using the vertices of the subtree  $\mathcal{T}(j)$ .

#### Theorem

The column sparsity pattern  $col_L\{j\}$  of the Cholesky factor L of the matrix A is equal to the adjacency set of vertices of the subtree  $\mathcal{T}(j)$  in  $\mathcal{G}(A)$ , that is,

$$col_L\{j\} = adj_{\mathcal{G}(A)}\{\mathcal{T}(j)\}.$$
(29)

#### Proof.

If  $i \in col_L\{j\}$  then  $j \in row_L\{i\}$  and Theorem on necessary and sufficient condition for the fill-in implies  $j \in anc_T\{k\}$  for some k such that  $a_{ik} \neq 0$ . That is,  $i \in adj_G\{T(j)\}$ . Conversely,  $i \in adj_G\{T(j)\}$  implies that in row i the entry in column j of L is nonzero. Thus,  $j \in row_L\{i\}$  and hence  $i \in col_L\{j\}$ .  $\Box$ 

Algorithm can be used to compute the column counts and the column sparsity patterns because when *j* is added to *row*<sub>L</sub>{*i*} at line 8, *i* can be added to *col*<sub>L</sub>{*j*}. This does not generally obtain the column sparsity patterns sequentially.

To derive an approach that does compute them sequentially consider

$$col_L\{j\} = \left(adj_{\mathcal{G}(A)}\{j\} \bigcup_{\{k \mid k \in \mathcal{T}(j)\}} col_L\{k\}\right) \setminus \{1, \dots, j\}.$$

Using the column replication, this can be significantly simplified

$$col_L\{j\} = \left(adj_{\mathcal{G}(A)}\{j\} \bigcup_{\{k \mid j=parent(k)\}} col_L\{k\}\right) \setminus \{1,\ldots,j\}.$$
 (30)

- The following algorithm constructs the sparsity pattern of each column *j* of *L* as the union of the sparsity pattern of column *j* of *A* (*adj*<sub>G(A)</sub>{*j*}) and the patterns of the children of *j* in *T*(*A*).
- Here *child*{*j*} denotes the set of children *j*.
- Because any child k of j satisfies k < j, the j-th outer step has the information needed to compute the sparsity pattern described by (30). Observe that T(A) does not need to be input.

### Algorithm (Determining sparsity patterns of columns of L)

**Input:** A with symmetric sparsity pattern and its undirected graph  $\mathcal{G}(A)$ . **Output:** Column sparsity patterns  $col_L\{j\}$  of the Cholesky factor L of A ( $1 \le j \le n$ ).

```
1: for j = 1 : n do
2: child\{j\} = \emptyset
                                                                                           Initialisation
3: end for
4: for j = 1 : n do
                                                                      \triangleright Loop over the columns of L
5:
        col_L\{j\} = adj_{\mathcal{G}(A)}\{j\} \setminus \{1, \dots, j-1\}
     for k \in child\{j\} do
                                                                  ▷ Unifying child structures in (30)
6:
7:
            col_L\{j\} = col_L\{j\} \cup col_L\{k\} \setminus \{j\}
8:
        end for
9:
        if col_L\{j\} \neq \emptyset then
10:
             l = \min\{i \mid i \in col_L\{j\}\}
11:
             child\{l\} = child\{l\} \cup \{j\}
                                                        Parent of j detected using Corollary ??
12:
         end if
13: end for
```

### **Topological orderings**

- The outer loop in Algorithm to find column sparsity patterns does not have to be performed in the strict order j = 1,...,n.
- What is necessary is that for each step *j*, the column sparsity pattern for each child of *j* has already been computed.
- An ordering of the vertices in a tree (and, more generally, in a DAG) is a topological ordering if, for all *i* and *j*, *j* ∈ desc<sub>T</sub>{*i*} implies *j* < *i*
- Observation above confirms that the ordering of vertices in the elimination tree  $\mathcal{T}$  is a topological ordering.
- A new topological ordering of  $\mathcal{T}$  defines a relabelling of its vertices corresponding to a symmetric permutation of *A*.

### **Topological orderings**



Figure: Two topological orderings of an elimination tree.

### **Topological orderings**

The sparsity patterns of the Cholesky factors of A and  $PAP^{T}$  can be different but the following result shows that the amount of fill-in is the same.

#### Theorem

Let  $S{A}$  be symmetric. If *P* is the permutation matrix corresponding to a topological ordering of the elimination tree T of *A* then the filled graphs of *A* and  $PAP^{T}$  are isomorphic.

### **Topological orderings**

- There are many topological orderings of  $\mathcal{T}$ .
- One class is obtained using the depth-first search.
- This algorithm searches all the components of  ${\mathcal T}$  starting at their root vertices.
- In this case, once vertex *i* has been visited, all the vertices of the subtree T(*i*) are visited immediately after *i* and *i* is labelled as the last vertex of T(*i*).

### Postordering

- A topological ordering of T is a postordering if the vertex set of any subtree T(i) (i = 1,...,n) is a contiguous sublist of 1,...,n.
- Unless additional rules on how vertices are selected are imposed, a postordering is generally not unique.
- One possible postordering is defined by the depth-first search, and this is apparently not unique.

### Postordering



Figure: An example to illustrate the non uniqueness of postorderings of an elimination tree.

### Leaf vertices of row subtrees

- Leaf vertices of row subtrees play a key role in graph algorithms related to sparse Cholesky factorizations.
- They can be used to find the skeleton matrix described above.
- See the next theorem

#### Theorem

Let the elimination tree  $\mathcal{T}$  of A be postordered. Let the column indices of the nonzeros in the strictly lower triangular part of row i of A be  $c_1, \ldots, c_s$  with  $s \ge 1$  and  $0 < c_1 < \ldots < c_s < i$ . Then  $c_t$  is a leaf vertex of the row subtree  $\mathcal{T}_r(i)$  if and only if

$$t = 1$$
 or  $(1 < t \leq s \text{ and } c_{t-1} \notin \mathcal{T}(c_t)).$ 

### Proof.

 $c_1$  is always a leaf vertex of  $\mathcal{T}_r(i)$ . If this is not the case then there exists a directed path from some vertex  $k, k \neq c_1$  to *i* via  $c_1$  such that  $k \in \mathcal{T}_r(i)$  and  $a_{ik} \neq 0$ . Postordering of  $\mathcal{T}$  implies  $k < c_1$ . This is a contradiction because  $c_1$  is the index of the first nonzero in row *i*.

Consider now t > 1. Assume that  $c_{t-1} \in \mathcal{T}(c_t)$  and that  $c_t$  is a leaf vertex of  $\mathcal{T}_r(i)$ . Row replication implies any  $k \in anc_{\mathcal{T}}\{c_{t-1}\}$  such that  $c_{t-1} \leq k < i$  satisfies  $l_{ik} \neq 0$ . Because  $\mathcal{T}$  is postordered,  $c_{t-1} \leq k \leq c_t$  and there is at least one  $k < c_t$  satisfying this inequality. It follows that  $k = c_{t-1}$ . Because k belongs to  $\mathcal{T}_r(i)$ ,  $c_t$  cannot be a leaf vertex of  $\mathcal{T}_r(i)$ , which is a contradiction.

Conversely, assume for t > 1 that  $c_{t-1} \notin \mathcal{T}(c_t)$  and  $c_t$  is not a leaf vertex of  $\mathcal{T}_r(i)$ . From the second part of the assumption and the fact that  $c_t \in \mathcal{T}_r(i)$  it follows that there is at least one leaf

### Leaf vertices of row subtrees

### Corollary

Under the assumptions of the previous Theorem,  $c_t$  is a leaf vertex of  $\mathcal{T}_r(i)$  if and only if

$$t = 1$$
 or  $(1 < t \le s \text{ and } c_{t-1} < c_t - |\mathcal{T}(c_t)| + 1).$ 

- Subtree sizes can be easily computed bottom up.
- Correctness of next Algorithm is guaranteed because *parent* defines a topological ordering of *T*.

Subtree sizes - needed to get the leaf vertices of row subtrees

### Algorithm (Find the sizes of subtrees T(i) of T)

Input: Elimination tree  $\mathcal{T}$  described by the vector parent. Output: Subtree sizes  $|\mathcal{T}(i)|$  ( $1 \le i \le n$ ).

1:  $|\mathcal{T}(1:n)| = 1$ 2: **for** i = 1: n - 1 **do** 3: k = parent(i)4:  $|\mathcal{T}(k)| = |\mathcal{T}(k)| + |\mathcal{T}(i)|$ 

Leaf vertices of row subtrees by columns

• Theorem relaxes the condition that the entries in the rows of *A* are sorted by increasing column indices. This allows the leaf vertices of the row subtrees to be determined by columns.

#### Theorem

Consider the elimination tree  $\mathcal{T}$  of A. Vertex j is a leaf vertex of some row subtree of  $\mathcal{T}$  if and only if there exists  $i \in adj_{\mathcal{G}(A)}\{j\}$ ,  $j < i \leq n$ , such that  $i \notin adj_{\mathcal{G}(A)}\{k\}$  for all  $k \in \mathcal{T}(j) \setminus \{j\}$ .

#### Proof.

Assume that for some  $i \in anc_{\mathcal{T}}\{j\}$  vertex j is a leaf vertex of  $\mathcal{T}_r(i)$ . That is,  $i \in adj_{\mathcal{G}(A)}\{j\}$ , i > j. Suppose there exists k in  $\mathcal{T}(j) \setminus \{j\}$  such that  $i \in adj_{\mathcal{G}(A)}\{k\}$ . Then all the ancestors of  $k, k \leq i$ , in particular j, belong to  $\mathcal{T}_r(i)$  and j cannot be a leaf vertex of  $\mathcal{T}_r(i)$ . This is a contradiction.

Conversely, assume that j is not a leaf vertex of any row subtree of T and that there exists

 $i \in adj_{\mathcal{G}(A)}\{j\}, j < i \le n$ , such that  $i \notin adj_{\mathcal{G}(A)}\{k\}$  for all  $k \in \mathcal{T}(j) \setminus \{j\}$ . Because j is not a leaf vertex of any such  $\mathcal{T}_r(i)$ , Theorem on necessary and sufficient fill-in conditions implies that there exists  $k \in \mathcal{T}(j) \setminus \{j\}$  such that  $a_{ik} \neq 0$  which gives a contradiction and completes the  $2^{28/609}$ 

### Leaf vertices of row subtrees

- To find leaf vertices of row subtrees of *T*, Algorithm uses a marking scheme based on Theorem above and exploits the postordering of *T*.
- The auxiliary vector *prev\_nonz* stores the column indices of the most recently encountered entries in the rows of the strictly lower triangular part of *A*.

### Leaf vertices of row subtrees

### Algorithm (Find leaf vertices of row subtrees of T)

**Input:** A with a symmetric sparsity pattern and a corresponding postordered elimination tree T.

Output: Logical vector *isleaf* with entries set to true for leaf vertices of row subtrees.

1: isleaf(1:n) = false, prev nonz(1:n) = 02: Compute  $|\mathcal{T}(1:n)|$ ▷ Use Algorithm 7.5 *3:* for j = 1 : n do  $\triangleright$  Loop over the columns of A 4: for i such that i > j and  $a_{ij} \neq 0$  do  $\triangleright$  Row index in strictly lower triangular part of A 5: k = prev nonz(i)Column index of most recently seen entry in row i 6: if  $k < j - |\mathcal{T}(j)| + 1$  then 7: ▷ j is a leaf vertex by Corollary ?? isleaf(j) = true8: end if 9: prev nonz(i) = j $\triangleright$  Flag *j* as the most recently seen entry in row *i* 10: end for 11: end for

### Blocks

- Blocks are absolutely crucial to compute efficiently on contemporary computers: we need as much data as possible for a unit of data transfer inside memory hierarchy.
- In BLAS terminology:

 $z = x + \alpha y \longrightarrow Z = X + \alpha Y (vector opes)$ 

in general: saxpy  $\longrightarrow$  dgemm

- But we have sparse matrices. It is not so straightforward to split their nonzeros into blocks.
- In fact, we need to reorder them in order to get blocks.
  - Application-based blocks in discretized systems.
  - Graph-based strategies which can be very fast.
  - ▶ But we need to optimize the block structure of *L*: supernodes.
  - Help: again our good friend, the elimination tree.

## Sparse Cholesky factorization - components





#### Supernodes and the assembly tree

- Because of column replication, the columns of *L* generally become denser as the Cholesky factorization proceeds.
- To exploit this, we require the concept of supernodes. The idea: group together columns with the same sparsity structure, so that they can be treated as a dense matrix.
- Let  $1 \le s, t \le n$  with  $s + t 1 \le n$ . A set of contiguously numbered columns of *L* with indices  $S = \{s, s + 1, \dots, s + t 1\}$  is a supernode of *L* if

$$col_L\{s\} \cup \{s\} = col_L\{s+t-1\} \cup \{s, \dots, s+t-1\},$$
 (31)

and S cannot be extended for s > 1 by adding s - 1 or for s + t - 1 < n by adding s + t.

- Because *S* cannot be extended it is a maximal subset of column indices. In graph terminology, a supernode is a maximal clique of contiguous vertices of  $\mathcal{G}(L + L^T)$ .
- A supernode may contain a single vertex.

### Supernodes





Figure: An example to illustrate supernodes in *L*. The first supernode comprises columns 1 and 2, the second columns 3 and 4, and the third columns 5 to 8.

### Supernodes

- The supernodal elimination or assembly tree is defined to be the reduction of the elimination tree that contains only supernodes.
- Each vertex of the elimination tree is associated with one elimination and a single integer (the index of its parent) is needed.
- Associated with each vertex of the assembly tree is an index list of the row indices of the nonzeros in the columns of the supernode. These implicitly define the sparsity pattern of *L*.

Demonstrating the difference between the elimination and assembly trees is given in Figure 28. Here the elimination tree is postordered and there are 5 supernodes: {1, 2}, 3, 4, 5, {6, 7, 8, 9}. For supernode 1 that comprises columns 1 and 2, the row index list is {1, 2, 8, 9}.



#### Convenient characterization of supernodes

#### Theorem

The set of columns of *L* with indices  $S = \{s, s + 1, ..., s + t - 1\}$  is a supernode of *L* if and only if it is a maximal set of contiguous columns such that s + i - 1 is a child of s + i for i = 1, ..., t - 1 and

$$|col_L\{s\}| = |col_L\{s+t-1\}| + t - 1.$$
 (32)

### Proof.

Let *S* be a supernode. For  $i, j \in S$  with i > j we have  $i \in col_L\{j\}$ . This implies that in the postordered elimination tree the vertex i = j + 1 is the parent of j for  $j = s, \ldots, s + t - 2$ . Moreover, from Observation 26, for any  $i, j \in S$  with i > j,  $i \in col_L\{j\}$  implies  $col_L\{j\} \setminus \{1, \ldots, i\} \subseteq col_L\{i\}$ . Therefore,

$$|col_L\{s+i\}| \ge |col_L\{s+i-1\}| - 1, \quad i = 1, \dots, t-1,$$
(33)

with equality if and only if

 $col_L\{s+i\}=col_L\{s+i-1\}\setminus\{s+i\},$ 

238/609

### Supernodes

- Enhance the efficiency of sparse factorizations and sparse triangular solves because they enable floating-point operations to be performed on dense submatrices rather than on individual nonzeros, thus improving memory hierarchy utilization and allowing the use of highly efficient dense linear algebra kernels (such as Level 3 BLAS kernels).
- Columns within a supernode are numbered consecutively but they can be numbered within the supernode in any order without changing the number of nonzeros in *L* (assuming the corresponding rows are permuted symmetrically). On some architectures, particularly those using GPUs, this freedom can be exploited to improve the factorization efficiency.
- Supernode amalgamation to achieve better efficiency.

### Fundamental supernodes

- In practice, fundamental supernodes are easier to work with in the numerical factorization.
- Let  $1 \le s, t \le n$  with  $s + t 1 \le n$ . A maximal set of contiguously numbered columns of L with indices  $S = \{s, s + 1, \dots, s + t - 1\}$  is a fundamental supernode if for any successive pair i - 1 and i in the list, i - 1 is the only child of i in  $\mathcal{T}$  and  $col_L\{i\} \cup \{i\} = col_L\{i - 1\}$ . s is termed the starting vertex.
- The difference between the sets of supernodes and fundamental supernodes is normally not large, with the latter having (slightly) more blocks in the resulting partitioning of *L*.
- Note that fundamental supernodes are independent of the choice of the postordering of T.



Figure: A matrix *A* and its postordered elimination tree  $\mathcal{T}$  for which the set of supernodes  $\{1, 2\}$  and  $\{3, 4, 5, 6\}$  and the set of fundamental supernodes  $\{1, 2\}, \{3, 4\}$  and  $\{5, 6\}$  are different. The filled entries in  $\mathcal{S}\{L + L^T\}$  are denoted by *f*.

Fundamental supernodes and leaves of row subtrees

#### Theorem

Assume T is postordered. Vertex s is the starting vertex of a fundamental supernode if and only if it has at least two child vertices in T or it is a leaf vertex of a row subtree of T.

### Proof.

If *s* has at least two child vertices then, from the definition of a fundamental supernode, it must be the starting vertex of a fundamental supernode. Assume that, for some i > s, s is a leaf vertex of  $\mathcal{T}_r(i)$ . If *s* is also a leaf vertex of  $\mathcal{T}$  then *s* is a starting vertex of a supernode. The remaining case is *s* having only one child. Because  $\mathcal{T}$  is postordered, this child must be s - 1. Theorem of necessary and sufficient conditions for the fill-in implies  $a_{is} \neq 0$  and  $a_{i,s-1} = 0$ , that is,  $i \in col_L\{s\}$  and  $i \notin col_L\{s-1\}$ . It follows that

$$\mathcal{S}\{L_{s-1:n,s-1}\} \subsetneqq \mathcal{S}\{L_{s:n,s}\} \cup \{s-1\},\$$

and vertices s and s-1 cannot belong to the same supernode. Hence, s is the starting vertex of a new fundamental supernode.

Conversely, assume that *s* is the starting vertex of a fundamental supernode *S*. If *s* has no child vertices or at least two child vertices, the result follows. If *s* has exactly one child vertex.  $^{242/609}$ 

### Complexity of symbolic operations

- Because fundamental supernodes are characterized by their starting vertices, they can be found by modifying Algorithm to incorporate marking leaf vertices of the row subtrees and vertices with at least two child vertices.
- Once the elimination tree has been computed, the complexity is O(n + nz(A)).
- The computation can be made even more efficient by using the skeleton graph  $\mathcal{G}(A^-)$ .
Supernodes and efficient computation

#### • the loop over rows has no indirect addressing: (dense BLAS1)



Supernodes and efficient computation

#### • the loop over rows has no indirect addressing: (dense BLAS1)



Supernodes and efficient computation

- the loop over rows has no indirect addressing: (dense BLAS1)
- the loop over columns of the updating supernode can be unrolled to save memory references (dense BLAS2)



Supernodes and efficient computation

- the loop over rows has no indirect addressing: (dense BLAS1)
- the loop over columns of the updating supernode can be unrolled to save memory references (dense BLAS2)



## Supernodes and efficient computation

- the loop over rows has no indirect addressing: (dense BLAS1)
- the loop over columns of the updating supernode can be unrolled to save memory references (dense BLAS2)
- parts of the updating supernode can be used for blocks of updated supernode (dense BLAS3)





# Outline



#### Numerical factorization: blocks and panels needed

- Positive definite (SPD) matrix is factorizable (strongly regular) and (in exact arithmetic) its Cholesky factorization  $A = LL^T$  exists.
- Because efficient implementations of sparse Cholesky factorizations rely heavily on exploiting dense blocks, we first consider algorithms for the Cholesky factorization of dense matrices that can be applied to such blocks.
- The first one is an in-place algorithm because *L* can overwrite the lower triangular part of *A* (thus reducing memory requirements if *A* is no longer required).

#### In-place dense factorization

# Algorithm (In-place dense left-looking Cholesky factorization)

Input: Dense SPD matrix A. Output: Factor L such that  $A = LL^T$ .

1: for j = 1 : n do 2:  $L_{j:n,j} = A_{j:n,j}$ 3: for k = 1 : j - 1 do 4:  $L_{j:n,j} = L_{j:n,j} - L_{j:n,k} l_{jk}$ 5: end for 6:  $l_{jj} = (l_{jj})^{1/2}$ 7:  $L_{j+1:n,j} = L_{j+1:n,j}/l_{jj}$ 8: end for

> Only the lower triangular part of A is required

▷ Update column *j* using previous columns

Overwrite the diagonal entry with its square root
 Scale off-diagonal entries in column j

#### Reorganized using block panels (left-looking)

# Algorithm (In-place dense left-looking panel Cholesky factorization)

Input: Dense SPD matrix A with nb panels. Output: Factor L such that  $A = LL^T$ .

1: for 
$$jb = 1 : nb$$
 do  
2:  $L_{jb:nb,jb} = A_{jb:nb,jb}$   
3: for  $kb = 1 : jb - 1$  do  
4:  $L_{jb:nb,jb} = L_{jb:nb,jb} - L_{jb:nb,kb} L_{jb,kb}^T$   $\triangleright$  Update block column  $jb$   
5: end for  
6: Compute in-place factorization of  $L_{jb,jb}$   $\triangleright$  Overwrite  $L_{jb,jb}$  with its Cholesky factor  
7:  $L_{jb+1:nb,jb} = L_{jb+1:nb,j} L_{jb,jb}^{-T}$   $\triangleright$  Dense triangular solve  
8: end for

## Reorganized using panels (right-looking)

# Algorithm (In-place dense right-looking panel Cholesky factorization)

Input: Dense SPD matrix A in the form with nb panels. Output: Factor L such that  $A = LL^{T}$ .

1: for jb = 1 : nb do 2:  $L_{ib:nb,ib} = A_{ib:nb,ib}$ 3: end for 4: for jb = 1 : nb do 5: Compute in-place factorization of  $L_{ib,ib}$  $\triangleright$  Overwrite  $L_{jb,jb}$  with its Cholesky factor 6:  $L_{jb+1:nb,jb} = L_{jb+1:nb,j} L_{ib,jb}^{-T}$ Dense triangular solve 7: for kb = jb + 1 : nb do 8.  $L_{kb:nb,kb} = L_{kb:nb,kb} - L_{kb:nb,ib} L_{kb,ib}^{T}$ g. end for 10: end for

#### Large panels split to blocks, right-looking

# Algorithm (In-place dense right-looking block Cholesky factorization)

Input: Dense SPD matrix A in the form (2) with  $nb \times nb$  blocks. Output: Factor L such that  $A = LL^T$ .

1: for jb = 1 : nb do 2:  $L_{ib:nb,ib} = A_{ib:nb,ib}$ 3: end for 4: for jb = 1 : nb do 5: ▷ Task factorize(jb) Compute in-place factorization of  $L_{ib,ib}$ 6: for ib = jb + 1 : nb do 7:  $L_{ib,jb} = L_{ib,jb} L_{ib,jb}^{-T}$ ▷ Task solve(ib, jb) 8: for kb = jb + 1 : ib do 9:  $L_{ib,kb} = L_{ib,kb} - L_{ib,ib} L_{kb,ib}^T$ ▷ Task update(ib, jb, kb) 10: end for 11: end for 12: end for

#### Tasks for efficient parallelization (still dense case)

- The panel and block descriptions of the factorization enable efficient parallelization. The three main block operations, which are called tasks, are factor(*jb*), solve(*ib*, *jb*) and update(*ib*, *jb*, *kb*).
- There are the following dependencies between the tasks.
   factorize(jb) depends on update(jb, kb, jb) for all kb = 1,..., jb 1.
   solve(ib, jb) depends on update(ib, kb, jb) for all kb = 1,..., jb 1, and factorize(jb).
   update(ib, jb, kb) depends on solve(ib, kb), solve(jb, kb).
- A dependency graph (DAG) can be used to schedule the tasks.

## Sparse Cholesky factorizations

- Several classes of algorithms that implement sparse Cholesky factorizations.
- Their major differences relate to how they schedule the computations.
- First discussed: straightforward extension of the dense Cholesky factorization.

## Sparse Cholesky factorizations

The entries of L satisfy the relationship

$$L_{j+1:n,j} = \left(A_{j+1:n,j} - \sum_{k=1}^{j-1} L_{j+1:n,k} l_{jk}\right) / l_{jj} \quad \text{with} \quad l_{jj} = \left(a_{jj} - \sum_{k=1}^{j-1} l_{jk}^2\right)^{\frac{j}{2}}$$

#### Theorem

The numerical values of the entries in column j > k of L depend on the numerical values in column k of L if and only if  $l_{jk} \neq 0$ .

#### Sparse Cholesky factorizations

 Straightforward sparse: S{L} already determined, static storage formats based, for example, on compressed columns and/or rows.

Algorithm (Simplified sparse left-looking Cholesky)

Input: SPD matrix A and sparsity pattern  $S\{L\}$ . Output: Factor L such that  $A = LL^T$ .

1: $l_i$	$a_j = a_{ij}$ for all $(i, j) \in \mathcal{S}\{L\}$	Filled entries in L are initialised to zero
2: fc	or $j = 1:n$ do	
3:	for $k \in \{k < j     l_{jk}  eq 0\}$ do	
4:	for $i \in \{i \geq j     l_{ik}  eq 0\}$ do	
5:	$l_{ij} = l_{ij} - l_{ik}l_{jk}$	
6:	end for	
7:	end for	
8:	$l_{jj} = (l_{jj})^{1/2}$	
9:	for $i \in \{i > j     l_{ij}  eq 0\}$ do	
10:	$l_{ij} = l_{ij} / l_{jj}$	
11:	end for	
12: e	end for	

## Sparse Cholesky factorizations

## Algorithm (Simplified sparse right-looking Cholesky)

Input: SPD matrix A and sparsity pattern  $S\{L\}$ . Output: Factor L such that  $A = LL^T$ .

```
1: For all (i, j) \in \mathcal{S}\{L\} set l_{ij} = a_{ij}
2: for j = 1 : n do
3: l_{ii} = (l_{ii})^{1/2}
4: for i \in \{i > j | l_{ij} \neq 0\} do
5:
             l_{ii} = l_{ii}/l_{ii}
6:
      end for
7:
      for k \in \{k > j \mid l_{kj} \neq 0\} do
8:
             for i \in \{i \ge k | l_{ij} \neq 0\} do
9:
                 l_{ik} = l_{ik} - l_{ij} l_{kj}
10:
              end for
11:
          end for
12: end for
```

 $\triangleright$  Filled entries in L are initialised to zero

## Sparse Cholesky factorizations

- The theoretical background based on the elimination tree  $\mathcal{T}$  enables the dependencies to be searched for efficiently.
- $\mathcal{T}$  allows:
  - ▶ row (or column) counts of *L* to be computed
  - storage allocated
  - supernodes
- In practice, it can be beneficial to split large supernodes into smaller panels to better conform to computer caches.

Block column (left-looking) algorithm: notes on implementation

• The following theorem shows that we need to go through rows of columns of *L* computed so far within a block algorithm.

#### Theorem

Let j > k. Numerical values of entries in  $L_{*j}$  depend on the values of entries in  $L_{*k}$  iff  $l_{jk} \neq 0$ .

• Note that to get the sparsity patterns of columns we need less.

Block left-looking algorithm: notes on implementation

- Construction columns (block columns) one by one.
- Going through the rows can be simulated by linked lists as we saw in CSC-CSR matvec.
- Plan to mention this again when discussing approximate factorizations.

#### Sparse Cholesky factorizations

- An alternative for sparse matrices is to compute *L* one row at a time. This is sometimes called an up-looking factorization.
- Asymptotically optimal. Difficult to incorporate high level BLAS.
- The following relation holds for the *i*-th row of L

$$L_{i,1:i-1}^T = L_{1:i-1,1:i-1}^{-1} A_{1:i-1,i} \quad \text{with} \quad l_{ii}^2 = a_{ii} - L_{i,1:i-1} L_{i,1:i-1}^T.$$

The application of  $L_{1:i-1,1:i-1}^{-1}$  can be implemented by solving the triangular system

$$L_{1:i-1,1:i-1}y = A_{1:i-1,i},$$

and setting  $L_{i,1:i-1}^T = y$ .

#### Triangular systems and sparsity pattern of L

#### Theorem

Consider a sparse lower triangular matrix L and the DAG  $\mathcal{G}(L^T)$  with vertex set  $\{1, 2, ..., n\}$  and edge set  $\{(j \longrightarrow i) \mid l_{ij} \neq 0\}$ . The sparsity pattern  $\mathcal{S}\{y\}$  of the solution y of the system Ly = b is the set of all vertices reachable in  $\mathcal{G}(L^T)$  from  $\mathcal{S}\{b\}$ .

## Proof.

From factorization Algorithm and assuming the non-cancellation assumption, we see that (a) if  $b_i \neq 0$  then  $y_i \neq 0$  and (b) if for some  $j < i, y_j \neq 0$  and  $l_{ij} \neq 0$  then  $y_i \neq 0$ . These two conditions can be expressed as a graph transversal problem in  $\mathcal{G}(L^T)$ . (a) adds all vertices in  $\mathcal{S}\{b\}$  to the set of visited vertices and (b) states that if vertex j has been visited then all its neighbours in  $\mathcal{G}(L^T)$  are added to the set of visited vertices. Thus  $\mathcal{S}\{y\} = \mathcal{R}each(\mathcal{S}\{b\}) \cup \mathcal{S}\{b\}$ .

#### Triangular systems and sparsity pattern of L

Figure illustrates the sparsity patterns of a lower triangular matrix L and vector b together with  $\mathcal{G}(L^T)$ . The vertices that are reachable from  $\mathcal{S}\{b\} = \{2, 4\}$  are 5 and 6 and thus  $\mathcal{S}\{y\} = \{2, 4, 5, 6\}$ .



Figure: An example to illustrate L, b and  $\mathcal{G}(L^T)$ .

#### Triangular systems and sparsity pattern of L



 The only nonzero of the right-hand side implies the three nonzeros in the solution (34

Triangular systems and sparsity pattern of L





### Up-looking Cholesky - comments

- Sparse row Cholesky factorization is based on the repeated solution of triangular linear systems.
- Theorem above can be used to determine the sparsity pattern of row *i* at Step 3, that is, by finding all the vertices that are reachable in  $\mathcal{G}(L_{1:j-1,1:j-1}^T)$  from the set  $\{i \mid a_{ij} \neq 0, i < j\}$ .
- A depth-first search of G(L<sup>T</sup><sub>1:j-1,1:j-1</sub>) determines the vertices in the row sparsity patterns in topological order, and performing the numerical solves in that order correctly preserves the numerical dependencies.
- Another option is to find the row subtrees using  $\mathcal{T}(A)$ .

## Algorithm (Sparse up-looking Cholesky factorization)

**Input:** SPD matrix A. **Output:** Factor L such that  $A = LL^T$ .

- 1:  $l_{11} = (a_{11})^{1/2}$
- 2: for i = 2: n do
- *3:* Find  $S\{L_{i,1:i-1}\}$

4: 
$$L_{i,1:i-1}^T = L_{1:i-1,1:i-1}^{-1} A_{1:i-1,i}$$

5:  $l_{ii} = a_{ii} - L_{i,1:i-1}L_{i,1:i-1}^T$ 

6:  $l_{ii} = (l_{ii})^{1/2}$ 

7: end for

 $\triangleright$  Sparsity pattern of  $L_{i,1:i-1}$  $\triangleright$  Sparse triangular solve

#### Supernodal algorithms

- Sophisticated supernodal algorithms: efficiency in parallel computational environments.
- Assume supernodal structure. Arithmetic of dense trapezoidal matrices. Blocks termed a nodal matrix.



Figure: An illustration of a supernode (left), the corresponding nodal matrix (centre), and the nodal matrix with two panels (right).

## Comments on (right-looking) supernodal processing

- Once a supernode is ready to be factorized, a dense Cholesky factorization of the block on the diagonal of the nodal matrix is performed
- Then a triangular solve is performed with the computed factor and the rectangular part of the nodal matrix.
- Iterate over ancestors of the supernode in the assembly tree.
- For each parent, the rows of the current supernode for each parent columns are identified. The outer product of those rows and the subdiagonal part of the supernode (update operations).
- The resulting matrix can be held in a temporary buffer. The rows and columns of this buffer are matched against indices of the ancestors and are added to them in a sparse scatter operation. For efficiency, the updates may use panels so that the temporary buffer remains in cache.

DAG-based approach and the sparse case Each nodal matrix is subdivided into blocks.



Figure: An illustration of a blocked nodal matrix with two block columns. The first block on the diagonal is triangular and the second one is trapezoidal. The task factorize\_block is illustrated on the left and in the centre; the task solve\_block is illustrated on the right.

#### DAG-based approach and the sparse case

#### Different splitting into tasks.

• The key difference: distinguish: update\_internal (the same nodal matrix), update\_between (different nodal matrices).

factorize\_block( $L_{diag}$ ) Computes the dense Cholesky factor  $L_{diag}$  of the block on the diagonal (leftmost plot). If the block is trapezoidal, the factorization is followed by a triangular solve of its rectangular part  $L_{rect} = L_{rect}L_{diag}^{-T}$  (centre). solve\_block( $L_{dest}$ ) Performs a triangular solve of an off-diagonal block  $L_{dest}$  of the form  $L_{dest} = L_{dest}L_{diag}^{-T}$  (rightmost). update\_internal ( $L_{dest}, L_r, L_c$ ) Performs the update  $L_{dest} = L_{dest} - L_r L_c^T$ , where  $L_{dest}, L_r$  and  $L_c$  belong to the same nodal matrix.

update\_between  $(L_{dest}, L_r, L_c)$  Performs the update  $L_{dest} = L_{dest} - L_r L_c^T$ , where  $L_r$  and  $L_c$  belong to the same nodal matrix and  $L_{dest}$  belongs to a different nodal matrix.

#### DAG-based approach and the sparse case

- Again, the tasks are partially ordered and a task DAG is used to capture the dependencies.
- For example, the updating of a block of a nodal matrix from a block column of *L* that is associated with a descendant of the supernode has to wait until all the relevant rows of the block column are available.
- At each stage of the factorization, tasks will be executing (in parallel) while others are held (in a stack or pool of tasks) ready for execution.

Alternative way: the multifrontal method

#### Theorem

Let *A* be SPD and let  $\mathcal{T}$  be its elimination tree. The numerical values of entries in column *k* of the Cholesky factor *L* of *A* only affect the numerical values of entries in column *i* of *L* for  $i \in anc_{\mathcal{T}}\{k\}$   $(1 \le k \le n-1)$ .

#### Proof.

Setting  $S^{(1)}=A,$  for  $k\geq 2$  the  $(n-k+1)\times (n-k+1)$  Schur complement  $S^{(k)}$  can be expressed as

$$S^{(k)} = S^{(k-1)}_{k:n,k:n} - \begin{pmatrix} l_{k,k-1} \\ \vdots \\ l_{n,k-1} \end{pmatrix} \left( l_{k,k-1} \dots l_{n,k-1} \right) = S^{(k-1)}_{k:n,k:n} - L_{k:n,k-1} L^T_{k:n,k-1}.$$
(35)

Theorem on replication implies that all nonzero entries  $l_{ik}$  in column k of L explicitly used in the update are such that  $i \in anc_T\{k\}$ . Considering the Cholesky factorization as a sequence of Schur complement updates, only columns i with  $i \in anc_T\{k\}$  can be influenced numerically by form

Alternative way: the multifrontal method

- The computation of subsequent Schur complements by adding individual updates is straightforward; the multifrontal method employs further modifications and enhancements of this basic concept.
- Because the vertices of  $\mathcal{T}$  are topologically ordered, the order in which the updates are applied progresses up the tree from the leaf vertices to the root vertex.
- This allows the computation of  $S^{(k)}$  to be rewritten as

$$S^{(k)} = A_{k:n,k:n} - \sum_{j \in \mathcal{T}(k) \setminus \{k\}} L_{k:n,j} L_{k:n,j}^T,$$

emphasizing the role of  $\mathcal{T}$ .

 In place of Schur complements, the multifrontal method uses frontal matrices connected to subtrees of T.

#### Alternative way: the multifrontal method

- Assume  $k, k_1, \ldots, k_r$  are the row indices of the nonzeros in column k of L.
- The frontal matrix  $F_k$  of the *k*-th subtree T(k) of T is the dense  $(r+1) \times (r+1)$  matrix defined by

$$F_{k} = \begin{pmatrix} a_{kk} & a_{kk_{1}} & \dots & a_{kk_{r}} \\ a_{k_{1}k} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ a_{k_{r}k} & 0 & \dots & 0 \end{pmatrix} - \sum_{j \in \mathcal{T}(k) \setminus \{k\}} \begin{pmatrix} l_{kj} \\ l_{k_{1}j} \\ \vdots \\ l_{k_{r}j} \end{pmatrix} (l_{kj} \ l_{k_{1}j} \dots \ l_{k_{r}j})$$
(36)

One step of the Cholesky factorization of Fk can be written as

$$F_{k} = \begin{pmatrix} l_{kk} & 0 & \dots & 0 \\ l_{k_{1}k} & & & \\ \vdots & & I \\ l_{krk} & & \end{pmatrix} \begin{pmatrix} 1 & & \\ & V_{k} \end{pmatrix} \begin{pmatrix} l_{kk} & l_{k_{1}k} & \dots & l_{krk} \\ 0 & & & \\ \vdots & & I & \\ 0 & & & \end{pmatrix} (\mathbf{1} \\ V_{k} \end{pmatrix} \begin{pmatrix} l_{kk} & l_{k_{1}k} & \dots & l_{krk} \\ \vdots \\ l_{krk} \end{pmatrix} (l_{kk} & l_{k_{1}k} & \dots & l_{krk}) + \begin{pmatrix} 0 & & \\ & V_{k} & \\ & & V_{k} \end{pmatrix},$$
(38)

- *V<sub>k</sub>* is termed a generated element (it is also sometimes called an update matrix or a contribution block).
- The name "generated element" is because the multifrontal method has its origins in the simpler frontal method, which uses a single frontal matrix.

• Equating the last r rows and columns yields

$$V_{k} = -\sum_{j \in \mathcal{T}(k)} \begin{pmatrix} l_{k_{1}j} \\ \vdots \\ l_{k_{r}j} \end{pmatrix} \left( l_{k_{1}j} \ \dots \ l_{k_{r}j} \right).$$
(39)

Assume that  $c_j$  (j = 1, ..., s) are the children of k in  $\mathcal{T}$ . The set  $\mathcal{T}(k) \setminus \{k\}$  is the union of disjoint sets of vertices in the subtrees  $\mathcal{T}(c_j)$ . Each of these subtrees is represented in the overall update by the generated element. Thus,  $F_k$  can be written in an recursive form as follows

$$F_{k} = \begin{pmatrix} a_{kk} & a_{kk_{1}} & \dots & a_{kk_{r}} \\ a_{k_{1}k} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ a_{k_{r}k} & 0 & \dots & 0 \end{pmatrix} \Leftrightarrow V_{c_{1}} \Leftrightarrow \dots \Leftrightarrow V_{c_{s}}.$$
(40)
- Here, the operation ↔ denotes the addition of matrices that have row and column indices belonging to subsets of the same set of indices (in this case, k, k<sub>1</sub>,..., k<sub>r</sub>); entries that have the same row and column indices are summed. This is referred to as the extend-add operator.
- Adding a row and column of *A* and the generated elements into a frontal matrix is called the assembly.
- A variable is fully summed if it is not involved in any rows and columns of *A* that have still to be assembled or in a generated element.
- Once a variable is fully summed, it can be eliminated.
- A key feature of the multifrontal method is that the frontal matrices and the generated elements are compressed and stored without zero rows and columns as small dense matrices.
- Integer arrays used to maintain a mapping of the local indices to the global indices of *A* and its factors.
- Symmetry allows only the lower triangular part of these matrices to be held.

### Algorithm outlines the basic multifrontal method.

### Algorithm (Basic multifrontal Cholesky factorization)

**Input:** SPD matrix A and its elimination tree. **Output:** Factor L such that  $A = LL^T$ .

- 1: for k = 1 : n do
- 2: Assemble the frontal matrix  $F_k$  using (40)  $\triangleright$  Only the lower triangle is needed
- 3: Perform a partial Cholesky factorization of  $F_k$  using (37) to obtain column k of L and the generated element  $V_k$
- 4: end for

We have the following observation.

### Observation

Each generated element  $V_k$  is used only once to contribute to a frontal matrix  $F_{parent(k)}$ . Furthermore, the index list for the frontal matrix  $F_k$  is the set of row indices of the nonzeros in column k of the Cholesky factor L.

- In practical implementations, efficiency is improved by using the assembly tree because it allows more than one elimination to be performed at once.
- This is outlined here. kb is used to index the frontal matrix on the kb-th step (1 ≤ kb ≤ nsup).

## Algorithm (Multifrontal Cholesky factorization using the assembly tree)

**Input:** SPD matrix A and its assembly tree. **Output:** Factor L such that  $A = LL^T$ .

- 1: nelim = 0 ▷ nelim is the number of eliminations performed 2: for kb = 1 : nsup do ▷ nsup is the number of supernodes
- 3: Assemble the frontal matrix  $F_{kb}$ ; let *l* be the number of fully summed variables in  $F_{kb}$ .
- 4: Perform a block partial Cholesky factorization of  $F_{kb}$  to obtain columns nelim + 1 to nelim + l of L and the generated element  $V_{kb}$
- 5: nelim = nelim + l
- 6: end for

As an example, consider the matrix and its assembly tree given above. The nsup = 5 supernodes are  $\{1, 2\}, 3, 4, 5, \{6, 7, 8, 9\}$  and so variables 1 and 2 can be eliminated together on the first step. Assembling rows/columns 1 and 2 of the original matrix, the frontal matrix  $F_1$  and generated element  $V_1$  have the structure

where f denotes fill-in entries (only the lower triangular entries are stored in practice). Similarly,

$$F_{2} = \begin{array}{c} 3 & 4 & 8 \\ 3 & \left( \begin{array}{c} * & * & * \\ * & * & * \\ 8 & \left( \begin{array}{c} * & * & * \\ * & * & * \end{array} \right), \quad V_{2} = \begin{array}{c} 4 & \left( \begin{array}{c} * & * \\ * & * \end{array} \right).$$

The frontal matrix  $F_3$  and generated element  $V_3$  are given by

$$F_{3} = \frac{4}{7} \begin{pmatrix} 4 & 7 & 8 & 7 & 8 \\ * & * & * \\ * & * & * \\ * & * & * \end{pmatrix} \Leftrightarrow V_{2}, \quad V_{3} = \frac{7}{8} \begin{pmatrix} * & f \\ f & * \end{pmatrix}.$$

Then

$$F_{4} = \begin{cases} 5 & 7 & 8 \\ 7 & \left( \begin{array}{c} * & * & * \\ * & * \\ 8 & \left( \begin{array}{c} * & * & * \\ * & * \\ * & * \end{array} \right), \quad V_{4} = \frac{7}{8} \begin{pmatrix} * & f \\ f & * \end{pmatrix},$$

and, finally, with kb = 5 we have

- An important implementation detail is how and where to store the generated elements.
- The partial factorization of  $F_{kb}$  at supernode kb can be performed once the partial factorizations at all the vertices belonging to the subtree of the assembly tree with root vertex kb are complete.
- If the vertices of the assembly tree are ordered using a depth-first search, the generated elements required at each stage are the most recently computed ones amongst those that have not yet been assembled.
- This makes it convenient to use a stack.

- The memory demands of the multifrontal method can be very large. Auxiliary storage can be used.
- The ordering of the children of a vertex in the assembly tree can significantly affect the required stack size.
- Tree and node parallelism.







### Multifrontal method: (example 1)

$$F_{1} = \begin{bmatrix} 1 & 6 & 8 & 9 \\ 6 & 8 & * & * \\ 8 & 8 \\ 9 & * & * \end{bmatrix}, \quad V_{1} = \begin{bmatrix} 6 & 8 & 9 \\ 8 & 8 \\ * & * & * \\ 8 & 9 \\ * & * & * \end{bmatrix},$$

Here  $V_1$  is dense and f denotes fill-in entries. Similarly, we have

$$F_{2} = \begin{pmatrix} 2 & 4 & 7 & & & & & & \\ 2 & \begin{pmatrix} * & * & * \\ * & & \\ * & & \\ * & & \end{pmatrix}, \quad V_{2} = \begin{pmatrix} 4 & 7 & & & & & & & & & \\ 4 & 7 & \begin{pmatrix} * & f \\ f & * \\ \end{pmatrix}, \quad F_{3} = \begin{pmatrix} 3 & & & & & & & \\ 5 & & & & \\ * & & \\ * & & & \\ * & & \\$$

The sparsity pattern of the frontal matrix  $F_4$  is then

and an an

285/609

### Multifrontal method

### Theorem

Let  $\mathcal{T}$  be postordered. Assume each computed generated element  $(V_x)$  is pushed onto a stack. Then when constructing the (frontal matrix)  $F_j$ , the required generated elements are on the top of the stack. They can be directly popped from the stack and assembled into  $F_j$ .

**Proof sketch:** 1) Vertices of each subtree of the postordered  $\mathcal{T}$  form an interval. 2) Denote  $c_l$ , l = 1, ..., s children of j in  $\mathcal{T}$ . 3) Each  $c_l$  is the root of a subtree  $\mathcal{T}(c_l)$ . 4) Once the frontal matrix  $F_{c_l}$  for a leaf of  $\mathcal{T}(c_l)$  is constructed, all its children have been processed and the generated  $V_{c_l}$  is pushed onto the stack. 5) That is, all subtrees  $\mathcal{T}_{c_l}$ , l = 1, ..., s are fully assembled into the generated elements before  $F_j$  can be constructed. 6) If  $F_j$  is ready to be assembled (step j), the s generated  $V_{c_l}$ , l = 1, ..., s are on the top of the stack.







## Direct methods: Multifrontal method







### Multifrontal method: (example 2)



stack

### Multifrontal method: summary

- Right-looking (submatrix) method
- Does not form the Schur complement directly. Instead, the updates are moved to a stack as dense matrices and used when needed.
- The processing order is based on the elimination tree
- We will see that in order to have the needed updates at the stack top, postordering is needed.
- Specific postorderings used to minimize the needed amount of memory.
- Now example, properties repeated once more later.

Multifrontal method: assumptions and properties

- We do need to have the entries from the stack readily available.
- ullet ightarrow elimination tree should be postordered
- Arithmetic of dense matrices
- Connection with the frontal method (later) is relatively week.
- One of the most important methods for the sparse direct factorization.

### Multifrontal method: postorderings memory issues



- First case: Maximum stack size may be  $1 \times 1+2 \times 2+3 \times 3+4 \times 4$
- Second case: Maximum stack size may be  $4 \times 4$
- Conclusion: Even postorderings can be very different with respect to algorithmic/architectural needs

## Outline

Sparse LU factorization

### 296/609

### Graphs and LU factorization

- Structural changes during a sparse Cholesky factorization described by graphs.
- In particular, using elimination tree.
- For LU: more possibilities to capture structure by graphs
- But, often factorizability has to be assumed: much harder combine pivoting and structural predictions.

# Sparse LU factorization of generally nonsymmetric matrices

### LU factorization and DAGs



• **directed acyclic graphs** capture structure of the factors. We use  $G(L^T)$  (*L* by columns) and G(U) (*U* by rows).



### LU factorization and DAGs

- The first graph model uses the elimination DAGs associated with *L* and *U*: see previous slide.
- Observation, generalizes predictions to the nonsymmetric case.

### Observation

If i > j and  $u_{ji} \neq 0$  then the column replication principle states

$$S{L_{i:n,j}} \subseteq S{L_{i:n,i}},$$

that is, the pattern of column j of L (rows i to n) is replicated in the pattern of column i of L.

• Analogously, if i > j and  $l_{ij} \neq 0$  then the row replication principle states

 $\mathcal{S}\{U_{j,i:n}\} \subseteq \mathcal{S}\{U_{i,i:n}\},\$ 

that is, the pattern of row j of U (columns i to n) is replicated in the pattern of row i of U.



Figure: An illustration of the column and row replication principles of sparse LU factorizations. Left: the matrix A. Centre: showing column replication. Right: also row replication. Filled entries not involved in the demonstration and resulting from subsequent steps of the factorization are denoted in black.

## Sparse LU

### Basic sparse LU factorization

• Assumed that *A* is factorizable so that pivoting is not needed (we will remind this only sometimes).

### Algorithm (Basic sparse LU factorization)

Input: Nonsymmetric and factorizable matrix  $A = L_A + D_A + U_A$ . Output: LU factorization A = LU.

1:  $L = I + L_A$ ▷ Strictly lower triangular part of A 2:  $U = D_A + U_A$ ▷ Diagonal plus strictly upper triangular part of A 3: for k = 1 : n - 1 do 4: for  $i \in \{i > k \mid l_{ik} \neq 0\}$  do 5:  $l_{ik} = l_{ik}/u_{kk}$ 6:  $U_{i,i:n} = U_{i,i:n} - U_{k,i:n}l_{ik}$  $\triangleright$  Update row i of U 7: end for 8: for  $j \in \{j > k \mid u_{kj} \neq 0\}$  do 9:  $L_{i+1:n,i} = L_{i+1:n,i} - L_{i+1:n,k} u_{ki}$  $\triangleright$  Update column j of L 10: end for 11: end for

### **Recursive replications**

The following theorem formulates the recursive column replication and the replication of nonzeros along rows of L using directed paths in  $\mathcal{G}(U)$ .

### Theorem

Assume that for some k < j there is a directed path  $k \xrightarrow{\mathcal{G}(U)} j$ . Then

$$\mathcal{S}\{L_{j:n,k}\} \subseteq \mathcal{S}\{L_{j:n,j}\}.$$
(42)

Moreover, if  $l_{ik} \neq 0$  for some i > j then  $l_{is} \neq 0$  for all vertices s on this path.

• An analogous result holds for the rows of U and directed paths in  $\mathcal{G}(L^T)$ .

### Generalizing the necessary and sufficient condition for fill-in

### Theorem

If  $a_{ij} = 0$  and i > j then there is a filled entry  $l_{ij} \neq 0$  if and only if there exists k < j such that  $a_{ik} \neq 0$  and there is a directed path  $k \xrightarrow{\mathcal{G}(U)} j$ .

#### Theorem

If  $a_{ij} = 0$  and i < j then there is a filled entry  $u_{ij} \neq 0$  if and only if there exists k < i such that  $a_{kj} \neq 0$  and there is a directed path  $k \xrightarrow{\mathcal{G}(L^T)} i$ .

## Sparse LU

Consider the path  $1 \rightarrow 3 \rightarrow 5 \rightarrow 6$  in  $\mathcal{G}(U)$ . Its existence implies the fill-in in *L*, first in column 3, then in columns 5 and 6. Similarly, the path  $2 \rightarrow 4 \rightarrow 5 \rightarrow 6$  in  $\mathcal{G}(L^T) \Rightarrow$  fill-in at (4,7), (5,7) and (6,7) in *U*.



### Transitive reduction

- To employ  $\mathcal{G}(L^T)$  and  $\mathcal{G}(U)$  in efficient algorithms, they need to be simplified. Transitive reductions are sparser and preserve reachability within the graphs.
- A subgraph  $\mathcal{G}^0 = (\mathcal{V}, \mathcal{E}^0)$  is a **transitive reduction** of  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  if the following conditions hold:
  - (*T*1) there is a path from vertex i to vertex j in  $\mathcal{G}$  if and only if there is a path from i to j in  $\mathcal{G}^0$  (reachability condition), and
  - (T2) there is no subgraph with vertex set  $\mathcal{V}$  that satisfies (T1) and has fewer edges (minimality condition).
- A transitive reduction is unique for a DAG, as shown in the following theorem and illustrated below.

## Sparse LU factorization of generally nonsymmetric matrices

### Transitive reduction

• Transitive reduction is the loop interconnecting vertices 1, 2 and 3.



#### Theorem

*Transitive reduction of a* **directed acyclic graph** *is unique and is a subgraph of G.* 

## Sparse LU factorization of generally nonsymmetric matrices

### Transitive reduction



## Sparse LU

### Theorem

Let  $\mathcal{G}$  be a DAG. The transitive reduction  $\mathcal{G}^0$  of  $\mathcal{G}$  is unique and is the subgraph that has an edge for every path in  $\mathcal{G}$  and has no proper subgraph with this property.



Figure: Example to show the transitive reduction of a DAG.  $\mathcal{G}$  is on the left, its transitive reduction  $\mathcal{G}^0$  is in the centre, and one possible  $\mathcal{G}'$  that is equireachable with  $\mathcal{G}$  is on the right.

### Transitive reduction

If  $S{A}$  is symmetric, the role of the transitive reduction is played by the elimination tree.

### Theorem

If A is symmetrically structured then the transitive reduction of the DAG  $\mathcal{G}(L^T)$  (=  $\mathcal{G}(U)$ ) is the elimination tree  $\mathcal{T}(A)$ .

## Sparse LU

#### **Transitive reduction**



Figure: The sparsity patterns of L + U of a symmetrically structured A, the DAG  $\mathcal{G}(L^T)$  (left) and the elimination tree  $\mathcal{T}(A)$  (right). Straightforward to see that  $\mathcal{T}(A)$  is the transitive reduction of  $\mathcal{G}(L^T)$ .

## Sparse LU

### Transitive reduction may be expensive to obtain

- Obtaining the exact transitive reduction of a DAG can be expensive. Instead, approximate reductions that drop the minimality condition may be computed.
- A directed graph  $\mathcal{G}'$  with the same vertex set as  $\mathcal{G}$  that satisfies condition (T1) is said to be **equireachable** with  $\mathcal{G}$ .
- This is something in between the DAG and transitive reduction. Of course, due to the reachability, a lot of theoretical results are satisfied.

### Theorem

Assume  $\mathcal{G}'$  is equireachable with  $\mathcal{G}(U)$  and for some k < j there is a directed path  $k \xrightarrow{\mathcal{G}'} j$ . Then the replication theorem can use the reduced DAGs. Moreover, if  $l_{ik} \neq 0$  for some i > j then  $l_{is} \neq 0$  for all vertices s on the directed path.
#### Equireachability and efficiency

- Equireachability enables sparse triangular linear systems to be solved more efficiently being sparser.
- The necessary and sufficient conditions for the fill-in from above:

#### Theorem

If  $a_{ij} = 0$  and i > j then there is a filled entry  $l_{ij} \neq 0$  if and only if there exists k < j such that  $a_{ik} \neq 0$  and a directed path  $k \xrightarrow{\mathcal{G}'(U)} j$ , where  $\mathcal{G}'(U)$  is equireachable with  $\mathcal{G}(U)$ .

#### Theorem

If  $a_{ij} = 0$  and i < j then there is a filled entry  $u_{ij} \neq 0$  if and only if there exists k < i such that  $a_{kj} \neq 0$  and a directed path  $k \xrightarrow{\mathcal{G}'(L^T)} i$ , where  $\mathcal{G}'(L^T)$  is equireachable with  $\mathcal{G}(L^T)$ .

#### Equireachability: example

Figure 36 depicts  $\mathcal{G}(U)$  and  $\mathcal{G}'(U)$  for the matrix in Figure above.



Figure: The DAG  $\mathcal{G}(U)$  for the matrix from Figure **??** (left) and  $\mathcal{G}'(U)$  which is equireachable with  $\mathcal{G}(U)$  (right).

#### Column sparsity patterns (for *L*)

Standard description of the sparsity patterns of the columns of L can be obtained from the Schur complement as follows:

$$\mathcal{S}\{L_{j:n,j}\} = \mathcal{S}\{A_{j:n,j}\} \bigcup_{k < j, u_{kj} \neq 0} \mathcal{S}\{L_{j:n,k}\}, \quad 1 \le j \le n.$$

Theorem on the next slide implies that not all the terms in this union are needed to obtain  $S\{L_{j:n,j}\}$ .

#### Column sparsity patterns (for L)

#### Theorem

If  $\mathcal{G}'(U)$  is equireachable with  $\mathcal{G}(U)$  then

 $\mathcal{S}\{L_{j:n,j}\} = \mathcal{S}\{A_{j:n,j}\} \bigcup_{(k \to j) \in \mathcal{E}(\mathcal{G}'(U))} \mathcal{S}\{L_{j:n,k}\}, \quad 1 \le j \le n.$ (43)

#### Proof.

Consider an edge  $(k \to j)$  in  $\mathcal{G}(U)$  but not in  $\mathcal{G}'(U)$ . Repeatedly applying replication results along the directed path  $k \xrightarrow{\mathcal{G}'(U)} j$ , we see that  $L_{j:n,k}$  is contained in the right-hand side of the predicted structure and therefore  $\mathcal{S}\{L_{j:n,j}\}$  is contained in the right-hand side of this structure as well. Because the right-hand side of the formula is trivially contained in the left-hand side, the result follows.

#### Row sparsity patterns (for U)

An analogous result holds for the rows of U.

#### Theorem

If  $\mathcal{G}'(L)$  is equireachable with  $\mathcal{G}(L)$  then

$$\mathcal{S}\{U_{i,i:n}\} = \mathcal{S}\{A_{i,i:n}\} \bigcup_{(k \to i) \in \mathcal{E}(\mathcal{G}'(L^T))} \mathcal{S}\{U_{k,i:n}\}, \quad 1 \le i \le n$$

As an example, consider the matrix above. Because  $(3 \rightarrow 5)$  is the only edge of  $\mathcal{G}'(U)$  in the union on the right-hand side of (43),  $\mathcal{S}\{L_{5:7,5}\}$  is given by

$$S{L_{5:7,5}} = S{A_{5:7,5}} \cup S{L_{5:7,3}}.$$

We can see this from the graph  $\mathcal{G}'(U)$  in Figure that demonstrates equireachability (top right).

#### Factorization by bordering

- Factorization by bordering can be used to obtain  $S{L}$  by rows and  $S{U}$  by columns.
- Assume the sparsity patterns of the first *k* − 1 rows of *L* and the first *k* − 1 columns of *U* (1 < *k* ≤ *n*) have been computed.
- At step k, the matrix  $A_{1:k,1:k}$  is

$$\begin{pmatrix} A_{1:k-1,1:k-1} & A_{1:k-1,k} \\ A_{k,1:k-1} & a_{kk} \end{pmatrix} = \begin{pmatrix} L_{1:k-1,1:k-1} & 0 \\ L_{k,1:k-1} & 1 \end{pmatrix} \begin{pmatrix} U_{1:k-1,1:k-1} & U_{1:k-1,k} \\ 0 & u_{kk} \end{pmatrix}$$
(44)

• Equating terms for the (2,1) block, row k of L satisfies

$$L_{k,1:k-1}U_{1:k-1,1:k-1} = A_{k,1:k-1},$$

or, equivalently, if y denotes the off-diagonal part of the column k of  $L^T$  then it is the solution of the lower triangular system

$$U_{1:k-1,1:k-1}^T y = A_{k,1:k-1}^T.$$

#### Factorization by bordering

- The sparsity pattern of y is the set of all vertices reachable in the DAG  $\mathcal{G}(U_{1:k-1,1:k-1})$  (or in a graph that is equireachable with it) from the nonzeros in  $A_{k,1:k-1}$ .
- Similarly, equating terms in (44) for the (1,2) block, column k of U satisfies

$$L_{1:k-1,1:k-1}U_{1:k-1,k} = A_{1:k-1,k}.$$

- Again, its sparsity pattern can be determined by searching the DAG  $\mathcal{G}(L_{1:k-1,1:k-1}^T)$ .
- The diagonal entry  $u_{kk}$  is then computed as  $a_{kk} L_{k,1:k-1}U_{1:k-1,k}$ .
- Determining the sparsity patterns of *L* and *U* and computing their numerical values is coupled: computation of the factors needs be mutually interleaved because computing part of one requires information from a part of the other.

#### Pruning of the elimination DAGs

- The matrix in the centre is the same as the one on the left except that the entries in positions (4, 6) and (6, 4) have been removed (that is, pruned).
- Both matrices have the same sets of reachable vertices in  $\mathcal{G}(L^T)$  and  $\mathcal{G}(U)$ . This suggests how to find  $\mathcal{G}'(L^T)$  and  $\mathcal{G}'(U)$  that are equireachable with  $\mathcal{G}(L^T)$  and  $\mathcal{G}(U)$ , respectively.

	1	<b>2</b>	3	4	5	6		1	<b>2</b>	3	4	5	6		1	<b>2</b>	3	4	5	6
1	/*		*			* \	1	/*		*			* \	1	/*		*			
2	*	*				)	<b>2</b>	*	*					2	*	*				)
3			*			*	3			*			*	3			*			*
4		*		*	*	*	4		*		*	*		4		*		*	*	
5				*	*	*	5				*	*	*	5				*	*	*
6	/*			*	*	*/	6	/*				*	*/	6					*	*/

Figure: An example of symmetric pruning. On the left is  $S{L+U}$ . In the centre is the reduced sparsity pattern obtained by symmetric pruning. On the right is the reduced sparsity pattern that results from symmetric path pruning.

#### Pruning of the elimination DAGs

#### Theorem

If for some j < s both  $l_{sj} \neq 0$  and  $u_{js} \neq 0$ , then there are no edges  $(j \rightarrow k)$  with k > s in the transitive reductions of  $\mathcal{G}(U)$  and  $\mathcal{G}(L^T)$ .

#### Proof.

Let  $(j \to k)$  be an edge of  $\mathcal{G}(U)$ , that is,  $u_{jk} \neq 0$ . Because  $l_{sj} \neq 0$  and  $u_{jk} \neq 0$  implies that  $u_{sk} \neq 0$ , there is a path  $j \to s \to k$  in  $\mathcal{G}(U)$  and the edge  $(j \to k)$  does not belong to the transitive reduction of  $\mathcal{G}(U)$ . The result for  $\mathcal{G}(L^T)$  can be seen analogously.  $\Box$ 

The theorem implies that if for some s > 1 there are edges

$$j \xrightarrow{\mathcal{G}(L^T)} s \text{ and } j \xrightarrow{\mathcal{G}(U)} s,$$

then all edges  $(j \to k)$  in  $\mathcal{G}(U)$  and  $\mathcal{G}(L^T)$  with k > s can be pruned. The resulting DAGs  $\mathcal{G}'(U)$  and  $\mathcal{G}'(L^T)$  have fewer edges and are equireachable with  $\mathcal{G}(U)$  and  $\mathcal{G}(L^T)$ , respectively. The removal of redundant edges based on the Theorem is called **symmetric pruning**. Other ways of pruning of the elimination DAGs

If for some s > 1 there are paths

$$j \xrightarrow{\mathcal{G}(L^T)} s \text{ and } j \xrightarrow{\mathcal{G}(U)} s,$$

then for all k > s symmetric path pruning removes the edges  $(j \rightarrow k)$  from  $\mathcal{G}(U)$  and  $\mathcal{G}(L^T)$ .

Consider again previous Figure. In the centre: the sparsity pattern after symmetric pruning. On the right: the reduced sparsity pattern that results from symmetric path pruning. The edge (1 → 6) is not required in G'(L<sup>T</sup>) or G'(U) because there are paths

$$1 \xrightarrow{\mathcal{G}(L^T)} 2 \xrightarrow{\mathcal{G}(L^T)} 4 \xrightarrow{\mathcal{G}(L^T)} 5 \xrightarrow{\mathcal{G}(L^T)} 6 \quad \text{and} \quad 1 \xrightarrow{\mathcal{G}(U)} 3 \xrightarrow{\mathcal{G}(U)} 6$$

#### Another graph model: nonsymmetric elimination tree

- The elimination DAGs  $\mathcal{G}(L)$  and  $\mathcal{G}(U)$  can be combined into a single structure called the **nonsymmetric elimination tree** in which edges are replaced by paths.
- This can be advantageous because it is more compact.
- If  $S{A}$  is symmetric then its elimination tree is defined in terms of the mapping

$$parent(j) = \min\{i \mid i > j \text{ and } l_{ij} \neq 0\}.$$

The condition  $l_{ij} \neq 0$  is equivalent to  $i \xrightarrow{\mathcal{G}(L)} j \xrightarrow{\mathcal{G}(L^T)} i$ .

In the nonsymmetric case, the definition can be generalized using directed paths

$$parent(j) = \min\{i \mid i > j \text{ and } i \xrightarrow{\mathcal{G}(L)} j \xrightarrow{\mathcal{G}(U)} i\}.$$
 (45)

#### Example of the next slide

• Vertices 6, 8 and 10 are the only ones with cycles of the form

$$i \xrightarrow{\mathcal{G}(L)} 2 \xrightarrow{\mathcal{G}(U)} i.$$

Namely,

$$6 \xrightarrow{\mathcal{G}(L)} 2 \xrightarrow{\mathcal{G}(U)} 5 \xrightarrow{\mathcal{G}(U)} 6, \quad 8 \xrightarrow{\mathcal{G}(L)} 2 \xrightarrow{\mathcal{G}(U)} 8$$

and

$$10 \xrightarrow{\mathcal{G}(L)} 6 \xrightarrow{\mathcal{G}(L)} 2 \xrightarrow{\mathcal{G}(U)} 10.$$

• In this example, parent(2) = 6.



-. . .

#### Locality: way to find the nonsymmetric elimination tree

• Theorem can be regarded as a generalization of corollary on the standard symmetric elimination tree.

#### Theorem

Let A be a nonsymmetric matrix. i = parent(j) if and only if i > j and i is the smallest vertex that belongs to the same strong component of  $\mathcal{G}(A_{1:i,1:i})$  as vertex j.

- The result is employed in Algorithm below to find the elimination tree.
- The complexity of finding the strong components of a digraph with *m* edges and *n* vertices is O(n + m) time. Hence, the complexity of Algorithm is O(nz(A)n). More sophisticated approaches with complexity  $O(nz(A)\log n)$  exist.

#### Locality: way to find the nonsymmetric elimination tree

## Algorithm (**Basic computation of the elimination tree for nonsymmetric** *A*)

Input: Digraph  $\mathcal{G}(A)$ . Output: The elimination tree given by the mapping parent.

```
1: parent(1:n) = 0
2: for i = 1 : n do
3:
       Find the vertex set \mathcal{V}_C of the strong component of \mathcal{G}(A_{1:i,1:i}) that contains i
4:
     for i \in \mathcal{V}_C \setminus \{i\} do
5:
           if parent(j) = 0 then
6:
               parent(j) = i
7:
           end if
8:
     end for
9:
       parent(i) = 0
10: end for
```

#### Illustration using the matrix from above

- The main loop sets the first nonzero value in the array *parent* when i = 3 because this is the first *i* for which the set  $\mathcal{V}_C \setminus \{i\}$  is non empty; it is equal to  $\{1\}$  and thus parent(1) = i = 3.
- For i = 4, the vertex set  $\{1, 3, 4\}$  forms a strong component of  $\mathcal{G}(A_{1:4,1:4})$  and so parent(3) = 4.
- For i = 5, the single vertex  $\{5\}$  is a strong component of  $\mathcal{G}(A_{1:5,1:5})$  and, therefore, 5 is not a parent of any other vertex (it is a leaf vertex).
- G(A<sub>1:6,1:6</sub>) has two strong components with vertex sets {1,3,4} and {2,5,6}. i = 6 belongs to the second of these and thus the algorithm sets parent(j) = i = 6 for j = 2 and 5.

#### Another graph model: column elimination tree

# • An attractive idea for constructing $S\{L+U\}$ and subsequently computing the LU factorization is based on using the **column** elimination tree $T(A^TA)$ .

#### Theorem

Assume all the diagonal entries of *A* are nonzero and let  $\widehat{L}\widehat{L}^T$  be the Cholesky factorization of  $A^TA$ . Then for any row permutation matrix *P* such that PA = LU

$$\mathcal{S}\{L+U\} \subseteq \mathcal{S}\{\widehat{L}+\widehat{L}^T\}.$$

- Holds for any row permutation matrix *P* applied to *A*. This allows partial pivoting.
- The following result states that  $\mathcal{T}(A^T A)$  represents the potential dependencies among the columns in an LU factorization.
- For strong Hall matrices no tighter prediction is possible from the sparsity structure of *A*.

#### Theorem

If PA = LU is any factorization of A with partial pivoting then the following hold.

- **()** If vertex *i* is an ancestor of vertex *j* in  $\mathcal{T}(A^T A)$  then  $i \ge j$ .
- 2 If  $l_{ij} \neq 0$  then vertex *i* is an ancestor of vertex *j* in  $\mathcal{T}(A^T A)$ .
- If  $u_{ij} \neq 0$  then vertex j is an ancestor of vertex i in  $\mathcal{T}(A^T A)$ .

Suppose in addition that A is a strong Hall matrix. If l = parent(k) in  $\mathcal{T}(A^T A)$ , then there are values of the nonzero entries of A for which  $u_{kl} \neq 0$ .

Figures below illustrate the differences in the sparsity patterns of A and  $A^T A$  and of their factors; the corresponding elimination trees are also given.





Figure: Both figures: the sparsity patterns of A and L + U (top) and of  $A^T A$  and  $\hat{L} + \hat{L}^T$ , where  $A^T A = \hat{L}\hat{L}^T$  (bottom). Filled entries are denoted by f. The corresponding elimination trees are also given.

#### The column elimination tree

- A potential problem with the column elimination tree:
- $S{A^TA}$  can have significantly more entries than  $S{L+U}$ . An extreme example is when *A* has one or more dense rows because  $A^TA$  is then fully dense.

#### Supernodes and LU

- Supernodes group together columns of the factors with the same nonzero structure, allowing them to be treated as a dense submatrix for storage and computation.
- For nonsymmetric matrices, supernodes are harder to characterize.
- The need to incorporate pivoting means it may not be possible to predict the sparsity structures of the factors before the numerical factorization and they must be identified on the fly.
- More ways to define supernodes.
- Cholesky solver: fundamental supernodes are made contiguous by symmetrically permuting the matrix according to a postordering of its elimination tree; this does not change the sparsity of the Cholesky factor.
- For nonsymmetric A, before the numerical factorization,  $\mathcal{T}(A^T A)$  can be constructed and the columns of A then permuted according to its postordering to bring together supernodes.

#### Theorem

Let A have column elimination tree  $\mathcal{T}(A^T A)$ . Let p be a permutation vector such that if  $p_i$  is an ancestor of  $p_j$  in  $\mathcal{T}(A^T A)$  then  $i \ge j$ . Let Pbe the permutation matrix corresponding to p and let  $\widehat{A} = PAP^T$ . Then  $\mathcal{T}(\widehat{A}^T \widehat{A})$  is isomorphic to  $\mathcal{T}(A^T A)$ ; in particular, relabelling each vertex i of  $\mathcal{T}(\widehat{A}^T \widehat{A})$  as  $p_i$  yields  $\mathcal{T}(A^T A)$ . If, in addition,  $\widehat{A} = \widehat{L}\widehat{U}$  is an LU factorization without pivoting then  $P^T \widehat{L}P$  and  $P^T \widehat{U}P$  are lower triangular and upper triangular matrices, respectively, so that  $A = (P^T \widehat{L}P)(P^T \widehat{U}P)$  is also an LU factorization.

- In practice, for many matrices the average size of a supernode is only 2 or 3 columns and many comprise a single column.
- Larger artificial supernodes may be created by merging vertex j with its parent vertex i in  $\mathcal{T}(A^T A)$  if the subtree rooted at i has fewer than some chosen number of vertices.

#### Multifrontal LU

- The multifrontal method can be generalized to nonsymmetric *A* by modifying the definitions of the frontal matrices and generated elements to conform to an LU factorization.
- Natural generalizations to rectangular frontal and generated element matrices do not simultaneously satisfy the statements from above. Rewritten for the LU factorization:
  - (a) Each generated element  $V_j$  is used only once to contribute to a frontal matrix.
  - (b) The row and column index lists for the rectangular frontal matrix  $F_j$  correspond to the nonzeros in column  $L_{j:n,j}$  and nonzeros in row  $U_{j,j:n}$ , respectively.
- An approach that satisfies (a) can be based on the sparsity pattern of S{A + A<sup>T</sup>} and storing some explicit zeros if S{A} is not symmetric. It is then analogous to the symmetric multifrontal method.
- Performs well if  $S{A}$  is close to symmetric.

#### Multifrontal LU

 An approach that satisfies (b) and not necessarily (a) splits the generated elements into smaller ones that are embedded into further rectangular frontal matrices. We illustrate this using the example from above.

where \* are entries in A and filled entries in L + U are denoted by f.

#### Multifrontal LU

$$F_{1} = \begin{array}{c} 1 & 3 \\ 3 \\ * \\ * \\ * \\ * \end{array} \right), V_{1} = \begin{array}{c} 3 \\ 2 \\ 3 \\ * \\ * \end{array} \right)$$

To construct  $F_2$  that satisfies (b) we can only use part of  $V_1$ . Because  $a_{13} \neq 0$ , the sparsity pattern of column 1 is replicated in that of column 3 of the factors. The entry in position (2,3) belongs to  $F_2$ . The row replication implies that the remaining entries contribute to  $F_3$  and so we split  $V_1$  into two as follows

$$V_1^2 = 2 \ (f), \quad V_1^3 = \frac{3}{8} \ {\binom{*}{f}}, \quad V_1 = V_1^2 \Leftrightarrow V_1^3,$$

where  $\Leftrightarrow$  is the extend-add operator and  $V_1^2$  and  $V_1^3$  contribute to  $F_2$ and  $F_3$ , respectively.

#### Multifrontal LU Then $F_2$ and the corresponding generated element $V_2$ are

$$F_{2} = \begin{pmatrix} 2 & 5 & 8 & 10 \\ 2 & \begin{pmatrix} * & * & * & * \\ * & & & \\ * & & & \end{pmatrix} \Leftrightarrow V_{1}^{2} = \begin{pmatrix} 2 & 3 & 5 & 8 & 10 \\ 2 & \begin{pmatrix} * & f & * & * & * \\ * & & & & \\ * & & & & \end{pmatrix}, V_{2} = \begin{pmatrix} 3 & 5 & 8 & 10 \\ 4 & f & f & f \\ 5 & f & & f \end{pmatrix}.$$

Consider the following splitting of  $V_2$ 

$$V_2 = {}^6_8 \begin{pmatrix} f \\ f \end{pmatrix} \Leftrightarrow {}^6_8 \begin{pmatrix} f \\ * & f \end{pmatrix} \equiv V_2^3 \Leftrightarrow V_2^5 \Leftrightarrow V_2^6.$$

The next frontal matrix is

$$F_{3} = {3 \atop 4} \begin{pmatrix} 3 & 4 \\ * & * \\ * & * \end{pmatrix} \Leftrightarrow V_{1}^{3} \Leftrightarrow V_{2}^{3} = {3 \atop 6} \begin{pmatrix} * & * \\ * & * \\ 6 \\ 8 \end{pmatrix}, \quad V_{3} = {4 \atop 6} \begin{pmatrix} * \\ f \\ f \\ f \end{pmatrix}.$$

The subsequent steps can be described in a similar way.

#### Multifrontal LU

• Theorem expresses the nested relationship between the nonsymmetric multifrontal method and the nonsymmetric elimination tree.

#### Theorem

Assume A is a general nonsymmetric matrix and t = parent(k) in  $\mathcal{T}(A)$ . Then

$$\mathcal{S}\{L_{t:n,k}\} \subseteq \mathcal{S}\{L_{t:n,t}\} \text{ and } \mathcal{S}\{U_{k,t:n}\} \subseteq \mathcal{S}\{U_{t,t:n}\}.$$

#### Proof.

Because *t* is the parent of *k*, by definition  $t \xrightarrow{\mathcal{G}(L)} k \xrightarrow{\mathcal{G}(U)} t$ . If  $u_{ij} \neq 0$  then a multiple of column *i* is added to column *j* during the LU factorization. Thus, by a simple induction argument, for each *j* on the path  $k \xrightarrow{\mathcal{G}(U)} t$ , we must have  $\mathcal{S}\{L_{j:n,k}\} \subseteq \mathcal{S}\{L_{j:n,j}\}$ . In particular, this holds for column *t*. The second part follows by a similar argument using the path  $t \xrightarrow{\mathcal{G}(L)} k$ .  $\Box$ 

#### **Multifrontal LU**

- The theorem shows that the parent relationship in the nonsymmetric elimination tree guarantees that both row and column replications can be applied at the same time.
- Hence all entries of the submatrices of the generated element  $V_k$  with indices greater than or equal to parent(k) can be added to  $V_{parent(k)}$  using the operation  $\Leftrightarrow$ .
- To illustrate this, consider again the 10 × 10 example above for which *parent*(1) = 3. Theorem guarantees that V<sub>1</sub> can be embedded into F<sub>3</sub> because S{L<sub>3:n,1</sub>} ⊆ S{L<sub>3:n,3</sub>} and S{U<sub>1,3:n</sub>} ⊆ S{U<sub>3,3:n</sub>}.

#### Preprocessing for LU

- Consider the case when A does not have a full transversal (that is, it has one or more zeros on the diagonal).
- For numerical stability and to reduce the number of permutations required during the factorization, it can be useful to permute *A* before the factorization begins to put large nonzero entries on the diagonal.
- Given a graph G = (V, E), an edge subset M ⊆ E is called a matching (or assignment) if no two edges in M are incident to the same vertex.
- The cardinality of a matching is the number of edges in it. A maximum cardinality matching (or maximum matching) is a matching of maximum cardinality. A matching is perfect if all the vertices are matched.

- A **bipartite graph** is an undirected graph whose vertices can be partitioned into two disjoint sets such that no two vertices within the same set are adjacent, that is, each set is an **independent** set.
- Let the  $n \times n$  matrix A have entries  $\{a_{ij'}\}$ . Associated with A is a bipartite graph defined as a triple  $\mathcal{G}_b = (\mathcal{V}_{row}, \mathcal{V}_{col}, \mathcal{E})$ , where the row vertex set  $\mathcal{V}_{row} = \{i \mid a_{ij'} \neq 0\}$  and the column vertex set  $\mathcal{V}_{col} = \{j' \mid a_{ij'} \neq 0\}$  correspond to the rows and columns of A and there is an (undirected) edge  $(i, j') \in \mathcal{E}$  if and only if  $a_{ij'} \neq 0$ .
- In Figure below we use prime to distinguish between the independent set of row vertices and the independent set of column vertices, that is, *i* denotes a row vertex and *i'* denotes a column vertex.



Figure: A sparse matrix its bipartite graph  $\mathcal{G}_b$  (left), perfect matching.

- If *A* is structurally nonsingular, a matching  $\mathcal{M}$  in  $\mathcal{G}_b$  is perfect if it has cardinality *n*.
- A perfect matching defines an  $n \times n$  permutation matrix Q with entries  $q_{ij}$  given by

$$q_{ij} = \begin{cases} 1, & \text{if } (j, i') \in \mathcal{M}, \\ 0, & \text{otherwise.} \end{cases}$$

• Both *QA* and *AQ* have the matching entries on the (zero-free) diagonal.

• *Q* and the column permuted matrix *AQ* for the example in Figure above.



Figure: The permutation matrix Q and the column permuted matrix AQ corresponding to the matrix above. The matched entries are on the diagonal of AQ.

#### Augmenting paths

- If a perfect matching exists, it can be found using augmenting paths.
- A path P in a graph is an ordered set of edges in which successive edges are incident to the same vertex. P is called an *M*-alternating path if the edges of P are alternately in *M* and not in *M*.
- An *M*-alternating path is an *M*-augmenting path in *G<sub>b</sub>* if it connects an unmatched column vertex with an unmatched row vertex. Note that the length of an *M*-augmenting path is an odd integer.

• Let  $\mathcal M$  and  $\mathcal P$  be subsets of  $\mathcal E$  and define the symmetric difference

 $\mathcal{M} \oplus \mathcal{P} := (\mathcal{M} \setminus \mathcal{P}) \cup (\mathcal{P} \setminus \mathcal{M}),$ 

• This is set of edges that belongs to either  $\mathcal{M}$  or  $\mathcal{P}$  but not to both.

- If *M* is a matching and *P* is an *M*-augmenting path, then *M* ⊕ *P* is a matching with cardinality |*M*|+1.
- Growing the matching in this way is called augmenting along  $\mathcal{P}$ .

#### Theorem

A matching  $\mathcal{M}$  in an undirected graph is a maximum matching if and only if there is no  $\mathcal{M}$ -augmenting path
## Algorithm (Maximum matching algorithm)

*Input:* An undirected graph. *Output:* Output maximum matching.

Find an initial matching M
 while there exists a M-augmenting path P do
 M = M ⊕ P
 end while

 $\triangleright \textit{ For example, } \mathcal{M} = \emptyset$ 

 $\triangleright$  Augment the matching along  ${\cal P}$ 

# Sparse LU

#### Augmenting paths: demonstration

- On the left is a bipartite graph with a matching with cardinality 5. An augmenting path  $2 \Longrightarrow 3' \Longrightarrow 3 \Longrightarrow 4' \Longrightarrow 4 \Longrightarrow 2'$  shown.
- Augmenting the matching along this path, the cardinality of the matching increases to 6 and M ⊕ P is a perfect matching.



Figure: Search for a perfect matching using augmenting paths.

# Sparse LU

### Weighted matchings

- While the maximum matching algorithm finds a permutation of *A* such that the permuted matrix has nonzero diagonal entries, there are more sophisticated variations that aim to ensure the absolute values of the diagonal entries of the permuted matrix (or their product) are in some sense large.
- The problem: given an  $n \times n$  matrix A, find a matching of the rows to the columns such that the product of the matched entries is maximized.
- That is, find a permutation vector *q* that maximizes

$$\prod_{i=1}^{n} |a_{iq_i}|. \tag{46}$$

• Define a matrix C corresponding to A with entries  $c_{ij'} \ge 0$  as follows:

$$c_{ij'} = \begin{cases} \log(\max_i |a_{ij'}|) - \log |a_{ij'}|, & \text{if } a_{ij'} \neq 0\\ \infty, & \text{otherwise.} \end{cases}$$
(47)

 It is straightforward to see that finding a *q* that solves the problem is equivalent to finding a *q* that minimizes

$$\sum_{i=1}^{n} |c_{iq_i}|,\tag{48}$$

- This is equivalent to finding a minimum weight perfect matching in an edge weighted bipartite graph.
- This is a well-studied problem and is known as the bipartite weighted matching or linear sum assignment problem.

### Weighted matchings: formally

- If  $\mathcal{G}_b = (\mathcal{V}_{row}, \mathcal{V}_{col}, \mathcal{E})$  is the bipartite graph associated with A then let  $\mathcal{G}_b(C) = (\mathcal{V}_{row}, \mathcal{V}_{col}, \mathcal{E})$  be the corresponding weighted bipartite graph in which each edge  $(i, j') \in \mathcal{E}$  has a weight  $c_{ij'} \ge 0$ .
- The weight (or cost) of a matching  $\mathcal{M}$  in  $\mathcal{G}_b(C)$ , denoted by  $csum(\mathcal{M})$ , is the sum of its edge weights; i.e.,

$$csum(\mathcal{M}) = \sum_{(i,j')\in\mathcal{M}} c_{ij'}.$$

A perfect matching *M* in *G<sub>b</sub>(C)* is said to be a minimum weight perfect matching if it has smallest possible weight i.e., csum(*M*) ≤ csum(*M*) for all possible perfect matchings *M*.

## Weighted matchings: shortest augmenting paths

- The key concept for finding a minimum weight perfect matching is that of a **shortest augmenting path**.
- An *M*-augmenting path *P* starting at an unmatched column vertex is called **shortest** if

$$csum(\mathcal{M}\oplus\mathcal{P})\leq csum(\mathcal{M}\oplus\widehat{\mathcal{P}})$$

for all other possible  $\mathcal M\text{-augmenting paths }\widehat{\mathcal P}$  starting at the same column vertex.

#### Weighted matchings: primal-dual algorithm

A matching *M<sub>e</sub>* is **extreme** if and only if there exist *u<sub>i</sub>* and *v<sub>j'</sub>* (which are termed **dual variables**) satisfying

- This is employed by the MC64 algorithm.
- The dual variables can be used for scaling.
- The algorithm starts with a feasible solution and corresponding extreme matching and then proceeds to iteratively increase its cardinality by one by constructing a sequence of shortest augmenting paths until a perfect extreme matching is found.

# Sparse LU

- Can be more efficient if a large initial matching is used.
- For example, Step 3 can be replaced by setting

 $u_i = \min\{c_{ij'} | j' \in S\{A_{i,1:n}\}\}$  for  $i \in \mathcal{V}_{row}$  and  $v_{j'} = \min\{c_{ij'} - u_i | i \in S\{A_{1:n,j'}\}\}$  for  $j' \in \mathcal{V}_{col}$ . In Step 4, an initial extreme matching can be determined from the edges for which

 $c_{ij'} - u_i - v_{j'} = 0.$ 

## Algorithm (Outline of the MC64 algorithm)

#### Input: Matrix A.

**Output:** A matching M and dual variables  $u_i, v_{j'}$ .

1: Define the weights  $c_{ij'}$  using (47)

2: Construct the weighted bipartite graph  $\mathcal{G}_b(C) = (\mathcal{V}_{row}, \mathcal{V}_{col}, \mathcal{E})$ 

3: Set  $u_i = 0$  for  $i \in \mathcal{V}_{row}$  and  $v_{j'} = \min\{c_{ij'} : (i, j') \in \mathcal{E}\}$  for  $j' \in \mathcal{V}_{col}$   $\triangleright$  Initial solution

4: Set 
$$\mathcal{M} = \{(i, j') | u_i + v_{j'}\}$$

- 5: while *M* is not perfect do
- 6: Find the shortest augmenting path  $\mathcal{P}$  with respect to  $\mathcal{M}$
- 7: Augment the matching  $\mathcal{M} = \mathcal{M} \oplus \mathcal{P}$
- 8: Update  $u_i$ ,  $v_{i'}$  so that (49) is satisfied for new  $\mathcal{M}$
- 9: end while

 $\triangleright$  make  $\mathcal{M}$  extreme

Initial extreme matching

- Number of potential problems with the MC64 algorithm.
  - The runtime is hard to predict and depends on the initial ordering of A.
  - It is a serial algorithm and as such it can represent a significant fraction of the total factorization time of a direct solver.
- Because the complexity of Step 6 of Algorithm is O((n + nz(A)) log n) and the complexity of Step 7 is O(n) and of Step 8 is O(n + nz(A), MC64 has a worst-case complexity of O(n(n + nz(A)) log n).
- In practice, this bound is not achieved and the algorithm is widely used.

#### Extension to non-square matrices

- The **Dulmage-Mendelsohn decomposition** is based on matchings and it generalizes the block triangular form.
- It comprises row and column permutations P and Q such that

$$PAQ = \begin{array}{ccc} \mathcal{C}_{1} & \mathcal{C}_{2} & \mathcal{C}_{3} \\ \mathcal{R}_{1} & A_{1} & A_{4} & A_{6} \\ 0 & A_{2} & A_{5} \\ \mathcal{R}_{3} & 0 & 0 & A_{3} \end{array} \right).$$
(50)

Here  $A_1$  is an  $m_1 \times n_1$  underdetermined matrix ( $m_1 < n_1$  or  $m_1 = n_1 = 0$ ),  $A_2$  is an  $m_2 \times m_2$  square matrix and  $A_3$  is an  $m_3 \times n_3$  overdetermined matrix ( $m_3 > n_3$  or  $m_3 = n_3 = 0$ ).

•  $A_1^T$  and  $A_3$  are strong Hall matrices but  $A_2$  need not be a strong Hall matrix.

### DM decomposition: an example

An example decomposition for a  $10 \times 10$  matrix is given in Figure 43. Here  $\mathcal{R} = \{1, 2, \dots, 9\}$  and  $\mathcal{C} = \{2, 3, \dots, 10\}$ .



Figure: An example of a coarse Dulmage-Mendelsohn decomposition. The blue entries belong to the maximum matching.  $m_1 = 3$ ,  $m_2 = 4$ ,  $m_3 = 3$ ,  $n_1 = 4$ ,  $n_2 = 4$ ,  $n_3 = 2$ . Column 1 and row 10 are unmatched.

## Outline



Backward stability

- Practical computations are invariably based on finite precision arithmetic.
- Computational algorithm

$$z = g(d)$$

for computing z as a function g of given data d.

- The algorithm is said to be backward stable if the computed solution <sup>2</sup> is the exact solution of <sup>2</sup> = g(d + Δd), where the perturbation Δd is "small" for all possible inputs d. What is meant by small depends on the context.
- If *d* is based on physical measurements (necessarily inaccurate), Δ*d* is small if it is of the same or smaller absolute value as the inaccuracies in determining *d*.

## Forward and backward errors

- The minimum absolute value  $|\Delta d|$  among such perturbations is called the (absolute) backward error (or, if divided by |d|, the relative backward error).
- The absolute and relative errors of  $\hat{z}$  are called forward errors, to distinguish them from the backward error.
- Backward stability is a property of the computational algorithm and to compute solutions with a small backward error we need to consider stable algorithms.

## Ill-conditioning

- A related concept is ill-conditioning.
- We say that the problem z = g(d) is ill-conditioned if small perturbations in the data d can lead to large changes in  $\hat{z}$ .
- The condition number measures how sensitive the output of a function is to its input. Ill-conditioning, which is measured in terms of the condition number, is a property of the problem.

## **III-conditioning**

 Provided the backward error, forward error, and the condition number are defined in a consistent manner, the following approximate inequality holds:

forward error  $\lesssim$  condition number  $\times$  backward error.

- This says that the computed solution to an ill-conditioned problem can have a large forward error because even if the computed solution has a small backward error, this error can be amplified by a large condition number.
- By preprocessing the problem it may be possible to improve its conditioning.

### Backward error result for LU

 $\epsilon$  denotes the machine precision.

## Theorem

Let the computed LU factorization of a matrix A be  $A + \Delta A = \hat{L}\hat{U}$ . The perturbation  $\Delta A$  that results from using finite precision arithmetic satisfies

$$||\Delta A||_{\infty} \le n O(\epsilon) ||\widehat{L}||_{\infty} ||\widehat{U}||_{\infty} + O(\epsilon^2).$$
(51)

Moreover, the computed solution  $\hat{x}$  of the linear system Ax = b satisfies  $(A + \Delta' A)\hat{x} = b$  with

$$||\Delta' A||_{\infty} \le n O(\epsilon) ||\widehat{L}||_{\infty} ||\widehat{U}||_{\infty} + O(\epsilon^2).$$
(52)

## Pivoting to improve stability

- At step k of GE, the computed diagonal entry a<sup>(k)</sup><sub>kk</sub> is termed the pivot (1 ≤ k < n).</li>
- Gaussian elimination breaks down if a zero pivot is encountered.
- Provided *A* is nonsingular, row interchanges can be incorporated to prevent this happening
- The systematic use of row permutations is called partial pivoting.
- If  $|a_{kk}^{(k)}|$  is very small (compared to other entries in the active submatrix) then it can cause difficulties in finite precision arithmetic because the absolute value of the corresponding computed multiplier  $l_{ik} = a_{ik}^{(k)}/a_{kk}^{(k)}$  can then be very large.
- Partial pivoting can be used to ensure  $|l_{ik}| \leq 1$ , that is, the rows of A that have not yet been pivoted on can be permuted so that the new pivot satisfies

$$\max_{i>k} |a_{ik}^{(k)}| \le |a_{kk}^{(k)}|.$$

## Pivoting to improve stability

• If  $P_k$  is the row permutation at stage k and  $P = P_{n-1}P_{n-2} \dots P_1$ then the computed factors of PA satisfy

$$||\widehat{L}||_{\infty} \leq n \quad \text{and} \quad ||\widehat{U}||_{\infty} \leq n \, \rho_{growth} ||A||_{\infty},$$

where the growth factor  $\rho_{growth}$  is defined to be

$$\rho_{growth} = \max_{i,j,k} \left( \left| a_{ij}^{(k)} \right| / \left| a_{ij} \right| \right).$$
(53)

• The bounds (51) and (52) can be rewritten as

 $||\Delta A||_{\infty} \le n^3 \, \rho_{growth} \, O(\epsilon) \, ||A||_{\infty}, \quad ||\Delta' A||_{\infty} \le n^3 \, \rho_{growth} \, O(\epsilon) \, ||A||_{\infty}.$ 

- In practice, these bounds are pessimistic and the actual errors are typically much smaller.
- Because backward stability of an LU factorization is influenced both by the initial ordering of *A* and the pivoting strategy, it is said to be conditionally backward stable.

## Stability and ill-conditioning

## Stability of Cholesky

- Pivoting not needed.
- Cholesky factorization of *A* is unconditionally backward stable. It enables the stable computation of the solution of the corresponding linear system.

## Theorem

Let the computed Cholesky factorization of a SPD matrix A be  $A + \Delta A = \widehat{L}\widehat{L}^T$ . The perturbation  $\Delta A$  that results from using finite precision arithmetic satisfies

 $||\Delta A||_{\infty} \le n^2 O(\epsilon) ||A||_{\infty}.$ 

Moreover, the computed solution  $\hat{x}$  of the linear system Ax = b satisfies  $(A + \Delta' A)\hat{x} = b$  with

$$||\Delta' A||_{\infty} \le n^2 O(\epsilon) ||A||_{\infty}.$$

#### Using the inverse instead of factorization

- No such stability results: The computed inverse is typically not the exact inverse of a nearby matrix  $A + \Delta A$  for any small perturbation  $\Delta A$ .
- Impractical to compute and store  $A^{-1}$ , regardless of how sparse A is: see below.

## Stability and ill-conditioning

Using the inverse instead of factorization

#### Theorem

If A is irreducible then the sparsity pattern  $S{A^{-1}}$  of its inverse is fully dense.

## Proof.

Without loss of generality, assume A is factorizable. For if not, there is a permutation matrix P such that the LU factorization of the row permuted matrix PA is factorizable. In this case, consider PA instead of A because for any permutation matrix P the inverse  $(PA)^{-1}$  is fully dense if and only if A is fully dense. Consider the matrix K of order 2n

$$K = \begin{pmatrix} A & I_n \\ I_n & 0 \end{pmatrix}$$

After applying *n* elimination steps to  $K = K^{(1)}$ , the order *n* active submatrix of  $K^{(n+1)}$  is  $-A^{-1}$ . Consider entry  $(A^{-1})_{ij}$   $(1 \le i, j \le n)$ . Because *A* is irreducible and the off-diagonal (1, 2) and (2, 1) blocks of *K* are equal to the identity matrix, there is a directed path  $i \Longrightarrow j$  in  $\mathcal{G}(K)$  such that the indices of all the intermediate vertices on the path are less than or equal to *n*. Theorem on fill paths and the non-cancellation assumption imply  $(A^{-1})_{ij} \ne 0$ . It follows that  $a_{69/609}$ 

## Using the inverse instead of factorization

• Theorem above implies that entries of  $A^{-1}$  correspond to paths in  $\mathcal{G}(A)$  even when A is not irreducible.

## Corollary

If *A* is factorizable then  $(A^{-1})_{ij} \neq 0$  ( $1 \leq i, j \leq n$ ) if and only if there exists a path  $i \xrightarrow{\mathcal{G}(A)} j$ .

### Pivoting strategies for dense matrices

Convention here: all the quantities (such as a<sup>(k)</sup><sub>ij</sub>) are the computed quantities.

#### Partial pivoting

- Partial pivoting interchanges rows at each step of the factorization to select the entry of largest absolute value in its column as the next pivot.
- If partial pivoting is used, it is straightforward to show that the growth factor satisfies

$$\rho_{growth} \le 2^{n-1}.$$

• Can be achieved in nontrivial cases. But, generally pessimistic, particularly when *n* is very large.

### Complete pivoting

- A much smaller bound can be obtained if complete (or full) pivoting is used: complete pivoting chooses the pivot to be the largest entry (in absolute value) in the active submatrix,
- That is, at stage k the pivot  $a_{kk}^{(k)}$  is chosen so that

$$\max_{i \ge k, j \ge k} |a_{ij}^{(k)}| \le |a_{kk}^{(k)}|.$$

Then

$$\rho_{growth} \le n^{1/2} (2.3^{1/2}.4^{1/3} \dots n^{1/(n-1)})^{1/2}.$$
(54)

- Can be expensive.
- Relaxations in practice.

## Rook pivoting

- Rook pivoting: more restrictive than. partial pivoting but cheaper than complete pivoting
- The pivot is chosen to be the largest entry in its row *and* its column, that is,

$$\max_{i>k} \left( |a_{ik}^{(k)}|, |a_{ki}^{(k)}| \right) \le |a_{kk}^{(k)}|.$$

- In practice, the cost is usually a small multiple of the cost of partial pivoting and significantly less than that of complete pivoting.
- The growth factor for rook pivoting satisfies

$$\rho_{growth} \le 1.5 \, n^{(3/4)\log n}$$

## $2\times 2~\mathrm{pivoting}$

- When the matrix *A* is symmetric but indefinite, it may not be possible to select pivots from the diagonal (for example, if all the diagonal entries of *A* are zero).
- If rows of A are permuted (so that off-diagonal entries are selected as pivots) then symmetry is destroyed, which means an LU factorization must be performed and this essentially doubles the cost of the factorization in terms of both storage and operation counts.
- Symmetry can be preserved by extending the notion of a pivot to  $2 \times 2$  blocks.

## Stability and ill-conditioning

### Symmetric indefinite

Consider

$$A = \begin{pmatrix} \delta & 1 \\ 1 & 0 \end{pmatrix}.$$

- If  $\delta = 0$ , an LDLT factorization in which D is a diagonal matrix does not exist.
- If  $\delta \ll 1$  then an LDLT factorization with D diagonal is not stable because  $\rho_{growth} = 1/\delta$ .
- If the LDLT factorization is generalized to allow D to be a block diagonal matrix with 1 × 1 and 2 × 2 blocks on the diagonal then a factorization is obtained that preserves symmetry and is nearly as stable as an LU factorization.

$$A = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix} = LDL^{T}.$$

Here D has one  $1 \times 1$  block and one  $2 \times 2$  block.

## Symmetric indefinite

- Rook pivoting can be extended to include 2 × 2 pivots. At each stage, an iterative procedure searches for an entry that is simultaneously the largest in absolute value in row *i* and column *j* of the active submatrix A<sup>(k)</sup>.
- This entry is used to build a symmetric  $2 \times 2$  pivot; the search terminates prematurely if a suitable  $1 \times 1$  pivot is found, that is, a pivot that satisfies a threshold test.
- The standard choice for the threshold comes from requiring the same potential maximal growth in the absolute values of the entries of the partially eliminated matrix that results from either two consecutive 1 × 1 pivots or one 2 × 2 pivot.
- It can be shown that the appropriate choice is  $(1 + \sqrt{17})/8$ .

$$\rho_{growth} < 3n\sqrt{2 \ 3^{1/2} 4^{1/3} \dots n^{1/(n-1)}},$$

## Stability and ill-conditioning

## Pivoting strategies for sparse matrices

• Preserving sparsity needed.

• Threshold partial pivoting: the pivot is chosen to satisfy

$$\max_{i>k} |a_{ik}^{(k)}| \le \gamma^{-1} |a_{kk}^{(k)}|, \tag{55}$$

where  $\gamma \in (0,1]$  is a chosen threshold parameter.

It is straightforward to see that

$$\max_{i} |a_{ij}^{(k)}| \le (1 + \gamma^{-1}) \max_{i} |a_{ij}^{(k-1)}|,$$

and

$$\max_{i} |a_{ij}^{(k)}| \le (1 + \gamma^{-1})^{nz_j} \max_{i} |a_{ij}|,$$

where  $nz_j$  is the number of off-diagonal entries in the *j*-th column of the U factor. Furthermore,

$$\rho_{growth} \le (1 + \gamma^{-1})^{nz_{cmax}},$$

where  $nz_{cmax} = \max_j nz_j \le n-1$ .

Pivoting strategies for sparse matrices

- A threshold can also be incorporated into rook pivoting. The pivot must then be at least γ times the absolute value of any other entry in its row and column of the active submatrix.
- In the symmetric case, if pivots are selected from the diagonal (to preserve symmetry), threshold partial pivoting is the same as threshold rook pivoting.
- Threshold rook pivoting has the potential to limit growth more successfully than threshold partial pivoting.

## Threshold $2\times 2$ pivoting and symmetry

- If *A* is a symmetric matrix then standard fill-reducing ordering algorithms (explained later) and the symbolic factorization phase employ only the sparsity pattern of *A*.
- In general, if *A* is indefinite, during the numerical factorization it is necessary to modify the chosen elimination order to maintain stability.
- If symmetry is to be preserved,  $1 \times 1$  and  $2 \times 2$  pivots are needed, resulting in an LDLT factorization in which *D* is a block diagonal matrix with  $1 \times 1$  and  $2 \times 2$  blocks.
- Limiting the size of the entries of L so that

$$|l_{ij}| \le \gamma^{-1} \tag{56}$$

for all i, j, together with a backward stable scheme for solving  $2 \times 2$  linear systems, suffices to show backward stability for the entire solution process.

#### Threshold $2 \times 2$ pivoting and symmetry

 In the sparse symmetric indefinite case, the stability test for a 1 × 1 pivot in column t of the active submatrix at stage k is the following standard threshold test

$$\max_{i \neq t, \, i \ge k} |a_{it}^{(k)}| \le \gamma^{-1} |a_{tt}^{(k)}|,\tag{57}$$

where  $\gamma$  is the threshold parameter. For a  $2\times 2$  pivot in rows and columns s and t the corresponding test is

$$\begin{pmatrix} a_{ss}^{(k)} & a_{st}^{(k)} \\ a_{st}^{(k)} & a_{tt}^{(k)} \end{pmatrix}^{-1} \begin{pmatrix} \max_{i \neq s, t; i \ge k} |a_{is}^{(k)}| \\ \max_{i \neq s, t; i \ge k} |a_{it}^{(k)}| \end{pmatrix} \le \gamma^{-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix},$$
(58)

where the absolute value of the matrix is interpreted element-wise.

## Not only bounding pivots

- In addition to bounding the size of the entries in *L*, the ability to stably apply the inverse of *D* to a vector is required.
- This is trivially the case for  $1 \times 1$  pivots, but for  $2 \times 2$  pivots it is necessary to check that the determinant  $|a_{ss}^{(k)}a_{tt}^{(k)} a_{st}^{(k)}a_{st}^{(k)}|$  is sufficiently large.
- A major difficulty when stability tests are incorporated into sparse factorizations is that a pivot satisfying the stability criteria may not exist.

## Stability and ill-conditioning

Stability and supernodes



Figure: An illustration of a simple nodal matrix. Pivot candidates are restricted to the square block  $L_{diag}$  on the diagonal.

- Pivots can only be chosen from the block  $L_{diag}$  on the diagonal (the block is square and symmetric)
- Entries in the off-diagonal block *L<sub>rect</sub>* are involved in the stability tests.
- Possibly delaying columns.

## Stability and supernodes

- All the off-diagonal entries in a block column must be fully updated before the block on the diagonal is factorized.
- The factorize\_block task and all the solve\_block tasks for a block column from the SPD case are combined into a single factorize\_column task.
- Fewer but larger tasks reducing the scope for parallelism.
- Delayed pivots arises in the multifrontal method: Frontal matrix F of order  $n_F$  of the form

$$F = \begin{pmatrix} F_{11} & F_{21}^T \\ F_{21} & F_{22} \end{pmatrix},$$
 (59)

where  $F_{11}$  is a  $p \times p$  matrix corresponding to the fully summed variables.

• Pivots can only be selected from  $F_{11}$  but the numerical values of the entries in  $F_{21}$  must be taken into account when testing for stability.
#### Sparse indefinite factorization

Algorithm (Simple partial sparse indefinite factorization)

*Input:* Symmetric indefinite matrix F of order  $n_F$  of the form (59) with  $F_{11}$  of order p; threshold  $\gamma \in (0, 0.5]$ . *Output:* Updated F; partial factors  $L_F$  and  $D_F$  and permutation  $P_F$ .

1:	$q = 0, t = 0$ $\triangleright q$ holds	he sum of the sizes (1 or 2) of the pivots chosen so far
2:	while $q < p$ do	
3:	find_pivot (piv_size)	⊳ See Algorithm 10.2
4:	if $(piv\_size = 0)$ exit while loop	Failed to find a pivot
5:	$q = q + piv\_size$	
6:	Update columns $q+1$ to $p$ of $F$	⊳ Right-looking
7:	end while	
8:	Apply updates to columns $p+1 \ {\rm to} \ n_F$	of $F$ $\triangleright$ Left-looking

### Algorithm (Find a pivot in F using threshold partial pivoting)

Input: F,  $L_F$ ,  $D_F$ ,  $P_F$ , p, q, t,  $\gamma$  are accessed from the environment of the call. Output: Pivot of size  $piv\_size$ ; columns  $q + 1:q + piv\_size$  of  $L_F$ ,  $D_F$ , updated  $P_F$  and t.

1: subroutine find_pivot (piv_size)				
$2: piv\_size = 0$				
3: for $test = 1 : p - q$ do				
4: $t = t$	$t + 1$ ; if $(t > p)$ set $t = q + 1$ $\triangleright$ Column t is search	hed for a pivot		
5: <b>if</b> (th	here is $s$ such that $q+1\leq s\leq t-1$ and $egin{pmatrix} f_{ss}&f_{st}\ f_{st}&f_{tt} \end{pmatrix}$ passes $2 imes 2$ piv	ot test) <b>then</b>		
<b>6</b> : p	$iv\_size = 2$			
7: S	Symmetrically permute rows and columns $q+1$ and $s$ of $F$	$\triangleright$ Update $P_F$		
8: S	Symmetrically permute rows and columns $q+2$ and $t$ of $F$	$\triangleright$ Update $P_F$		
9: C	Compute columns $q+1$ and $q+2$ of $D_F$ and $L_F$			
10: 1	return			
11: else	e if $(f_{tt} \text{ passes } 1 \times 1 \text{ pivot test})$ then			
12: p	$piv_size = 1$			
13: 5	Symmetrically permute rows and columns $q+1$ and $t$ of $F$	$\triangleright$ Update $P_F$		
14: (	Compute column $q + 1$ of $D_F$ and $L_F$			
15: r	return			
16: end	líf			
17: end for	<b>,</b>	385/609		

Algorithm

One step of full indefinite pivoting by Bunch and Parlett (1971) Set  $\alpha = (1 + \sqrt{17})/8 \approx 0.64$ Find  $a_{kk}$ : diagonal entry of maximum size Find  $a_{ij}$ : off-diagonal entry of maximum size (i < j)if  $|a_{kk}| \ge \alpha |a_{ij}|$  then use  $a_{kk}$  as  $1 \times 1$  pivot (ready for  $a_{kk} = 0$ ) else use  $\begin{pmatrix} a_{ii} & a_{ij} \\ a_{ji} & a_{jj} \end{pmatrix}$  as  $2 \times 2$  pivot end if

- Full pivoting: choosing entries of largest magnitudes: can be expensive.
- But a growth factor bound can be derived (only a slightly worse than for LU)

Algorithm

```
\alpha = (1 + \sqrt{17})/8 \approx 0.64, i = 1 (possible an initial preprocessing)
Find j \neq i such that a_{ii} = \max\{|a_{ki}|, k \neq i\} =: \lambda
if |a_{ii}| > \alpha \lambda then
     use a_{ii} as 1 \times 1 pivot
else
    \sigma = \max\{|a_{kj}|, k \neq j\}
    if |a_{ii}|\sigma > \alpha \lambda^2 then
         use a_{ii} as 1 \times 1 pivot
     else if |a_{ij}| \geq \alpha \sigma then
          use a_{ii} as a 1 \times 1 pivot
     else
   use \begin{pmatrix} a_{ii} & a_{ij} \\ a_{ji} & a_{jj} \end{pmatrix} as 2 \times 2 pivot
end if
                                                                                                        387/60
```

- This scheme shows why the previous algorithm (Bunch-Kaufmann) is useful to factorize sparse matrices
- The price for less searches expressed theoretically by large growth factor bound

- What if only f is used to decide such that  $|f| \ge \tau |\lambda|$ ?
- Again weaker bounds
- This can be practical if the processed matrix is not available in a current step (this can happen in multifrontal approaches). Implies rules like

### Algorithm

$$\begin{array}{l} \text{if } |d| \geq \alpha |\lambda| \text{ use } d \text{ as } 1 \times 1 \text{ pivot} \\ \text{if } |d\gamma| \geq \alpha |\lambda|^2 \text{: use } d \text{ as } 1 \times 1 \text{ pivot} \\ \text{if } |e| \geq \alpha |\gamma| \text{: use } e \text{ as } 1 \times 1 \text{ pivot} \\ \text{else} \\ \text{use } \begin{pmatrix} d & f \\ f & e \end{pmatrix} \text{ as } 2 \times 2 \text{ pivot} \end{array}$$

## Symmetric indefinite factorization



#### Relaxed and static pivoting

 If pivots are delayed during the numerical factorization then the data structures that were set up during the symbolic phase must be modified. This significantly complicates the development of general and symmetric indefinite sparse direct solvers.

If no candidate pivot satisfies the threshold test but the pivot that is nearest to satisfying it would satisfy it with a threshold  $\gamma_1 < \gamma$ , then provided  $\gamma_1$  is at least some chosen minimum value, relaxed pivoting accepts this pivot and reduces  $\gamma$  to  $\gamma_1$ .

With relaxed pivoting, delayed pivots can still occur and it may not be possible to use static data structures.

• A standard technique is to employ regularization (modification of *A*). This can avoid the need for an LDLT factorization in favour of a stable Cholesky factorization.

Special indefinite matrices that avoid pivoting

• Symmetric saddle point matrices: indefinite matrices of the form

$$A = \begin{pmatrix} G & R^T \\ R & -B \end{pmatrix},$$
 (60)

where  $G \in \mathbb{R}^{n_1 \times n_1}$  is a SPD matrix,  $B \in \mathbb{R}^{n_2 \times n_2}$  is a positive semidefinite matrix (including B = 0), and  $R \in \mathbb{R}^{n_2 \times n_1}$  with  $n_1 + n_2 = n$ .

• Define the permutation matrix P to be

$$P = [e_1, e_{n_1+1}, e_2, e_{n_1+2}, \dots e_{n_1}, e_n, e_{n_2+1}, \dots, e_{n_1}]^T$$

Then the permuted matrix  $PAP^T$  has a block form in which each entry  $A_{i,j}$  is a  $2 \times 2$  or  $2 \times 1$  or  $1 \times 2$  or  $1 \times 1$  block. In particular,

$$A_{i,i} = \begin{cases} \begin{pmatrix} g_{ii} & r_{ii} \\ r_{ii} & -b_{ii} \end{pmatrix}, & 1 \le i \le n_2 \\ b_{ii}, & n_2 + 1 \le i \le n_1 \end{cases}$$

• The following theorem shows that a 2 × 2 pivot updated by the Schur complement of a 1 × 1 pivot is nonsingular and vice versa.

#### Theorem

Let *A* be the symmetric saddle point matrix. Assume  $R = (R_1 \ R_2)$  is of full rank with  $R_1 \in \mathbb{R}^{n_2 \times n_2}$  nonsingular. Let  $G \in \mathbb{R}^{n_1 \times n_1}$  be SPD and partitioned conformally and let  $B \in \mathbb{R}^{n_2 \times n_2}$  be positive semidefinite matrix. If *A* is permuted to the form

$$\begin{pmatrix} G_{11} & R_1^T & G_{12} \\ R_1 & -B & R_2 \\ \hline G_{12}^T & R_2^T & G_{22} \end{pmatrix}$$

then the Schur complement of the symmetric indefinite matrix  $\begin{pmatrix} G_{11} & R_1^T \\ R_1 & -B \end{pmatrix}$  and the Schur complement of the SPD matrix  $G_{22}$  are nonsingular.

- Provided R is of full rank and  $R_1$  is nonsingular then the LDLT factorization of  $PAP^T$  exists, with  $2 \times 2$  pivots and  $1 \times 1$  pivots chosen from the diagonal blocks of  $PAP^T$  in any order.
- Assume all the  $2 \times 2$  pivots are selected ahead of the  $1 \times 1$  pivots. Then if B = 0 and  $|r_{ii}| \ge \max_{i \le j \le n_1} |r_{ij}|$   $(1 \le i \le n_2)$  then the growth factor is bounded by  $2^{2n_2}$ .
- A potential difficulty is that permutation matrices  $P_r$  and  $P_c$  are needed such that  $P_r R P_c = [R_1 \ R_2]$  with  $R_1$  nonsingular. If  $P_r$  and  $P_c$  can be constructed so that

$$P_r R P_c = \begin{pmatrix} R_{11} & R_{12} \\ & R_{22} \end{pmatrix}, \tag{61}$$

where  $R_{11}$  is upper triangular with nonzero diagonal entries then the permuted R is said to have a trapezoidal form.

#### Special indefinite matrices that avoid pivoting

● A simple case where *R* can be permuted this form is if it satisfies the following one-degree principle. Let *R* be of full rank and let *G<sub>b</sub>(R) = (V<sub>row</sub>, V<sub>col</sub>, E)* be the bipartite graph of *R* (Section **??**). *R* can be permuted to trapezoidal form if, for *k* = 1, 2, ..., *n*<sub>1</sub> − 1, the bipartite graph of *R<sup>(k)</sup>* has at least one vertex *j'<sub>k</sub> ∈ V<sub>col</sub>* of degree one, where *R<sup>(1)</sup> = R* and *R<sup>(k+1)</sup>* is obtained by removing from *R<sup>(k)</sup>* the column vertex *j'<sub>k</sub>* and its matched row index *i<sub>k</sub>* together with all edges involving *j'<sub>k</sub>* or *i<sub>k</sub>*.

- To illustrate the above mentioned case: consider the  $6 \times 8$  matrix R and its associated bipartite graph  $\mathcal{G}_b(R)$ .
- The first column vertex with degree one is 2'; it is matched with the row vertex 4. Deleting 2' and 4 removes edges {(4, 2'), (4, 3'), (4, 5'), (4, 6'), (4, 8')}. Column vertex 3' now has degree one; it is matched with row vertex 6.
- Repeating the process gives a perfect matching  $\mathcal{M} = \{(4, 1'), (6, 3'), (1, 4'), (5, 5'), (2, 6'), (3, 8')\}$  together with row and column matched vertex sets  $\{4, 6, 1, 5, 2, 3\}$  and  $\{2', 3', 4', 5', 1', 6'\}$ , respectively, and permutation matrices  $P_r$  and  $P_c$  of order 6 and 8 can be defined to obtain the trapezoidal form in Figure 45.





- If after k ≥ 1 steps all columns of the reduced matrix R<sup>(k)</sup> have degree greater than 1, the permuted matrix has the form above, where R<sub>11</sub> is k × k upper triangular, R<sub>12</sub> is k × (n<sub>1</sub> − k) and the (n<sub>2</sub> − k) × (n<sub>1</sub> − k) block R<sub>22</sub> has columns of degree greater than one.
- $n_1 k$  steps of Gaussian elimination (with partial pivoting) can be applied to  $R_{22}$  to complete the transformation of R to trapezoidal form.

#### Solving ill-conditioned problems

 Ill-conditioning is connected to the input data: a problem is ill-conditioned if small changes in the data can lead to large changes in the solution. Assume for the general linear system Ax = b that A and b are perturbed by ΔA and Δb, respectively, and the corresponding perturbation of the solution x is Δx, so that the perturbed problem

$$(A + \Delta A)(x + \Delta x) = b + \Delta b \tag{62}$$

has been solved. The perturbations in A and b may include both data uncertainty and algorithmic errors. Rearranging (62), we obtain

$$A\Delta x = \Delta b - \Delta A - \Delta A \Delta x.$$

• Premultiplying by  $A^{-1}$  and considering *any* norm  $\|.\|$  and the corresponding subordinate matrix norm yields

 $\|\Delta x\| \le \|A^{-1}\| \left( \|\Delta b\| + \|\Delta A\| \|x\| + \|\Delta A\| \|\Delta x\| \right).$ 

It follows that

 $(1 - ||A^{-1}|| ||\Delta A||) ||\Delta x|| \le ||A^{-1}|| (||\Delta b|| + ||\Delta A|| ||x||)$ 

and, provided  $\|A^{-1}\|\,\|\Delta A\|<1,$  this gives the following bound on the absolute error

$$\|\Delta x\| \le \frac{\|A^{-1}\|}{\|A^{-1}\| \|\Delta A\|} (\|\Delta b\| + \|\Delta A\| \|x\|).$$

• Dividing by ||x|| and using  $||b|| \le ||A|| \, ||x||$ , yields the relative error bound

$$\|\Delta x\|/\|x\| \le \frac{\kappa(A)}{1 - \kappa(A)} \|\Delta A\|/\|A\| \left( \|\Delta A\|/\|A\| + \|\Delta b\|/\|b\| \right).$$
(63)

### Solving ill-conditioned problems

Here

$$\kappa(A) = \|A\| \, \|A^{-1}\| \tag{64}$$

is the condition number of the matrix A.

• The inequality shows that the condition number is a relative error magnification factor. If we have a stable algorithm then a neighbouring problem has been solved, that is,

 $\|\Delta A\| / \|A\| + \|\Delta b\| / \|b\|$ 

is small. This ensures an accurate solution if  $\kappa(A)$  is small. A large condition number means that A is close to being singular ( $\kappa(A)$  tends to infinity as A tends to singularity).

- Observe that the condition number is very dependent on the scaling of *A*. Furthermore, κ(*A*) takes no account of the right-hand side vector *b* or the fact that small entries of *A* (including zeros) may be known within much smaller tolerances than larger entries.
- If the matrix norm is that induced by the Euclidean norm (that is! 1609

#### Iterative refinement

 Iterative refinement can be used to overcome matrix ill-conditioning and improve the accuracy of the computed solution. It may also be used after relaxed or static pivoting.

Algorithm (Iterative refinement of the computed solution of Ax = b)

**Input:** The vector *b* and matrix *A*. **Output:** A sequence of approximate solutions  $x^{(0)}, x^{(1)}, \ldots$ 

1: Solve  $Ax^{(0)} = b$ 2: for k = 0, 1, ... do 3: Compute  $r^{(k)} = b - Ax^{(k)}$ 4: Solve  $A \, \delta x^{(k)} = r^{(k)}$ 5:  $x^{(k+1)} = x^{(k)} + \delta x^{(k)}$ 6: end for  $\triangleright x^{(0)}$  is the initial computed solution

▷ Residual on iteration k
 ▷ Solve correction equation

### Iterative refinement

- Iterative refinement terminates when either the norm of the residual vector  $r^{(k)}$  is sufficiently close to zero that the corresponding correction  $\delta x^{(k)}$  is very small or the chosen maximum number of iterations is reached.
- A possible approach is to switch to using the computed factors as a preconditioner for a Krylov subspace solver.

#### Iterative refinement

- In fixed precision refinement, all computations use the same precision. In mixed precision iterative refinement, the most expensive parts of the computation (the LU factorization of *A* and solving the correction equation) are performed in single precision and the residual computation in double precision.
- Holding the factors in single precision substantially reduces the memory required and the amount of data movement. The use of half precision (16 bit) arithmetic is also a possibility, assuming it is considerably faster than single precision, with a proportional saving in energy consumption.

### Theorem

One step of single precision iterative refinement enough for obtaining componentwise relative backward error to the order of  $O(\epsilon)$  under weaker assumptions. Strong bound for the error norm using double precision iterative refinement.

### Scaling to reduce ill-conditioning

- An important way to decrease the condition number is by scaling *A* before the numerical factorization begins.
- Consider two nonsingular  $n \times n$  diagonal matrices  $S_r$  and  $S_c$ . Diagonal scaling of the system Ax = b transforms it to

$$S_r A S_c y = S_r b, \qquad y = S_c^{-1} x.$$
 (66)

#### Theorem

Let the matrix A be SPD and let  $D_A$  be the diagonal matrix with entries  $a_{ii}$  ( $1 \le i \le n$ ). Then for all diagonal matrices D with positive entries

$$\kappa(D_A^{-1/2} A D_A^{-1/2}) \le n z_{rmax} \, \kappa(D^{-1/2} A D^{-1/2}),$$

where  $nz_{rmax}$  is the maximum number of entries in a row of A.

### Equilibration scaling

- Finding an appropriate scaling is an open question, but a number of heuristics have been proposed.
- An obvious choice is to seek to balance entries of the scaled matrix  $S_r A S_c$  to have approximately equal absolute values. This is called (approximate) equilibration scaling.

Algorithm 10.7 presents an iterative procedure for computing such a scaling.

### Algorithm (Equilibration scaling in the infinity norm)

Input: The matrix A and convergence tolerance  $\delta > 0$ . Output: Diagonal scaling matrices  $S_r$  and  $S_c$ .

$$\begin{array}{l} 1: \ B^{(1)} = A, \ D^{(1)} = I, \ E^{(1)} = I \\ 2: \ \text{for } k = 1, 2, \dots \ \text{do} \\ 3: \quad Compute \ \|B_{i,1:n}^{(k)}\|_{\infty} \ \text{and} \ \|B_{1:n,i}^{(k)}\|_{\infty}, \ 1 \leq i \leq n \qquad \triangleright i \text{-th row and column of } B^{(k)} \\ 4: \quad \text{if } max_i \left\{ |1 - \|B_{i,1:n}^{(k)}\|_{\infty} || \right\} \leq \delta \ \text{ and } max_i \left\{ |1 - \|B_{1:n,i}^{(k)}\|_{\infty} || \right\} \leq \delta \ \text{ exit for loop} \\ 5: \quad R = diag \left( \sqrt{\|B_{i,1:n}^{(k)}\|_{\infty}} \right) \quad \text{and} \ C = diag \left( \sqrt{\|B_{1:n,i}^{(k)}\|_{\infty}} \right) \\ 6: \quad B^{(k+1)} = R^{-1}B^{(k)}C^{-1}, \quad D^{(k+1)} = D^{(k)}R^{-1}, \ E^{(k+1)} = E^{(k)}C^{-1} \\ 7: \ \text{end for} \\ 8: \ S_r = D^{(k+1)} \ \text{and} \ S_c = E^{(k+1)} \end{array}$$

- Our equilibration scaling preserves symmetry. In the nonsymmetric case, Algorithm yields the same results when applied to *A* and *A*<sup>T</sup> in the sense that the scaled matrix obtained for *A*<sup>T</sup> is the transpose of that for *A*.
- The infinity norm may be replaced by the 1-norm, resulting in a matrix whose row and column sums are exactly one (this is sometimes called a doubly stochastic matrix).
- It can be advantageous to combine the use of the infinity and one norms.

### Matching-based scalings

• Recall the problem of finding a permutation vector *q* that maximizes the product

$$\prod_{i=1}^{n} |a_{iq_i}|.$$

- The dual variables  $u_i$  and  $v_j$  computed by the MC64 algorithm that seeks to compute q can be used:
- Define the diagonal scaling matrices  $S_r$  and  $S_c$  to have entries

$$(S_r)_{ii} = \exp(u_i), \ (S_c)_{jj} = \exp(v_j - \log(\max_i |a_{ij}|)), \ 1 \le i, j \le n.$$

The entries of the scaled matrix  $S_r A S_c$  satisfy

$$|(S_r A S_c)_{ij}| \begin{cases} = 1, & \text{if } (i,j) \in \mathcal{M}, \\ \leq 1, & \text{otherwise}, \end{cases}$$

where  $\mathcal{M}$  is the maximum weighted matching computed by the MC64 algorithm.

Combining matching-based scalings and orderings

- The matching-based ordering and scaling can be used independently but they can also be combined.
- After scaling, if the matched entries are non symmetrically permuted onto the diagonal then because they are large, they provide good pivot candidates for an LU factorization.
- This approach is commonly used alongside static pivoting to obtain a factorization of a perturbed matrix, followed by iterative refinement to recover the solution to the original system.

• In the symmetric indefinite case, symmetry needs to maintained and so the objective is to symmetrically permute a large off-diagonal entry  $a_{ij}$  onto the subdiagonal to give a 2 × 2 block

 $\begin{pmatrix} a_{ii} & a_{ij} \\ a_{ij} & a_{jj} \end{pmatrix}$  that is potentially a good  $2 \times 2$  candidate pivot.

- Assume that a matching  $\mathcal{M}$  has been computed using the MC64 algorithm and let q be the corresponding permutation vector.
- Any diagonal entries that are in the matching are immediately considered as potential  $1 \times 1$  pivots and are held in a set  $M_1$ .
- A set M<sub>2</sub> of potential 2 × 2 pivots is then built by expressing q in terms of its component cycles. A cycle of length 1 corresponds to a<sub>ii</sub> in the matching. A cycle of length 2 corresponds to i and j, where a<sub>ij</sub> and a<sub>ji</sub> are both in the matching. k potential 2 × 2 pivots can be extracted from even cycles of length 2k or from odd cycles of length 2k + 1.
- In practice, most cycles in *q* are of length 1 or 2.

• Example for based on matching from MC64:

 $\mathcal{M} = \{(1,2), (2,5), (3,1), (4,4), (5,3)\}, \text{ which is nonsymmetric. } q \\ \text{has one cycle of length 4 } (1 \to 2 \to 5 \to 3 \to 1) \text{ and one of length} \\ \text{1, giving } \mathcal{M}_1 = \{(4,4)\} \text{ and } \mathcal{M}_2 = \{(1,2), (2,1), (3,5), (5,3)\}.$ 

Figure: An illustration of a symmetric matching for a symmetric indefinite matrix. On the left is the matching  $\mathcal{M}$  returned by the MC64 algorithm and in the centre is a symmetric matching  $\mathcal{M}_s$  obtained from  $\mathcal{M}$ . Entries in the matching are in blue. The pairs (i, j) = (1, 2) and (3, 5) are possible  $2 \times 2$  pivot candidates. On the right is the compressed matrix that results from combining rows and columns 1 and 2 and rows and columns 3 and 5.

# Outline



### Complexity of factorizations

- Starting by showing complexity of LU and Cholesky
- The complexity of the most critical steps in the factorization is highly dependent on the amount of fill-in, as can be seen from the following observation.

### Observation

The operations to perform the sparse LU factorization A = LU and the sparse Cholesky factorization  $A = LL^T$  are  $O(\sum_{j=1}^{n} |col_L\{j\}| |row_U\{j\}|)$  and  $O(\sum_{j=1}^{n} |col_L\{j\}|^2)$  respectively, where  $|row_U\{j\}|$  and  $|col_L\{j\}|$  are the number of off-diagonal entries in row j of U and column j of L, respectively.

- Key problem: minimizing the fill-in
- Our tools: permutations.
- The problem of finding a permutation to minimize fill-in is NP complete. Thus heuristics are used to determine orderings that limit the amount of fill-in; we refer to these as fill-reducing orderings.
- Frequently, this is done using the sparsity pattern  $S{A}$  alone,
- If the matrix is not SPD, additional permutations of *A* may be needed to make the matrix factorizable.

Two main classes of reorderings that work with  $\mathcal{S}{A}$  are commonly used.

Local orderings attempt to limit fill-in by repeated local decisions based on  $\mathcal{G}(A)$  (or a relevant quotient graph).

Global orderings consider the whole sparsity pattern of *A* and seek to find a permutation using a divide-and-conquer approach. Such methods are normally used in conjunction with a local fill-reducing ordering, as the latter generally works well for problems that are not really large.

- Assumed that A is irreducible. If not,
  - If S{A} is symmetric, the algorithms are applied to each component of G(A) independently and n is then the number of vertices in the component.
  - ▶ If *S*{*A*} is nonsymmetric, we assume that *A* is in block triangular form and the algorithms are used on the graph of each block on the diagonal.
- We also assume that *A* has no rows or columns that are (almost) dense. If so, such rows and/or columns should be treated independently.

- Historically, ordering the matrix *A* before using a direct solver to factorize it was generally cheap compared to the numerical factorization cost.
- It is not the case nowadays due to the development of the computational tools.
- In the symmetric case, the diagonal entries of A are required to be present in S{A} (thus zeros on the diagonal are included in the sparsity structure). The aim is to limit fill-in in the L factor of an LL<sup>T</sup> (or LDL<sup>T</sup>) factorization of A.
- Two greedy heuristics are the minimum degree (MD) criterion and the local minimum fill (MF) criterion.

### Minimum fill-in (MF) criterion

- One way to reduce fill-in is to use a local minimum fill-in (MF) criterion:
  - Select as the next variable in the ordering one that will introduce the least fill-in in the factor at that step.
- This is sometimes called the minimum deficiency approach.
- MF can produce good orderings, its cost is often considered to be prohibitive.
- An approximate variant (AMF).
# Basic minimum degree (MD) algorithm

- The best-known and most widely-used greedy heuristic for limiting fill-in.
- It seeks to find a permutation such that at each step of the factorization the number of entries in the corresponding column of *L* is approximately minimized.
- Less expensive to compute compared to that used by the minimum fill-in criterion.
- The MD algorithm can be implemented using  $\mathcal{G}(A)$  and it can predict the required factor storage without generating the structure of *L*.
- At step k, the number of off-diagonal nonzeros in a row or column of the active submatrix is the current degree of the corresponding vertex in the elimination graph G<sup>k</sup>. The algorithm selects a vertex of minimum current degree in G<sup>k</sup> and labels it v<sub>k</sub>, i.e., next for elimination.

#### Minimum degree algorithm

### Algorithm (Basic minimum degree (MD) algorithm)

*Input:* Graph G of a symmetrically structured matrix. *Output:* A permutation vector p that defines a new labelling of the vertices of G.

1: Set  $\mathcal{G}^1 = \mathcal{G}$  and compute the degree  $deg_{\mathcal{G}^1}(u)$  of all  $u \in \mathcal{V}(\mathcal{G}^1)$ 

2: for k = 1 : n - 1 do

- 3: Compute  $mdeg = \min\{deg_{\mathcal{G}^k}(u) | u \in \mathcal{V}(\mathcal{G}^k)\} \triangleright mdeg$  is the current minimum degree
- 4: Choose  $v_k \in \mathcal{V}(\mathcal{G}^k)$  such that  $deg_{\mathcal{G}^k}(v_k) = mdeg$
- 5:  $p(k) = v_k$   $\triangleright v_k$  is the next vertex in the elimination order
- 6: Construct  $\mathcal{G}^{k+1}$  and update the current degrees of its vertices
- 7: end for

8:  $p(n) = v_n$  where  $v_n$  is the only vertex in  $\mathcal{G}^n$ 

#### Minimum degree algorithm

- The set of vertices adjacent to  $v_k$  in  $\mathcal{G}(A)$  is  $\mathcal{R}each(v_k, \mathcal{V}_k)$ , where  $\mathcal{V}_k$  is the set of k 1 vertices that have already been eliminated.
- If u ∈ Reach(v<sub>k</sub>, V<sub>k</sub>), u ≠ v<sub>k</sub>, then its updated current degree is |Reach(u, V<sub>k+1</sub>)|, where V<sub>k+1</sub> = V<sub>k</sub> ∪ v<sub>k</sub>.
- A tie-breaking strategy is needed when there is more than one vertex of current minimum degree.
- It is possible to construct artificial matrices showing that some systematic tie-breaking choices can lead to a large amount of fill-in but such behaviour is not typical.

#### Minimum degree algorithm



Figure: An illustration of three steps of the MD algorithm. The original graph  $\mathcal{G}$  and the elimination graphs  $\mathcal{G}^2$ ,  $\mathcal{G}^3$  and  $\mathcal{G}^4$  that result from eliminating vertex 2, then vertex 3 and then vertex 1. The red dashed lines denote fill edges.

## Minimum degree algorithm

- The construction of each elimination graph  $\mathcal{G}^{k+1}$  is central to the implementation.
- Because eliminating a vertex potentially creates fill-in, an efficient representation is needed.
- Moreover, recalculating the current degrees is time consuming.
- Using supervariables is a must.
- $\mathcal{G}_v$  denotes the elimination graph obtained from  $\mathcal{G}$  when vertex  $v \in \mathcal{V}(\mathcal{G})$  is eliminated.

# Indistinguishability (reminder)

# Definition

# Two different vertices $\mathbf{u}$ and $\mathbf{v}$ of G are called **indistinguishable** if

$$Adj_G(u) \cup \{u\} = Adj_G(v) \cup \{v\}.$$
(67)



#### Theorem

Let u and w be indistinguishable vertices in  $\mathcal{G}$ . If  $v \in \mathcal{V}(\mathcal{G})$  with  $v \neq u, w$ , then u and w are indistinguishable in  $\mathcal{G}_v$ .

# Proof.

Two cases must be considered. First, let  $u \notin adj_{\mathcal{G}}\{v\}$ . Then  $w \notin adj_{\mathcal{G}}\{v\}$  and if v is eliminated, the adjacency sets of u and w are unchanged and they remain indistinguishable in the resulting elimination graph  $\mathcal{G}_v$ . Second, let  $u, w \in adj_{\mathcal{G}}\{v\}$ . When v is eliminated, because u and w are indistinguishable in  $\mathcal{G}$ , their adjacency sets in  $\mathcal{G}_v$  will be modified in the same way, by adding the entries of  $adj_{\mathcal{G}}\{v\}$  that are not already in  $adj_{\mathcal{G}}\{u\}$  and  $adj_{\mathcal{G}}\{w\}$ . Consequently, u and w are indistinguishable in  $\mathcal{G}_v$ .  $\Box$ 

Figure demonstrates the two cases in the proof of Theorem above. Here, u and w are indistinguishable vertices in  $\mathcal{G}$ . Setting  $v \equiv v'$ illustrates  $u \notin adj_{\mathcal{G}}\{v\}$ . If v' is eliminated then the adjacency sets of uand w are clearly unchanged. Setting  $v \equiv v''$  illustrates  $u, w \in adj_{\mathcal{G}}\{v\}$ . In this case, if v'' is eliminated then vertices s and t are added to both  $adj_{\mathcal{G}}\{u\}$  and  $adj_{\mathcal{G}}\{w\}$ .



Figure: An example to illustrate the Theorem. u and w are indistinguishable vertices in  $\mathcal{G}$ ;  $adj_{\mathcal{G}}\{u\} = \{r, w, v''\}$  and  $adj_{\mathcal{G}}\{w\} = \{r, u, v''\}$ .

#### Theorem

Let u and w be indistinguishable vertices in  $\mathcal{G}$ . If w is of minimum degree in  $\mathcal{G}$  then u is of minimum degree in  $\mathcal{G}_w$ .

# Proof.

Let  $deg_{\mathcal{G}}(w) = mdeg$ . Then  $deg_{\mathcal{G}}(u) = mdeg$ . Indistinguishable vertices are always neighbours. Eliminating w gives  $deg_{\mathcal{G}_w}(u) = mdeg - 1$  because w is removed from the adjacency set of uand there is no neighbour of u in  $\mathcal{G}_w$  that was not its neighbour in  $\mathcal{G}$ . If  $x \neq w$  and  $x \in adj_{\mathcal{G}}\{u\}$ then the number of neighbours of x in  $\mathcal{G}_w$  is at least mdeg - 1. Otherwise, if  $x \notin adj_{\mathcal{G}}\{u\}$  then its adjacency set in  $\mathcal{G}_w$  is the same as in  $\mathcal{G}$  and is of the size at least mdeg. The result follows.  $\Box$ 

#### Indistinguishability

Theorem above is illustrated in Figure.



Figure: An illustration of Theorem. Vertices u and w are of minimum degree (with degree mdeg = 3) and are indistinguishable in  $\mathcal{G}$ . After elimination of w, the current degree of u is mdeg - 1 and the current degree of each of the other vertices is at most mdeg - 1. Therefore, u is of current minimum degree in  $\mathcal{G}_w$ . Note that vertices r and v are also of minimum degree and indistinguishable in  $\mathcal{G}$ ; they are not neighbours of w and their degrees do not change when w is eliminated.

- The results can be extended to more than two indistinguishable vertices, which allows indistinguishable vertices to be selected one after another in the MD ordering.
- This is referred to as mass elimination.
- Treating indistinguishable vertices as a single supervariable cuts the number of vertices and edges in the elimination graphs, which reduces the work needed for degree updates.
- The external degree of a vertex is the number of vertices adjacent to it that are not indistinguishable from it. Using this leads to algorithmic efficiency.

- A concept that is closely related to that of indistinguishable vertices is degree outmatching.
- This avoids computing the degrees of vertices that are known not to be of current minimum degree.
- Vertex w is said to be outmatched by vertex u if

 $adj_{\mathcal{G}}\{u\} \cup \{u\} \subseteq adj_{\mathcal{G}}\{w\} \cup \{w\}.$ 

• It follows that  $deg_{\mathcal{G}}(u) \leq deg_{\mathcal{G}}(w)$ .



Figure: An example G in which vertex w is outmatched by vertex u. v' is not a neighbour of u or w; vertex v'' is a neighbour of both u and w; v''' is a neighbour of w but not of u.

 Importantly, degree outmatching is preserved when vertex v ∈ G of minimum degree is eliminated, as stated in the following result.

#### Theorem

In the graph  $\mathcal{G}$  let vertex w be outmatched by vertex u and vertex v  $(v \neq u, w)$  be of minimum degree. Then w is outmatched in  $\mathcal{G}_v$  by u.

### Proof.

Three cases must be considered. First, if  $u \notin adj_{\mathcal{G}}\{v\}$  and  $w \notin adj_{\mathcal{G}}\{v\}$  then the adjacency sets of u and w in  $\mathcal{G}_v$  are the same as in  $\mathcal{G}$ . Second, if v is a neighbour of both u and w in  $\mathcal{G}$  then any neighbours of v that were not neighbours of u and w are added to their adjacency sets in  $\mathcal{G}_v$ . Third, if  $u \notin adj_{\mathcal{G}}\{v\}$  and  $w \in adj_{\mathcal{G}}\{v\}$  then the adjacency set of u in  $\mathcal{G}_v$  is the same as in  $\mathcal{G}$  but any neighbours of v that were not neighbours of w are added to the adjacency set of w in  $\mathcal{G}_v$ . In all three cases, w is still outmatched by u in  $\mathcal{G}_v$ .  $\Box$ 

- The three possible cases for v in the proof of Theorem are illustrated in Figure above by setting  $v \equiv v'$ , v'' and v''', respectively.
- If *w* is outmatched by *u* then it is not necessary to consider *w* as a candidate for elimination and
- all updates to the data structures related to *w* can be postponed until *u* has been eliminated.

- From Parter's rule, if vertex v is selected at step k then the elimination matrix that corresponds to  $\mathcal{G}^{k+1}$  contains a dense submatrix of size equal to the number of off-diagonal entries in row and column v in the matrix that corresponds to  $\mathcal{G}^k$ .
- For large matrices, creating and explicitly storing the edges in the sequence of elimination graphs is impractical and a more compact and efficient representation is needed.
- Each elimination graph can be interpreted as a collection of cliques, including the original graph G, which can be regarded as having |E| cliques, each consisting of two vertices (or, equivalently, an edge).

- Let  $\{\mathcal{V}_1, \mathcal{V}_2, \dots, \mathcal{V}_q\}$  be the set of cliques for the current graph and let v be a vertex of current minimum degree that is selected for elimination. Let  $\{\mathcal{V}_{s_1}, \mathcal{V}_{s_2}, \dots, \mathcal{V}_{s_t}\}$  be the subset of cliques to which v belongs. Two steps are then required.
  - **O** Remove the cliques  $\{\mathcal{V}_{s_1}, \mathcal{V}_{s_2}, \ldots, \mathcal{V}_{s_t}\}$  from  $\{\mathcal{V}_1, \mathcal{V}_2, \ldots, \mathcal{V}_q\}$ .
  - 2 Add the new clique  $\mathcal{V}_v = \{\mathcal{V}_{s_1} \cup \ldots \cup \mathcal{V}_{s_t}\} \setminus \{v\}$  into the set of cliques.

Hence

$$deg_{\mathcal{G}}(v) = |\mathcal{V}_v| < \sum_{i=1}^t |\mathcal{V}_{s_i}|,$$

and because  $\{\mathcal{V}_{s_1}, \mathcal{V}_{s_2}, \dots, \mathcal{V}_{s_t}\}$  can now be discarded, the storage required for the representation of the sequence of elimination graphs never exceeds that needed for  $\mathcal{G}(A)$ .

 The storage to compute an MD ordering is therefore known beforehand in spite of the rather dynamic nature of the elimination process.

- The index of the eliminated vertex can be used as the index of the new clique. This is called an element or enode (the terminology comes from finite-element methods), to distinguish it from an uneliminated vertex, which is termed an snode.
- A sequence of special quotient graphs  $\mathcal{G}^{[1]} = \mathcal{G}(A), \mathcal{G}^{[2]}, \dots, \mathcal{G}^{[n]}$  with the two types of vertices is generated in place of the elimination graphs.
- Each  $\mathcal{G}^{[k]}$  has n vertices that satisfy

$$\mathcal{V}(\mathcal{G}) = \mathcal{V}_{snodes} \cup \mathcal{V}_{enodes}, \qquad \mathcal{V}_{snodes} \cap \mathcal{V}_{enodes} = \emptyset,$$

where  $\mathcal{V}_{snodes}$  and  $\mathcal{V}_{enodes}$  are the sets of snodes and enodes, respectively.

 When v is eliminated, any enodes adjacent to it are no longer required to represent the sparsity pattern of the corresponding active submatrix and so they can be removed. This is called element absorption.

- Working with the special quotient graphs can be demonstrated by considering the computation of the vertex degrees.
- To compute the degree of an uneliminated vertex, the set of neighbouring snodes is counted.
- Then, if a neighbour of one of these snodes is an enode, its neighbours are also counted (avoiding double counting).
- More formally, if v ∈ V<sub>snodes</sub> then the adjacency set of v is the union of its neighbours in V<sub>snodes</sub> and the vertices reachable from v via its neighbours in V<sub>enodes</sub>.
- In this way, vertex degrees are computed by considering fill-paths
- Amalgamation improves this strategy: mass elimination.

# Cliques and quotient graphs: mass elimination model

# Definition

**Mass elimination graph**  $\Gamma$  of the graph G = (V, E) is a ordered triple  $(S, \mathcal{E}, E)$ , where  $S \cup \mathcal{E} = V, S \cap \mathcal{E} = \emptyset$  and  $E \subseteq \begin{pmatrix} S \\ 2 \end{pmatrix} \cup \begin{pmatrix} \mathcal{E}(\Gamma) \\ 2 \end{pmatrix}$  are its edges.

- Edge set  $\mathcal{E}$  captures eliminated vertices.
- Edge set *S* captures non-eliminated vertices.
- Neighbors of non-eliminated vertices are found as reachability sets.
- Search through the reachability sets can be pruned: → approximate minimum degree (AMD) algorithm.

#### Cliques and quotient graphs

- After elimination of 1, a new edge is added, getting a clique. The elimination of 2 creates no additional fill and  $\mathcal{G}^3$  represents the sparsity structure of the corresponding active submatrix  $A^{(3)}$ .
- Then, 1 is an enode, the fill edge is represented implicitly. After the second step, the enodes 1 and 2 can be amalgamated as well as snodes 3 and 4 being indistinguishable.



Figure: The top line shows  $\mathcal{G} = \mathcal{G}^1$ ,  $\mathcal{G}^2$  and  $\mathcal{G}^3$ . The bottom line shows the quotient graphs  $\mathcal{G}^{[2]}$  and  $\mathcal{G}^{[3]}$  after the first and second elimination steps. 440/609

# Multiple minimum degree (MMD) algorithm

- The multiple minimum degree (MMD) algorithm aims to improve efficiency by processing several independent vertices that are each of minimum current degree together in the same step, before the degree updates are performed.
- At each outer loop, t ≥ 1 denotes the number of vertices of minimum current degree that are mutually non-adjacent and so can be put into the elimination order one after another.
- An example follows.



Figure: The red (corner) vertices of G are each of degree 2 and are ordered consecutively during the first step of Algorithm 11.2.

# Multiple minimum degree (MMD) algorithm

# Algorithm (Basic multiple minimum degree (MMD) algorithm)

**Input:** Graph G of a symmetrically structured matrix. **Output:** A permutation vector p that defines a new labelling of the vertices of G.

1: Set k = 1,  $\mathcal{G}^1 = \mathcal{G}$  and compute the degree  $deg_{\mathcal{G}^1}(u)$  of all  $u \in \mathcal{V}(\mathcal{G}^1)$ 

- 2: while  $k \leq n$  do
- 3: Compute  $mdeg = \min\{deg_{\mathcal{G}^k}(u) \mid u \in \mathcal{V}(\mathcal{G}^k)\}$
- 4: Find all mutually non-adjacent  $v_j \in \mathcal{V}(\mathcal{G}^k)$ , j = 1, ..., t with  $deg_{\mathcal{G}^k}(v_j) = mdeg$
- 5: for j = 1 : t do 6:  $p(k) = v_j$

 $\triangleright$  Vertex  $v_j$  is the next vertex in the elimination order

- 7: k = k + 1
- 8: end for
- 9: if k < n then
- 10: Construct  $\mathcal{G}^{k+1}$  and update the current degrees of its vertices
- 11: end if
- 12: end while

# MD and MMD: complexity

- The complexity of the MD and MMD algorithms is  $O(nz(A)n^2)$  but because for MMD the outer loop of the algorithm update is performed fewer times, it can be significantly faster then MD.
- MMD orderings can also lead to less fill-in, possibly a consequence of introducing some kind of regularity into the ordering sequence.

#### Approximate minimum degree (AMD) algorithm

- The idea behind the widely-used approximate minimum degree (AMD) algorithm is to inexpensively compute an upper bound on a vertex degree in place of the degree, and to use this bound as an approximation to the external degree.
- The quality of the orderings obtained using the AMD algorithm are competitive with those computed using the MD algorithm and can surpass them.
- The complexity of AMD is O(nz(A)n) and, in practice, its runtime is typically significantly less than that of the MD and MMD approaches.

#### Minimizing the bandwidth and profile

- An alternative way of reducing the fill-in locally is to add another criterion to the relabelling of the vertices, such as restricting the nonzeros of the permuted matrix to specific positions.
- The most popular approach is to force them to lie close to the main diagonal.
- All fill-in then takes place between the first entry of a row and the diagonal or between the first entry of a column and the diagonal.
- This allows straightforward implementations of Gaussian elimination that employ static data structures.
- Here we again consider symmetric  $S{A}$ ; generalizations are possible.

#### The band and envelope of a matrix

Denote:

$$\eta_i(A) = \min\{j \mid 1 \le i \le j, \text{ with } a_{ij} \ne 0\}, \quad 1 \le i \le n,$$
 (68)

that is,  $\eta_i(A)$  is the column index of the first entry in the *i*-th row of *A*. Define

$$\beta_i(A) = i - \eta_i(A), \quad 1 \le i \le n.$$

The semibandwidth of A is

$$\max\{\beta_i(A) \mid 1 \le i \le n\},\$$

and the bandwidth is

$$\beta(A) = 2 * \max\{\beta_i(A) \mid 1 \le i \le n\} + 1.$$

The band of A is the following set of index pairs in A

$$band(A) = \{(i, j) \mid 0 < i - j \le \beta(A)\}.$$

## The band and envelope of a matrix

• The envelope is the set of index pairs that lie between the first entry in each row and the diagonal:

 $env(A) = \{(i, j) \mid 0 < i - j \le \beta_i(A)\}.$ 

- Note that the band and envelope of a sparse symmetrically structured matrix *A* include only entries of the strict lower triangular part of *A*.
- The envelope is easily visualized: picture the lower triangular part of *A*, and remove the diagonal and the leading zero entries in each row. The remaining entries (whether nonzero or zero) comprise the envelope of *A*.
- The profile of *A* is defined to be the number of entries in the envelope (the envelope size) plus *n*.

#### The band and envelope of a matrix: shape pushers



#### The band and envelope of a matrix

An illustrative example: Here  $\eta_1(A) = 1$ ,  $\beta_1(A) = 0$ ,  $\eta_2(A) = 1$ ,  $\beta_2(A) = 1$ ,  $\eta_3(A) = 2$ ,  $\beta_3(A) = 1$ , and so on.



Figure: Illustration of the band and envelope of a matrix A whose sparsity pattern is on the left. In the centre, the positions of band(A) are circled and on the right, the positions of env(A) are circled. The bandwidth is 5 and the envelope size and the profile are 7 and 14, respectively.

#### The band and envelope of a matrix

#### Static structures!

#### Theorem

If L is the Cholesky factor of A then

env(A) = env(L).

# Proof.

The proof uses mathematical induction on the principal leading submatrices of A of order k. The result is clearly true for k = 1 and k = 2. Assume it holds for  $2 \le k < n$  and consider the block factorization

$$\begin{pmatrix} A_{1:k,1:k} & u_{1:k} \\ u_{1:k}^T & \alpha \end{pmatrix} = \begin{pmatrix} L_{1:k,1:k} & 0 \\ v_{1:k}^T & \beta \end{pmatrix} \begin{pmatrix} L_{1:k,1:k}^T & v_{1:k} \\ 0 & \beta \end{pmatrix},$$

where  $\alpha$  and  $\beta$  are scalars. Equating the left and right sides,  $L_{1:k,1:k}v_{1:k} = u_{1:k}$ . Because  $u_j = 0$  for  $j < \eta_{k+1}(A)$  and  $u_{\eta_{k+1}} \neq 0$ , it follows that  $v_j = 0$  for  $j < \eta_{k+1}(A)$  and  $v_{\eta_{k+1}} \neq 0$ . This proves the induction step.  $\Box$ 

A straightforward corollary is that band(A) = band(L).

### The band and envelope of a matrix

- Finding a permutation *P* to minimize the band or profile of *PAP*<sup>T</sup> is again combinatorially hard and again heuristics are used to efficiently find an acceptable *P*.
- The popular Cuthill McKee (CM) approach chooses a suitable starting vertex *s* and labels it 1.
- Then, for i = 1, 2, ..., n 1, all vertices adjacent to vertex i that are still unlabelled are labelled successively in order of increasing degree, as described in Algorithm below.
- A very important variation is the Reverse Cuthill McKee (RCM) algorithm, which incorporates a final step in which the CM ordering is reversed.

#### Level-based orderings

# Algorithm (CM and RCM algorithms for band and profile reduction)

*Input:* Graph  $\mathcal{G}$  of a symmetrically structured matrix and a starting vertex *s*. *Output:* Permutation vectors  $p_{cm}$  and  $p_{rcm}$  that define new labellings of the vertices of  $\mathcal{G}(A)$ .

1: label(1:n) = false2:  $Compute adj_{\mathcal{G}}\{u\}$  and  $deg_{\mathcal{G}}(u)$  for all  $u \in \mathcal{V}(\mathcal{G})$ 3:  $k = 1, v_1 = s, p_{cm}(1) = v_1, label(v_1) = true$ 4: for i = 1: n - 1 do 5: for  $w \in adj_{\mathcal{G}}\{v_i\}$  with label(w) = false in order of increasing degree do 6:  $k = k + 1, v_k = w, p_{cm}(k) = v_k, label(v_k) = true$ 7: end for 8: end for 9: For the RCM ordering,  $p_{rcm}(i) = p_{cm}(n - i + 1), i = 1, 2, ..., n$ .

• The CM- and RCM-permuted matrices have the same bandwidth but the latter can decrease the envelope.



	1	<b>2</b>	3	4	5	6	7			3	7	1	5	<b>2</b>	4	6			6	4	<b>2</b>	5	1	7	3	
1	/*	*		*	*		* \		3	/*	*							6	/*		*					
2	*	*		*		*			7	*	*	*						4	(	*	*	*	*			١
3			*				*		1		*	*	*	*	*			<b>2</b>	*	*	*		*			I
4	*	*		*	*			,	5	l		*	*		*		,	5		*		*	*			l
5	*			*	*				<b>2</b>			*		*	*	*		1		*	*	*	*	*		I
6		*				*			4			*	*	*	*			7					*	*	*	I
7	(*		*				*/		6					*		*/		3						*	*/	/

#### Level-based orderings

The importance of the CM and RCM orderings is also expressed by:

## Theorem

Let *A* be symmetrically structured and irreducible. If *P* corresponds to the CM labelling obtained from Algorithm and *L* is the Cholesky factor of  $P^TAP$  then env(L) is full, that is, all entries of the envelope are nonzero.

- The full envelope of the Cholesky factor of the permuted matrix implies cache efficiency when performing the triangular solves once the factorization is complete.
- A crucial difference between profile reduction ordering algorithms and minimum degree strategies is that the former is based solely on *G*: the costly construction of quotient graphs is not needed.
- However, unless the profile after reordering is very small, there can be significantly more fill-in in the factor.
#### Level-based orderings

- Key to the success is the choice of the starting vertex *s*.
- A good candidate is a vertex for which the maximum distance between it and some other vertex in *G* is large.
- Formally, the eccentricity  $\epsilon(u)$  of the vertex u in the connected undirected graph  $\mathcal{G}$  is defined to be

$$\epsilon(u) = \max_{v \in \mathcal{V}} d(u, v),$$

where d(u, v) is the distance between the vertices u and v (the length of the shortest path between these vertices).

- The maximum eccentricity taken over all the vertices is the diameter of G (that is, the maximum distance between any pair of vertices). The endpoints of a diameter (also termed peripheral vertices) provide good starting vertices.
- The complexity of finding a diameter is  $O(n^3)$ : approximation (pseudo-preferal vertices) are needed.

#### Level-based orderings

- A heuristic algorithm is used to find pseudo-peripheral vertices. A commonly-used approach is based on level sets. A level structure rooted at a vertex *r* is defined as the partitioning of V into disjoint levels L<sub>1</sub>(*r*), L<sub>2</sub>(*r*),..., L<sub>h</sub>(*r*) such that

   (i) L<sub>1</sub>(*r*) = {*r*} and
   (ii) for 1 < *i* ≤ *h*, L<sub>i</sub>(*r*) is the set of all vertices that are adjacent to vertices in L<sub>i-1</sub>(*r*) but are not in L<sub>1</sub>(*r*), L<sub>2</sub>(*r*),..., L<sub>i-1</sub>(*r*).
- The level structure rooted at r may be expressed as the set  $\mathcal{L}(r) = \{\mathcal{L}_1(r), \mathcal{L}_2(r), \dots, \mathcal{L}_h(r)\}$ , where h is the total number of levels and is termed the depth.
- The level sets can be found using a breadth-first search that starts at the root *r*.

# Reorderings

#### Level-based orderings: GPS

# Algorithm (GPS algorithm to find pseudo-peripheral vertices)

1:	Construct $\mathcal{L}(r)$ and set $flag$	= false
2:	while $flag = false$ do	
3:	flag = true	
4:	for $i=1: \mathcal{L}(r) $ do	
5:	$w_i \in \mathcal{L}(r)$	$\triangleright$ Select vertex $w_i$ from last level set
6:	if $flag = true$ then	
7:	Construct $\mathcal{L}(w_i)$	
8:	if $depth(\mathcal{L}(w_i)) >$	$depth(\mathcal{L}(r))$ then
9:	flag = false	$\triangleright$ Flag that $w_i$ will be used as new initial vertex
10	cend if	
11.	end if	
12	end for	
13	: if $flag = true$ then	
14	: $s = r$ and $t = w_i$	$\triangleright s$ has been chosen; while loop will terminate algorithm
15	: else	
16	$r = w_i$	
17	end if	
18	end while	458/609

#### Level-based orderings

A simple example: starting with r = 2, after two passes through the while loop, the GPS algorithm returns s = 8 and t = 1 as pseudo-peripheral vertices.



Figure: An example to illustrate Algorithm 11.4 for finding pseudo-peripheral vertices. With root vertex r = 2, the first level set structure is  $\mathcal{L}(2) = \{\{2\}, \{1,3\}, \{4,5,7\}, \{6,8\}\}$ . Setting r = 8 at Step 16, the second level set structure is  $\mathcal{L}(8) = \{\{8\}, \{4,7\}, \{3,6\}, \{2,5\}, \{1\}\}$  and the algorithm terminates with s = 8 and t = 1.

#### Spectral orderings

• The spectral algorithm associates a positive semidefinite Laplacian matrix *L<sub>p</sub>* with the symmetric matrix *A* as follows:

$$(L_p)_{ij} = egin{cases} -1 & ext{if } i 
eq j ext{ and } a_{ij} 
eq 0, \\ deg_{\mathcal{G}}(i) & ext{if } i = j, \\ 0 & ext{otherwise.} \end{cases}$$

- An eigenvector corresponding to the smallest positive eigenvalue of the Laplacian matrix is called a Fiedler vector. If  $\mathcal{G}$  is connected,  $L_p$  is irreducible and the second smallest eigenvalue is positive.
- The vertices of  $\mathcal{G}$  are ordered by sorting the entries of the Fiedler vector into monotonic order. Applying the permutation symmetrically to A yields the spectral ordering.

#### Spectral orderings

• The use of the Fiedler vector for reordering *A* comes from considering the matrix envelope with the size

$$|env(A)| = \sum_{i=1}^{n} \beta_i = \sum_{i=1}^{n} \max_{\substack{k < i \\ (k,i) \in \mathcal{G}}} (i-k).$$

The asymptotic upper bound on the operation count for the factorization based on env(A) is

$$work_{env} = \sum_{i=1}^{n} \beta_i^2 = \sum_{i=1}^{n} \max_{\substack{k < i \\ (k,i) \in \mathcal{G}}} (i-k)^2.$$

Ordering the vertices using the Fiedler vector is closely related to minimizing  $weight_{env}$  over all possible vertex reorderings, where

$$weight_{env} = \sum_{i=1}^{n} \sum_{\substack{k < i \\ (k,i) \in \mathcal{G}}} (i-k)^2.$$

## Spectral orderings

- Thus, while minimizing the profile and envelope is related to the infinity norm, minimizing *weight*<sub>env</sub> is related to the Euclidean norm of the distance between graph vertices.
- Although computing the Fiedler vector can be computationally expensive it does have the advantage of easy vectorization and parallelization and the resulting ordering can give small profiles and low operation counts.

## Local fill-reducing orderings for nonsymmetric $S{A}$

- If  $S{A}$  is nonsymmetric then an often-used strategy is to apply the minimum degree algorithm (or one of its variants) or a band or profile-reducing ordering to the undirected graph  $\mathcal{G}(A + A^T)$ .
- This can work well if the symmetry index *s*(*A*) is close to 1. But if *A* is highly nonsymmetric, another approach is required.
- Markowitz pivoting generalizes the MD algorithm by choosing the pivot entry based on vertex degrees computed directly from the nonsymmetric  $S{A}$ ; the result is a nonsymmetric permutation.

# Reorderings

## Markowitz pivoting

At step k of the LU factorization, consider the (n − k + 1) × (n − k + 1) active submatrix, that is, the Schur complement S<sup>(k)</sup>. Let nz(row<sub>i</sub>) and nz(col<sub>j</sub>) denote the number of entries in row i and column j of S<sup>(k)</sup> (1 ≤ i, j ≤ n − k + 1). Markowitz pivoting selects as the k-th pivot the entry of S<sup>(k)</sup> that minimizes the Markowitz count given by the product

$$(nz(row_i) - 1)(nz(col_j) - 1).$$

- It can be described using a sequence of bipartite graphs of the active submatrices but here we use a matrix-based description that permutes *A* on the fly.
- Markowitz pivoting is generally incorporated into the numerical factorization phase of an LU solver, rather than being used to derive an initial reordering of *A*.

## Markowitz pivoting

Implementation of the algorithm requires access to the rows and the columns of the matrix.

# Algorithm (Markowitz pivoting)

*Input:* Matrix A with a nonsymmetric sparsity pattern. **Output:** A' = PAQ, where P and Q are permutation matrices chosen to limit fill in.

1: Set  $S^{(1)} = A$  and A' = A

2: for k = 1 : n - 1 do

- 3: Compute  $nz(row_i)$  and  $nz(col_j)$   $(1 \le i, j \le n k + 1)$
- 4: Find an entry  $s_{ij}^{(k)}$  of  $S^{(k)}$  that minimizes  $(nz(row_i) 1)(nz(col_j) 1))$
- 5: Permute the rows and columns so that  $s_{ij}^{(k)}$  is the (1,1) entry of the permuted  $S^{(k)}$
- 6: Compute Schur complement  $S^{(k+1)}$  of the permuted  $S^{(k)}$  with respect to its (1,1) entry 7: end for

#### Markowitz pivoting

Example: the first pivot is  $a_{24}$  with Markowitz count 1; it does not cause fill-in. The second pivot has Markowitz count 2 in  $S^{(2)}$ ; it results in one filled entry.



Figure: Illustration of Markowitz pivoting. The first and second pivots are circled. The sparsity pattern of  $A = S^{(1)}$  is on the left. In the centre is the sparsity pattern after permuting the pivot in position (2,4) to the (1,1) position of  $S^{(1)}$ . There is no fill-in after the first factorization step. On the right is the sparsity pattern after selecting the second pivot that has the original position (4,2) and permuting it to the (1,1) position of  $S^{(2)}$ . The resulting filled entry is denoted by f. Note that the nonsymmetric permutations transform the originally irreducible matrix into a reducible one.

# Local fill-reducing orderings for nonsymmetric $\mathcal{S}\{A\}$

- Markowitz pivoting as described here only considers the sparsity of *A* and the subsequent Schur complements.
- In practice, small pivots should be avoided.
- Practical implementations: relaxations: restriction to a limited number of rows and columns.
- Dynamic sparse formats like DS needed.

- Nested dissection (ND) is the most important and widely-used global ordering strategy for direct methods when S{A} is symmetric; it is particularly effective for ordering very large matrices.
- Identifying a small set of vertices V<sub>S</sub> (known as a vertex separator)
- If removed separates the graph into two disjoint subgraphs described by the vertex subsets B and W.
- The rows and columns belonging to B are labelled first, then those belonging to W and finally those in V<sub>S</sub>. The reordered matrix has the form

$$\begin{pmatrix} A_{\mathcal{B},\mathcal{B}} & 0 & A_{\mathcal{B},\mathcal{V}_{\mathcal{S}}} \\ 0 & A_{\mathcal{W},\mathcal{W}} & A_{\mathcal{W},\mathcal{V}_{\mathcal{S}}} \\ A_{\mathcal{B},\mathcal{V}_{\mathcal{S}}}^T & A_{\mathcal{W},\mathcal{V}_{\mathcal{S}}}^T & A_{\mathcal{V}_{\mathcal{S}},\mathcal{V}_{\mathcal{S}}} \end{pmatrix}.$$
 (69)

Definition

**Vertex separator** of an undirected G = (V, E) is subset *S* of its vertices such that the subgraph induced by  $V \setminus S$  has more components than *G*.

Induced reordering

$$A = \begin{pmatrix} A_{11} & 0 & A_{31}^T \\ 0 & A_{22} & A_{32}^T \\ A_{31} & A_{32} & A_{33} \end{pmatrix}$$
(70)

# Initial reordering

#### Global nested dissection orderings





# Reorderings

#### Global nested dissection orderings





472/609

( *	*	*	*									
*	*		*	*								
*		*	*							*	*	
*	*	*	*	*							*	
	*		*	*							*	*
					*	*		*		*		
					*	*	*	*	*	*	*	
						*	*		*		*	*
					*	*		*	*			
						*	*	*	*			
		*			*	*				*	*	
		*	*	*		*	*			*	*	*
				*			*				*	*

- Provided the variables are eliminated in the permuted order, no fill occurs within the zero off-diagonal blocks.
- If |V<sub>S</sub>| is small and |B| and |W| are similar, these zero blocks account for approximately half the possible entries in the matrix.
- The reordering can be applied recursively to the submatrices  $A_{\mathcal{B},\mathcal{B}}$  and  $A_{\mathcal{W},\mathcal{W}}$  until the vertex subsets are of size less than some prescribed threshold.
- Combinations with local reorderings.

# Algorithm (Nested dissection algorithm)

*Input:* Graph  $\mathcal{G}$  of a symmetrically structured matrix A and a partitioning algorithm PartitionAlg. *Output:* A permutation vector p that defines a new labelling of the vertices of  $\mathcal{G}$ .

1: recursive function (p = nested\_dissection(A, PartitionAlg)) 2: if dissection has terminated then Vertex subsets are smaller than some threshold 3:  $p = AMD(\mathcal{V}, \mathcal{E})$ Compute an AMD ordering 4: else 5: Use PartitionAlg( $\mathcal{V}, \mathcal{E}$ ) to obtain the vertex partitioning ( $\mathcal{B}, \mathcal{W}, \mathcal{V}_S$ ) 6:  $p_{\mathcal{B}} = nested \ dissection(A_{\mathcal{B},\mathcal{B}}, PartitionAlg)$ 7:  $p_{W} = nested_dissection(A_{W,W}, PartitionAlg)$ 8:  $p_{\mathcal{V}_S}$  is an ordering of  $\mathcal{V}_S$ Set  $p = \begin{pmatrix} p_{\mathcal{B}} \\ p_{\mathcal{W}} \\ p_{\mathcal{W}} \end{pmatrix}$ 9: 10: end if 11: end recursive function

#### Bordered forms

 The matrix (69) is an example of a doubly bordered block diagonal (DBBD) form. More generally, a matrix is said in DBBD form if it has the block structure

$$A_{DB} = \begin{pmatrix} A_{1,1} & & C_1 \\ & A_{2,2} & & C_2 \\ & & \dots & & \\ & & & A_{Nb,Nb} & C_{Nb} \\ & & & & R_{1} & R_2 & \dots & R_{Nb} & B \end{pmatrix},$$
(71)

- The blocks can have very different sizes. A nested dissection ordering can be used to permute a symmetrically structured matrix A to a symmetrically structured DBBD form  $(S\{R_i\} = S\{C_i^T\})$ .
- If S{A} is close to symmetric then nested dissection can be applied to S{A + A<sup>T</sup>}. In finite-element applications, the DBBD form corresponds to partitioning the underlying finite-element domain into non-overlapping subdomains.

# Reorderings

#### Bordered forms

- Coarse-grained parallel approaches aim to factorize the *A*<sub>*lb,lb*</sub> blocks in parallel before solving the interface problem that connects the blocks.
- The block factorization of *A*<sub>DB</sub> is

$$A_{DB} = \begin{pmatrix} L_1 & & & \\ & L_2 & & \\ & & \dots & & \\ & & & L_{Nb} \\ \widehat{R}_1 & \widehat{R}_2 & \dots & \widehat{R}_{Nb} & L_S \end{pmatrix} \begin{pmatrix} U_1 & & & \widehat{C}_1 \\ & U_2 & & & \widehat{C}_2 \\ & & \dots & & & \\ & & & U_{Nb} & \widehat{C}_{Nb} \\ & & & & U_S \end{pmatrix},$$

where

$$\widehat{R}_{lb} = R_{lb}U_{lb}^{-1}, \ \widehat{C}_{lb} = L_{lb}^{-1}C_{lb} \ (1 \le lb \le Nb), \ L_S U_S = B - \sum_{lb=1}^{Nb} \widehat{R}_{lb}\widehat{C}_{lb}.$$

• Here, for simplicity, no permutations emphasized.

#### Bordered forms

Algorithm (Coarse-grained parallel LU factorization using DBBD form)

Input: Matrix  $A_{DB}$  in DBBD form (71). Output: Block LU factorization.

1: Initialise S = B2: for lb = 1 : Nb do 3:  $A_{lb,lb} = L_{lb}U_{lb}$ 4:  $\hat{R}_{lb} = R_{lb}U_{lb}^{-1}$ 5:  $\hat{C}_{lb} = L_{lb}^{-1}C_{lb}$ 6: end for 7:  $S = S - \sum_{lb=1}^{Nb} \hat{R}_{lb}\hat{C}_{lb}$ 8:  $S = L_SU_S$ 

LU factorization of square block on diagonal
 Triangular solve for bottom-border blocks
 Triangular solve for right-border blocks

Assemble updates to interface block
 Factorize updated interface block (Schur complement)

# Reorderings

#### Bordered forms

- Factorization of each individual *A*<sub>*lb,lb*</sub> and solve steps can be parallelized.
- The assembly of the interface block *S* and its LU can be partially parallelized.
- S is generally significantly denser than the other blocks.
- If A is not SPD then factorizing the A<sub>lb,lb</sub> blocks without considering the entries in the border can potentially lead to stability problems. Consider the first step in factorizing A<sub>lb,lb</sub> and the threshold pivoting test for a sparse LU factorization. The pivot candidate (A<sub>lb,lb</sub>)<sub>11</sub> must satisfy

 $\max\{\max_{i>1} |(A_{lb,lb})_{i1}|, \max_{k} |(R_{lb})_{k1}|\} \le \gamma^{-1} |(A_{lb,lb})_{11}|,$ 

where  $\gamma \in (0,1]$  is the threshold parameter.

• Large entries in the row border matrix  $R_{lb}$  can prevent pivots being selected within  $A_{lb,lb}$ .

#### Singly bordered form

• Singly bordered block diagonal (SBBD) form

$$A_{SB} = \begin{pmatrix} A_{1,1} & & C_1 \\ & A_{2,2} & & C_2 \\ & & \ddots & & \\ & & & A_{Nb,Nb} & C_{Nb} \end{pmatrix},$$

 $A_{lb,lb}$  are rectangular  $m_{lb} \times n_{lb}$ . The linear system becomes

$$\begin{pmatrix} A_{1,1} & & C_1 \\ A_{2,2} & & C_2 \\ & \dots & & \vdots \\ & & & A_{Nb,Nb} & C_{Nb} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_{Nb} \\ x_I \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_{Nb} \end{pmatrix}, \quad (72)$$

 $x_{lb}$  is of length  $n_{lb}$ ,  $x_I$  is a vector of length  $n_I$  of interface variables, and the right-hand side vectors  $b_{lb}$  are of length  $m_{lb}$ , such that

$$\begin{pmatrix} A_{lb,lb} & C_{lb} \end{pmatrix} \begin{pmatrix} x_{lb} \\ x_I \end{pmatrix} = b_{lb}, \quad 1 \le lb \le Nb.$$

#### Singly bordered form

• A partial factorization of each block matrix is performed, that is,

$$\begin{pmatrix} A_{lb,lb} & C_{lb} \end{pmatrix} = P_{lb} \begin{pmatrix} L_{lb} \\ \bar{L}_{lb} & I \end{pmatrix} \begin{pmatrix} U_{lb} & \bar{U}_{lb} \\ S_{lb} \end{pmatrix} Q_{lb},$$
(73)

- Pivots can only be chosen from the columns of  $A_{lb,lb}$  because the columns of  $C_{lb}$  have entries in at least one other border block  $C_{jb}$   $(jb \neq lb)$ .
- The pivot candidate  $(A_{lb,lb})_{11}$  at the first elimination step must satisfy

$$\max_{i>1} |(A_{lb,lb})_{i1}| \le \gamma^{-1} |(A_{lb,lb})_{11}|,$$

#### Singly bordered form

# Algorithm (Coarse-grained parallel LU factorization and solve using SBBD form)

*Input:* Linear system in SBBD form (72). *Output:* Block LU factorization and computed solution x.

1: S = 0 and  $z_I = 0$ 2: for lb = 1 : Nb do Perform a partial LU factorization (73) of  $(A_{lb,lb}, C_{lb})$ . 3: Solve  $P_{lb} \begin{pmatrix} L_{lb} & \\ \bar{L}_{lb} & I \end{pmatrix} \begin{pmatrix} y_{lb} \\ \bar{y}_{lb} \end{pmatrix} = b_{lb}$ 4: 5:  $S = S + S_{lb}$  and  $z_I = z_I + \bar{y}_{lb}$  $\triangleright$  Assemble S and  $z_I$ 6: end for 7:  $S = P_{\circ}L_{\circ}U_{\circ}Q_{\circ}$  $\triangleright P_s$  and  $Q_s$  are permutation matrices 8: Solve  $P_s L_s y_I = z_I$  and then  $U_s Q_s x_I = y_I$ Forward then back substitution 9: for lb = 1 : Nb do Solve  $U_{lb} Q_{lb} x_{lb} = y_{lb} - \bar{U}_{lb} Q_{lb} x_{l}$ 10: 11: end for

#### Ordering to singly bordered form

- The objective is to permute *A* to an SBBD form with a narrow column border.
- One way to do this is to choose the number Nb > 1 of required blocks and use ND to compute a vertex separator V<sub>S</sub> of G(A + A<sup>T</sup>) such that removing V<sub>S</sub> and its incident edges splits G(A + A<sup>T</sup>) into Nb components.
- Then initialise the set S<sub>C</sub> of border columns to V<sub>S</sub> and let V<sub>1b</sub>, V<sub>2b</sub>,..., V<sub>Nb</sub> be the subsets of column indices of A that correspond to the Nb components and let n<sub>i,kb</sub> be the number of column indices in row i that belong to V<sub>kb</sub>. If

 $lb = \arg \max_{1 \le kb \le Nb} |n_{i,kb}|$  then row *i* is assigned to partition *lb*. All column indices in row *i* that do not belong to  $\mathcal{V}_{lb}$  are moved into  $\mathcal{S}_{\mathcal{C}}$ .

## Ordering to singly bordered form

# Algorithm (SBBD ordering of a general matrix)

Input: Matrix A, the number Nb > 1 of blocks, vertex separator  $\mathcal{V}_{\mathcal{S}}$  of  $\mathcal{G}(A + A^T)$ . **Output:** Vector block such that block(i) denotes the partition in the SBBD form to which row i is assigned  $(1 \le i \le n)$  and  $\mathcal{S}_{\mathcal{C}}$  is the set of border columns.

- 1: Initialise  $S_C = V_S$  and block(1:n) = 0
- 2: Initialise  $\mathcal{V}_{kb}$  to hold the column indices of A that correspond to component kb of  $\mathcal{G}(A + A^T)$ after the removal of  $\mathcal{V}_S$ ,  $1 \le kb \le Nb$
- 3: for each row i do
- 4: Add up the number  $n_{i,kb}$  of column indices belonging to  $V_{kb}$ ,  $1 \le kb \le Nb$
- 5: Find  $lb = \arg \max_{1 \le kb \le Nb} n_{i,kb}$
- 6: block(i) = lb
- 7: for each column index j in row i do
- 8: **if**  $j \in \mathcal{V}_{kb}$  and  $kb \neq lb$  then
- 9: Remove j from  $V_{kb}$  and add to  $S_C$
- 10: end if
- 11: end for
- 12: end for
- 13: Assign the rows i for which block(i) = 0 equally between the Nb partitions.
- 14: If some column  $j \in S_C$  has nonzero entries only in rows belonging to partition kb then  $_{484/609}$

# Complexity

- Overall time dominated by time for the factorization
- General dense matrices
  - Space:  $O(n^2)$
  - Time:  $O(n^3)$
- General sparse matrices
  - Space:  $\eta(L) = n + \sum_{i=1}^{n-1} (\eta(L_{*i}) 1)$
  - ► Time in the *i*-th step:  $\eta(L_{*i}) 1$  divisions,  $1/2(\eta(L_{*i}) 1)\eta(L_{*i})$  multiple-add pairs
  - Time totally:  $1/2 \sum_{i=1}^{n-1} (\eta(L_{*i}) 1)(\eta(L_{*i}) + 2)$





# Complexity

- Band schemes ( $\beta << n$ )
  - ▶ Space:  $O(\beta n)$
  - Time:  $O(\beta^2 n)$



## Complexity

- Profile/envelope schemes
  - Space:  $\sum_{i=1}^{n} \beta_i$
  - Frontwidth:  $\omega_i(A) = |\{k|k > i \land a_{kl} \neq 0 \text{ for some } l \leq i\}|$
  - Time:  $1/2 \sum_{i=1}^{n-1} \omega_i(A)(\omega_i(A) + 3)$



# From direct to iterative methods

## Complexity

• General sparse schemes can be analyzed in some cases

Nested dissection



#### Definition

 $(\alpha, \sigma)$  separation of a graph with n vertices: each its subgraph can be separated by a vertex separator S such that its size is of the order  $O(n^{\sigma})$  and the separated subgraphs components have sizes  $\leq \alpha n, 1/2 \leq \alpha < 1$ .

# From direct to iterative methods

## Complexity: Generalized nested dissection



Vertex separator S

- Planar graphs, 2D finite element graphs (bounded degree)
  - $\sigma = 1/2, \, \alpha = 2/3$
  - Space:  $O(n \log n)$
  - Time:  $O(n^{3/2})$
- 3D Finite element graphs
  - σ = 2/3
  - ▶ Space: *O*(*n*<sup>4/3</sup>)
  - ▶ Time: *O*(*n*<sup>2</sup>)
- Lipton, Rose, Tarjan (1979), Teng (1997).

# Outline

#### Algebraic preconditioning

- Finite precision arithmetic: computed factors are not exact.
- Moreover, the effort to obtain more accurate results can lead to complex coding and unavoidable inefficiencies magnified by modern computer architectures.
- Potential solution: intentionally relaxing the required accuracy of the computed factors.
- Simpler, cheaper, sparser approximate factorization of *A* (or of *A*<sup>-1</sup>): preconditioners.
- Using the preconditioner in combination with an iterative solver.
- an approximate factorization called an incomplete factorization to distinguish it from a complete factorization of a direct method.
- Used with iterative methods for solving Ax = b from their two main classes: stationary (relaxation) iterative methods and Krylov subspace methods.
• Stationary iterative methods work by splitting *A* as follows:

$$A = M - N,$$

• The matrix M is chosen to be nonsingular and easy to invert. An initial guess  $x^{(0)}$ , the iterations are then given by

$$x^{(k+1)} = M^{-1}Nx^{(k)} + M^{-1}b, \quad k = 0, 1, \dots$$
 (74)

This can be rewritten as

$$x^{(k+1)} = x^{(k)} + M^{-1}(b - Ax^{(k)}) = x^{(k)} + M^{-1}r^{(k)}, \quad k = 0, 1, \dots$$
(75)

where the vector  $r^{(k)} = b - Ax^{(k)}$  is the residual on the *k*-th iteration.

# Algebraic preconditioning

#### Stationary iterative methods

• By substituting  $b = r^{(k)} + Ax^{(k)}$  into  $x = A^{-1}b$ , we obtain

$$x = A^{-1}(r^{(k)} + Ax^{(k)}) = x^{(k)} + A^{-1}r^{(k)}$$

If *M* is used to approximate *A*, we again get the iteration above. • Further

$$r^{(k+1)} = b - A(x^{(k)} + M^{-1} r^{(k)}) = (I - AM^{-1}) r^{(k)} = \dots = (I - AM^{-1})^k$$
(76)

and if  $e^{(k)} = x - x^{(k)}$  is the error vector on iteration k then

$$e^{(k+1)} = M^{-1}N e^{(k)} = \dots = (M^{-1}N)^{k+1} e^{(0)} = (I - M^{-1}A)^{k+1} e^{(0)}.$$
(77)

• The matrix  $I - M^{-1}A$  or  $I - AM^{-1}$  is called the iteration matrix. In general, monitoring residuals is more practical.

#### Theorem

For any initial  $x^{(0)}$  and vector *b*, the stationary iteration converges if and only if the spectral radius of  $(I - M^{-1}A)$  is less than unity.

# Proof.

The spectral radius of an  $n \times n$  matrix C with eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_n$  is defined to be

$$\rho(C) = \max\{|\lambda_i| \mid 1 \le i \le n\}.$$
(78)

Furthermore, the sequence of matrix powers  $C^k$ , k = 0, 1, ..., converges to zero if and only if  $\rho(C) < 1$ . It follows from (77) that if the spectral radius of  $(I - M^{-1}A)$  is less than unity then the iteration (74) converges for any  $x^{(0)}$  and *b*. Conversely, the relation

$$x^{(k+1)} - x^{(k)} = (I - M^{-1}N)(x^{(k)} - x^{(k-1)}) = \dots = (I - M^{-1}N)^k M^{-1}(b - Ax^{(0)})$$

shows that if the iteration converges for any  $x^{(0)}$  and b then  $(I - M^{-1}N)^k v$  converges to zero for any v. Consequently,  $\rho(I - M^{-1}A)$  must be less than unity, and the result follows.

- Generally impractical to compute the spectral radius and sufficient conditions that guarantee convergence are used.
- A sufficient condition for convergence is  $||I M^{-1}A|| < 1$ .
- Consider splitting (diagonal, strict lower triangular, strict upper triangular)

$$A = D_A + L_A + U_A. \tag{79}$$

- For  $\omega > 0$  is a scalar parameter, classical methods include:
  - Richardson method:  $M = \omega^{-1}I$ ,
  - ► Jacobi and damped Jacobi methods:  $M = D_A$  and  $M = \omega^{-1}D_A$ ,
  - Gauss-Seidel and SOR methods:  $M = D_A + L_A$  and  $M = \omega^{-1}D_A + L_A$ .

Theorem

If  $A \in \mathbb{R}^{n \times n}$  is strongly diagonally dominant then Jacobi method and Gauss-Seidel method are convergent.

#### Theorem

If  $A \in \mathbb{R}^{n \times n}$  is symmetric with positive diagonal  $D_A$  then the Jacobi method is convergent iff A and  $2D_A - A$  are positive definite.

#### Theorem

If  $A \in \mathbb{R}^{n \times n}$  is symmetric and positive definite then the Gauss-Seidel method is convergent.

Non-stationary iterative methods are of the form

 $x^{(k+1)} = x^{(k)} + \omega^{(k)} M^{-1} r^{(k)}, \quad k = 0, 1, \dots$ 

where the  $\omega^{(k)}$  are scalars.

- In this class, Krylov subspace methods are the most effective.
- Given a vector y, the k-th Krylov subspace K<sup>(k)</sup>(A, y) generated by A from the vector y is defined to be

$$\mathcal{K}^{(k)}(A, y) = span(y, Ay, \dots, A^{k-1}y).$$

- Generate a sequence of approximate solutions  $x^{(k)} \in x^{(0)} + \mathcal{K}^{(k)}(A, r^{(0)})$  such that the norm of the corresponding residuals  $r^{(k)} \in \mathcal{K}^{(k+1)}(A, r^{(0)})$  converge to zero.
- SPD systems the conjugate gradient method (CG); nonsymmetric systems: GMRES, BiCG, no single method of choice.
- At each iteration only matrix-vector products with *A* (and possibly with *A*<sup>T</sup>) required.

- Powerful, if (and only if) combined with a preconditioner: the most widely-used class of preconditioned iterative methods.
- Unfortunately, for a given A, b and  $x^{(0)}$ , it is usually not possible to predict the rate of convergence.
- If A is a SPD matrix then for CG

$$||x - x^{(k)}||_A \le 2\left(\frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1}\right)^k ||x - x^{(0)}||_A,$$

where  $\kappa(A)$  is the spectral condition number.

• Can be highly pessimistic. Does not show the potential for CG to converge superlinearly or that the rate of convergence depends on the distribution of all the eigenvalues of *A*.

- For non-SPD matrices, less is known. Like to emphasize favourable properties: the minimal residual method (MINRES) for solving symmetric indefinite systems in exact arithmetic, has the norm of the residual monotonically decreasing.
- No general descriptive convergence theory is available for Krylov subspace methods for nonsymmetric systems (including GMRES).

 Preconditioning corresponds to the application of a matrix (or linear operator) to the original linear system to yield a different linear system that has more favourable properties. Consider the preconditioned linear system

$$M^{-1}Ax = M^{-1}b.$$
 (80)

Here  $M^{-1}$  is applied to *A* from the left. We say that *A* is preconditioned from the left and *M* is a left preconditioner. Analogously, the linear system can be preconditioned from the right

$$AM^{-1}y = b, \qquad x = M^{-1}y.$$
 (81)

The following result states that it is not possible to determine a priori which variant is the best.

## Theorem

Let  $\delta$  and  $\Delta$  be positive numbers. Then for any  $n \geq 3$  there exist nonsingular  $n \times n$  matrices A and M such that all the entries of  $M^{-1}A - I$  have absolute value less than  $\delta$  and all the entries of  $AM^{-1} - I$  have absolute values greater than  $\Delta$ .

• The choice between left and right preconditioning may be based on properties of the coupling of the preconditioner with the iterative method or on the distribution of the eigenvalues of *A*.

- The computed quantities readily available during a preconditioned iterative method depend on how the preconditioner is applied and this may influence the choice. These quantities may be used, for example, to decide when to terminate the iterations.
- An obvious advantage of right preconditioning is that in exact arithmetic, the residuals for the right preconditioned system are identical to the true residuals, enabling convergence to be monitored accurately.
- In some cases, the numerical properties of an implementation and/or the computer architecture may also play a part.

• For M in factorized form  $M = M_1M_2$ , two-sided (or split) preconditioning is an option. The iterative method then solves the transformed system

$$M_1^{-1}AM_2^{-1}y = M_1^{-1}b, \qquad x = M_2^{-1}y.$$
 (82)

• If A and M are SPD matrices then we would like the preconditioned matrix  $M_1^{-1}AM_1^{-T}$  to be SPD. However, it is not necessary to use a two-sided transformation with the preconditioned conjugate gradient (PCG) method because it can be formulated using the *M*-inner product in which the matrix  $M^{-1}A$  is self-adjoint.

#### Theorem

Let A and M be SPD matrices. Then  $M^{-1}A$  is self-adjoint in the M-inner product.

# Proof.

Self-adjointness is implied by the following chain of equivalences.

 $\langle M^{-1}Ax,y\rangle_M=\langle Ax,y\rangle=\langle x,Ay\rangle=\langle x,MM^{-1}Ay\rangle=\langle Mx,M^{-1}Ay\rangle=\langle x,M^{-1}Ay\rangle_M.$ 

- Left preconditioned CG with the *M*-inner product is mathematically equivalent to right preconditioned CG with the *M*<sup>-1</sup>-inner product.
- If *A* is symmetric but not PD, the PCG method can breakcdown (division by a zero quantity).

- An obvious goal: to achieve rapid convergence.
- The preconditioner should aim to reduce the condition number, but this is not necessarily sufficient to give fast convergence.
- For general matrices, despite the lack of theoretical guarantees regarding convergence, many useful preconditioners motivated by bounding the condition number of the preconditioned matrix.
- Choosing a preconditioner is often based on how costly it is to compute and on some indicators that potentially reflect its quality.
- In particular, the accuracy of a preconditioner *M* can be assessed using the norm of the error matrix

$$||E|| = ||M - A||,$$

and its stability can be measured using

 $\|M^{-1}E\| = \|I - M^{-1}A\| \quad \text{or} \quad \|EM^{-1}\| = \|I - AM^{-1}\|.$ 

 In some cases, the inverse M<sup>-1</sup> is computed directly. In this case we have an approximate inverse preconditioner.

#### Simple preconditioners

- The simplest preconditioner consists of the diagonal of the matrix  $M = D_A$ . This is known as the (point) Jacobi preconditioner.
- Block versions can be derived by partitioning V = {1, 2, ..., n} into mutually disjoint subsets V<sub>1</sub>, ..., V<sub>l</sub> and then setting

 $m_{ij} = \begin{cases} a_{ij} & \text{if } i \text{ and } j \text{ belong to the same subset } \mathcal{V}_k \text{ for some } k, 1 \\ 0 & \text{otherwise.} \end{cases}$ 

• The SSOR preconditioner, like the Jacobi preconditioner, can be derived from *A* without any work. If *A* is symmetric then using the notation (79), the SSOR preconditioner is defined to be

$$M = (D_A + L_A)D_A^{-1}(D_A + L_A)^T,$$
(83)

or, using a parameter  $0 < \omega < 2$ , as

$$M = \frac{1}{2 - \omega} (\frac{1}{\omega} D_A + L_A) (\frac{1}{\omega} D_A)^{-1} (\frac{1}{\omega} D_A + L_A)^T.$$

• Finding optimal value of  $\omega$  typically expensive.

#### The Eisenstat trick

- Generally cheaper to apply  $M^{-1}$  and A separately,
- But: consider the matrix splitting (79) and let *M* be given by

$$M = (D + L_A) \left[ D^{-1} (D + U_A) \right] = M_1 M_2,$$
(84)

where D is a nonsingular diagonal matrix.

- The SSOR matrix is one example in the symmetric case but more generally D ≠ D<sub>A</sub>.
- Using two-sided preconditioning, (82) becomes

$$A'y = M_1^{-1}AM_2^{-1}y = (D+L_A)^{-1}A[D^{-1}(D+U_A)]^{-1}y = (D+L_A)^{-1}b.$$
(85)

# Setting

$$\bar{L} = D^{-1}L_A, \ \bar{U} = D^{-1}U_A, \ \bar{A} = D^{-1}A, \ \text{and} \ \bar{b} = (I + \bar{L})^{-1}D^{-1}b,$$
  
we obtain

$$A' = (D+L_A)^{-1}A[D^{-1}(D+U_A)]^{-1} = [(D+L_A)^{-1}D]D^{-1}A[D^{-1}(D+U_A)]^{-1}$$
$$= [D^{-1}(D+L_A)]^{-1}D^{-1}A(I+D^{-1}U_A)^{-1} = (I+\bar{L})^{-1}\bar{A}(I+\bar{U})^{-1}.$$

#### The Eisenstat trick

That is, the system becomes

 $A'y = (I + \bar{L})^{-1}\bar{A}(I + \bar{U})^{-1}y = (I + \bar{L})^{-1}D^{-1}b = \bar{b}.$  (86)

If y solves (86) then the solution x of  $(I + \overline{U}) x = y$  solves Ax = b. • Further

$$A' = (I + \bar{L})^{-1} (I + \bar{L} + D^{-1}D_A - 2I + I + \bar{U})(I + \bar{U})^{-1}$$

 $= (I+\bar{L})^{-1} \left[ (I+\bar{L})(I+\bar{U})^{-1} + (D^{-1}D_A - 2I)(I+\bar{U})^{-1} + I \right]$ 

$$= (I + \bar{U})^{-1} + (I + \bar{L})^{-1} [(D^{-1}D_A - 2I)(I + \bar{U})^{-1} + I].$$

• Thus to compute  $z = A'w = (I + \overline{L})^{-1}\overline{A}(I + \overline{U})^{-1}w$  for a given w, it is necessary only to solve two triangular systems

 $(I + \overline{U}) z_1 = w$  followed by  $(I + \overline{L}) z_2 = (D^{-1}D_A - 2I) z_1 + w$ ,

and then set  $z = z_1 + z_2$ .

 This trick is not a preconditioner: it is a way of applying the preconditioner of special shape.

508/609

- The development of algebraic preconditioners has been closely connected to solving linear systems from the discretization of PDEs.
- Two-dimensional Poisson problem, discretized using a uniform regular grid, finite differences, zero Dirichlet conditions on the boundary, natural ordering.

$$A = \begin{pmatrix} 4 & -1 & -1 & & \\ -1 & 4 & -1 & -1 & & \\ -1 & 4 & -1 & -1 & & \\ -1 & 4 & -1 & -1 & & \\ -1 & -1 & 4 & -1 & -1 & \\ & -1 & -1 & 4 & -1 & \\ & & -1 & -1 & 4 & -1 \\ & & & -1 & -1 & 4 & -1 \\ & & & & -1 & -1 & 4 \end{pmatrix}.$$
 (87)

If the spatial discretization on the domain is characterized by the mesh parameter *h* then the size of *A* is inversely proportional to *h*.
 κ(A) depends asymptotically on h<sup>-2</sup>.

- Matrices with similar banded sparsity patterns with nonzeros on only a small number of subdiagonals arise from simple finite difference or finite element discretizations of other partial differential equations.
- Particular cases of special classes of matrices help to describe the theoretical background behind the discretized systems.
- Let the off-diagonal entries of the nonsingular matrix A be nonpositive (that is, a<sub>ij</sub> ≤ 0 for all i ≠ j). Then A is a (nonsingular) M-matrix if one of the following holds:
  - ► *A* + *D* is nonsingular for any diagonal matrix *D* with nonnegative entries;
  - all the entries of  $A^{-1}$  are nonnegative;
  - ▶ all principal minors of *A* are positive.
- The matrix above is an example of an M-matrix. A symmetric M-matrix is known as a Stieltjes matrix, and such a matrix is positive definite.

• The class of nonsingular H-matrices includes matrices coming from simple discretizations of convection-diffusion problems. The comparison matrix C(A) of A is defined to have entries

$$C(A)_{ij} = \begin{cases} -|a_{ij}|, & i \neq j, \\ |a_{ij}|, & i = j. \end{cases}$$

- If C(A) is a nonsingular M-matrix then A is a nonsingular H-matrix.
- A is diagonally dominant by rows if

$$\sum_{j=1, \ j \neq i}^{n} |a_{ij}| \le |a_{ii}|, \quad 1 \le i \le n.$$
(88)

• A is strictly diagonally dominant by rows if strict inequality holds for all *i*. A is (strictly) diagonally dominant by columns if  $A^T$  is (strictly) diagonally dominant by rows.

- *A* is said to be irreducibly diagonally dominant if it is irreducible and the inequalities are satisfied with strict inequality for at least one row *i*. If *A* is strictly diagonally dominant by rows or columns or is irreducibly diagonally dominant then it is nonsingular and factorizable.
- The class of diagonally dominant matrices is closely connected to that of nonsingular H-matrices. For example, the property that there exists a diagonal matrix *D* with positive entries such that *AD* is strictly diagonally dominant is equivalent to *A* being a nonsingular H-matrix.

#### Some special classes of matrices: once more

- Still theoretical assumptions are rather strong.
- Concept of special matrices

#### Theorem

Matrix A is called a regular M-matrix if  $a_{ij} \leq 0, i \neq j$ , is regular and  $A^{-1} \geq 0$ .

#### Theorem

A is a **H-matrix** if  $B = |D_A| - |A - D_A|$  is an *M*-matrix.

#### Many equivalent definitions

# Introduction to incomplete factorizations

- Incomplete factorizations fall into three main classes:
  - ► Threshold-based methods: locations of permissible fill-in are determined in conjunction with the numerical factorization of *A*; entries of the computed factors of absolute value greater than a prescribed threshold  $\tau > 0$  are dropped.
  - Memory-based methods in which the amount of memory available for the incomplete factorization is prescribed and only the largest entries in each row (or column) are retained.
  - Structure-based methods: an initial symbolic factorization phase determines the location of permissible entries using only S{A}. This allows the memory requirements to be determined before an incomplete numerical factorization is performed. The specified set of positions is called the target sparsity pattern.

Introduction to incomplete factorizations

- The basic dropping approaches can be combined and employed in conjunction with sparsifying *A* before the factorization commences.
- Sparsification of *A* after permuting reveals a block structure.

	1	<b>2</b>	3	4	5	6		2	4	1	3	5	6			2	4	1	3	5	6
1	(*	*	*				2	2 /*	*	*		*	* \		2	(*	*			*	* \
2	*	*	f	*	*	*	4	1   *	*	f		*	*		4	*	*			*	*
3	*	f	*	f	f	f		L ×	f	*	*	f	f	1	1			*	*		
4		*	f	*	*	*	:	3		*	*	f	f		3			*	*		
5		*	f	*	*	*	Ę	5 *	*	f	f	*	*		5	*	*			*	*
6		*	f	*	*	*/	(	3 \*	*	f	f	*	*/		6	/*	*			*	*/

Figure: Illustration of matrix sparsification. f denotes filled entries in the factors. On the left is the original matrix A with its filled entries; in the centre is the permuted matrix with its filled entries; on the right is the sparsified permuted matrix after dropping the entries of A in positions (1,3) and (3,1) (it has no filled entries).

- Dropping entries can lead to breakdown of the incomplete factorization, that is, a zero pivot may be encountered during the factorization (or a non positive pivot in the Cholesky case).
- It is only possible to predict when this will happen in special cases, as stated in the following theorem, which is a consequence of the fact that being an M-matrix or an H-matrix is preserved in the sequence of the Schur complements during the factorization.

# Theorem

Let *A* be a nonsingular *M*-matrix or *H*-matrix. If the target sparsity pattern of the incomplete factors contains the positions of the diagonal entries then the incomplete factorization of *A* does not break down.

- To illustrate the error accumulation in the incomplete factorization of an M-matrix using dropping, consider the following example.
- Let *E* be the error matrix. *E* is initialised to zero and at each stage of the factorization the dropped entries are added into it.
- After one step of the complete factorization of *A* the partially eliminated matrix *A*<sup>(2)</sup> is

$$A^{(2)} = \begin{pmatrix} 4 & -1 & -1 & & \\ 3.75 & -1 & -0.25 & -1 & & \\ -1 & 4 & & -1 & & \\ -0.25 & 3.75 & -1 & -1 & & \\ -1 & -1 & 4 & -1 & -1 & \\ & -1 & -1 & 4 & -1 & \\ & & -1 & -1 & 4 & -1 \\ & & & -1 & -1 & 4 & -1 \\ & & & & -1 & -1 & 4 & -1 \\ & & & & & -1 & -1 & 4 \end{pmatrix}$$

- Suppose the filled entries -0.25 in positions (2, 4) and (4, 2) are dropped. Then the values of the corresponding diagonal entries in the subsequent elimination matrices are larger than they would have been without any dropping.
- Furthermore, as all the off-diagonal nonzero entries are negative, for any target sparsity pattern the dropped entries are negative.
- The M-matrix property applies to all subsequent Schur complements, which implies that all the entries added into *E* are negative and so the absolute values of the entries in *E* grow as the factorization proceeds (the contributions can never cancel each other out).
- Thus, although the factorization does not break down, the growth in the error is potentially a problem for the accuracy of an incomplete factorization of an M-matrix.
- Modifying the diagonal entries of *A* is a common approach to avoid breakdown in an incomplete factorization.

# Algebraic preconditioning

#### Incomplete factorization breakdown

$$A = \begin{pmatrix} 3 & -2 & 2 \\ -2 & 3 & -2 \\ 2 & -2 & 8 \end{pmatrix}, \ L = \begin{pmatrix} 1 & & \\ -2/3 & 1 \\ & -6/5 & 1 \\ 2/3 & 4/5 & -2/3 & 1 \end{pmatrix}, \ D = \begin{pmatrix} 3 & & \\ 5/3 & & \\ & 3/5 & \\ & 16/3 \end{pmatrix}$$
$$\widetilde{L} = \begin{pmatrix} 1 & & & \\ -2/3 & 1 & & \\ & -6/5 & 1 \\ 2/3 & & -10/3 & 1 \end{pmatrix}, \ \widetilde{D} = \begin{pmatrix} 3 & & \\ 5/3 & & \\ & 3/5 & \\ & 0 \end{pmatrix}.$$

Figure: An example to illustrate breakdown. The matrix A and its square-root free factors are given together with the incomplete factors  $\tilde{L}$  and  $\tilde{D}$  that result from dropping the entry  $l_{24}$  during the factorization.  $\tilde{d}_{44} = 0$  means the incomplete factorization has broken down.

- Remedy: perturb the diagonal value causing the breakdown.
- Practice of making simple ad hoc modifications not very positive.
   If breakdown (or near-breakdown) occurs, it may be too late. 519/609

- An alternative and more effective strategy to avoid breakdown is to modify all the diagonal entries of *A* a priori and then compute an incomplete factorization of  $A + \alpha I$ , where the shift  $\alpha > 0$  is a scalar parameter.
- It is always possible to find  $\alpha$  such that  $A + \alpha I$  is nonsingular and diagonally dominant and is thus an H-matrix.
- However, being an H-matrix is not a necessary condition for a matrix to be factorizable and, in practice, much smaller values of α can provide incomplete factorizations for which ||*E*|| is small.
- A simple trial-and-error procedure for choosing a shift is given below.

Algorithm (Trial-and-error global shifted incomplete factorization)

Input: Matrix A, incomplete factorization algorithm, initial shift  $\alpha^{(0)}$ Output: Shift  $\alpha$  such that  $A + \alpha I \approx \widetilde{L}\widetilde{U}$ .

1: for 
$$k = 0, 1, 2, \dots$$
 do  
2:  $A + \alpha^{(k)}I \approx \widetilde{L}\widetilde{U}$   
3: If successful,  $\alpha = \alpha^{(k)}$  and return  
4:  $\alpha^{(k+1)} = 2\alpha^{(k)}$   
5: end for

> Perform incomplete factorization

- An alternative approach to avoid small pivots:to follow what is done in sparse direct solvers and incorporate partial or threshold pivoting within the incomplete factorization algorithm: preprocessing by reordering, scaling etc.
- One way to attempt to minimize the norm of the error matrix *E* is to select the pivot candidate to minimize the sum of the absolute values of the dropped (discarded) entries. However, this minimum discarded fill ordering is typically too expensive to be useful in practice.

# Factorizations as preconditioner components

 Sometimes (incomplete) factorizations are employed as components in the construction of more complex preconditioners. Here some possible approaches are briefly discussed.

# Incomplete factorization breakdown Polynomial preconditioning

• Polynomial preconditioning selects a polynomial  $\phi$  and applies a Krylov subspace method to solve either

$$\phi(A)Ax = \phi(A) b$$

(left preconditioning) or

$$A \phi(A) y = b, \quad x = \phi(A) y$$

(right preconditioning).  $\phi$  should be of small degree and chosen to enhance convergence.

- Consider the characteristic polynomial  $\phi_n(\mu) = \det(A \mu I)$  of A (det denotes the determinant).
- The Cayley-Hamilton theorem states that *A* satisfies its own characteristic equation so that

$$\phi_n(A) = \sum_{j=0}^n \beta_j A^j = 0,$$
524/609

# Polynomial preconditioning

• Provided A is nonsingular,

$$A^{-1} = (-1)^{n+1} \frac{1}{\det(A)} \sum_{j=1}^{n} \beta_j A^{j-1}.$$

A preconditioner can be constructed by taking the first k terms, possibly weighted by some suitable scalar coefficients, that is,

$$M^{-1} = \sum_{j=0}^k \gamma_j A^k.$$

# Polynomial preconditioning

- An important question is why such a preconditioner can help in the presence of the optimality properties of Krylov subspace methods.
- For example, at iteration k + 1 of the CG method,  $x^{(k+1)}$  satisfies

$$x^{(k+1)} = x^{(0)} + \phi_k(A) r^{(0)}, \ k = 0, 1, \dots,$$

where  $\phi_k$  is a monic polynomial of degree k. This polynomial is optimal in the sense that  $x^{(k+1)}$  minimizes

$$\|x - x^{(k+1)}\|_A^2.$$
(89)

- A preconditioner that is a polynomial in A cannot speed the convergence because the resulting iteration again forms the new x<sup>(k+1)</sup> as x<sup>(0)</sup> plus a polynomial in A times r<sup>(0)</sup>, and thus the same or a higher degree polynomial is needed to achieve the same value of the A-norm of the error.
- Consequently, the number of matrix-vector multiplications cannot decrease.

# Polynomial preconditioning

- Nevertheless, polynomial preconditioning can be useful for a number of reasons.
  - The polynomial can improve the eigenvalue distribution of the preconditioned matrix and result in a reduction in the number of iterations required for convergence (even though the overall complexity may increase).
  - It requires very little memory and its implementation can be straightforward.
  - It can decrease the number of synchronization points in iterative methods as represented by inner products. This is potentially important for message-passing parallel architectures.
### Polynomial preconditioning

- Even if only a small number of terms are used in approximating  $A^{-1}$ , a crucial issue is getting  $\gamma_0, \ldots, \gamma_k$ .
- A straightforward way of doing this: based on the Neumann series of a matrix C given by  $\sum_{j=0}^{+\infty} C^j$ , which is convergent if and only if  $\rho(C) < 1$ .
- In this case,

$$(I-C)^{-1} = \sum_{j=0}^{+\infty} C^j.$$
 (90)

- Now let  $\overline{M}$  be a nonsingular matrix and  $\omega > 0$  a scalar such that the matrix  $C = I \omega \overline{M}^{-1}A$  satisfies  $\rho(C) < 1$ .
- Using (90),

$$A^{-1} = \omega(\omega \bar{M}^{-1} A)^{-1} \bar{M}^{-1} = \omega (I - C)^{-1} \bar{M}^{-1} = \omega \left( \sum_{j=0}^{+\infty} C^j \right) \bar{M}^{-1}.$$

Polynomial preconditioning

• Truncating the summation gives as a possible preconditioner

$$M^{-1} = \omega \left(\sum_{j=0}^{k} C^{j}\right) \bar{M}^{-1}.$$

Observe that

$$I - M^{-1}A = I - \omega \left(\sum_{j=0}^{k} C^{j}\right) \bar{M}^{-1}A = I - \left(\sum_{j=0}^{k} C^{j}\right) (I - C) = C^{k+1},$$

which shows the positive effect of increasing k. If A and  $\overline{M}$  are SPD matrices then M can be used with the CG method preconditioned from the left because  $M^{-1}A$  is self-adjoint in the  $\overline{M}$ -inner product.

• Generalizations of the approach weight the powers of C in  $M^{-1}$  using additional scalars. The choice of  $\overline{M}$  is crucial for the effectiveness of the approach.

# Outline



# World of incomplete factorizations

- Direct factorizations may not be feasible (data structures and pivoting, operation counts, stability)
- Even by direct factorizations improving solution when using less accurate arithmetic (smaller  $\epsilon$ ) may be needed.
- The incomplete factors denoted here by  $\widetilde{L}$  and  $\widetilde{U}$ ;
- SPD case:  $\widetilde{U} = \widetilde{L}^T$ .
- We assume that the sparsity patterns of *A* and its incomplete factors always include the positions of the diagonal entries.
- Notation (other mentioned later): ILU(0) factorization (or an IC(0) factorization if A is SPD):  $S{\widetilde{L} + \widetilde{U}} = S{A}$ .

#### Exactness within the target sparsity pattern

#### Theorem

Consider the incomplete LU factorization  $A + E = \tilde{L}\tilde{U}$  with sparsity pattern  $S{\tilde{L} + \tilde{U}}$ . The entries of the error matrix E are zero at positions  $(i, j) \in S{\tilde{L} + \tilde{U}}$ .

### Proof.

The result clearly holds for j = 1. Let  $(i, j) \in S{\{\tilde{L} + \tilde{U}\}}$  and assume without loss of generality that i > j > 1. The (i, j) entry of  $\tilde{L}$  is computed as

$$\tilde{l}_{ij} = \left(a_{ij} - \sum_{k=1}^{j-1} \tilde{l}_{ik} \, \tilde{u}_{kj}\right) / \tilde{u}_{jj}$$

with the sums over k implying  $(i,k)\in \mathcal{S}\{\widetilde{L}+\widetilde{U}\}$  and  $(k,j)\in \mathcal{S}\{\widetilde{L}+\widetilde{U}\}.$  This gives

$$a_{ij}=\widetilde{L}_{i,1:j-1}\widetilde{U}_{1:j-1,j}+\widetilde{l}_{ij}\widetilde{u}_{jj}=\widetilde{L}_{i,1:j}\widetilde{U}_{1:j,j}=L_{i,1:j}U_{1:j,j},$$

and the corresponding entry of E is zero.

# Incomplete factorizations and patterns

- Theorem  $\Rightarrow$  extending  $S{\widetilde{L} + \widetilde{U}}$  gives a larger set of entries of A for which  $(E)_{ij} = 0$ .
- In some situations, there are straightforward ways to extend S{*ũ* + *ũ*}. In simple discretizations of a PDE may be a natural choice is to allow S{*ũ* + *ũ*} to include fill-in along a few additional diagonals within the band.

$$A = \begin{cases} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 2 & * & * & * & * & * \\ 3 & * & * & * & * & * \\ 5 & * & * & * & * & * \\ 6 & 7 & * & * & * & * \\ 7 & * & * & * & * & * \\ 8 & * & * & * & * & * \\ 7 & * & * & * & * & * \\ 8 & * & * & * & * & * \\ \end{array}$$

8

#### Incomplete factorizations and patterns



Figure: An  $8 \times 8$  banded sparse SPD matrix A and its graph  $\mathcal{G}(A)$ . The first three steps of a Cholesky factorization are shown. Filled entries are denoted by f.

### Crout incomplete factorizations

• The Crout variant: computes  $\widetilde{U}$  (by rows) and  $\widetilde{L}$  (by columns).

# Algorithm (Crout incomplete LU factorization)

Input: Matrix A and, optionally, a target sparsity pattern  $S{\widetilde{L} + \widetilde{U}}$ . Output: Incomplete LU factorization  $A \approx \widetilde{L}\widetilde{U}$ .

### Row incomplete factorizations

# Algorithm (Row incomplete LU factorization)

Input: Matrix A and, optionally, a target sparsity pattern  $S{\widetilde{L} + \widetilde{U}}$ . Output: Incomplete LU factorization  $A \approx \widetilde{L}\widetilde{U}$ .

#### Incomplete factorizations based on shortest fill-paths

 Entries of the factors that correspond to nonzero entries of A are assigned the level 0 while each potential filled entry in position (i, j) is assigned a level as follows:

$$level(i,j) = \min_{1 \le k < \min\{i,j\}} (level(i,k) + level(k,j) + 1).$$
(91)

- Given ℓ ≥ 0, during the factorization a filled entry is permitted at position (i, j) provided level(i, j) ≤ ℓ.
- The resulting level-based incomplete factorization is denoted by ILU(l) (or IC(l)); the basic row variant is given below.

## Algorithm (Level-based incomplete LU factorization)

1: Initialise level to 0 for nonzeros and diagonal entries of A and to n + 1 otherwise 2: for i = 1 : n do ▷ Loop over rows  $\tilde{l}_{ii} = 1, \tilde{L}_{i,1:i-1} = A_{i,1:i-1}$  and  $\tilde{U}_{i,i:n} = A_{i,i:n}$  $\triangleright$  Initialise row i of L and U 3: 4: for k = 1 : i - 1 such that  $level(i, k) < \ell$  do 5:  $\tilde{l}_{ik} = \tilde{l}_{ik}/\tilde{u}_{kk}$ 6: for j = k + 1 : i - 1 do 7:  $\tilde{l}_{ij} = \tilde{l}_{ij} - \tilde{l}_{ik} \tilde{u}_{kj}$  and update level(i, j)8: end for 9: for j = i : n do 10:  $\tilde{u}_{ij} = \tilde{u}_{ij} - \tilde{l}_{ik} \tilde{u}_{kj}$  and update level(i, j)11: end for 12: end for 13: for k = 1 : i - 1 do  $\triangleright$  Drop factor entries in row i for which level is too high 14: if  $level(i,k) > \ell$  then  $\tilde{l}_{ik} = 0$ 15: end for 16: for k = i : n do 17: if  $level(i,k) > \ell$  then  $\tilde{u}_{ik} = 0$ 18: end for 19: end for

### Incomplete factorizations based on shortest fill-paths

- Figure depicts S{*L̃* + *L̃*<sup>T</sup>} for the IC(ℓ) factorization of A from the discretized Laplace equation on a square grid and for a matrix with a more general symmetric sparsity structure.
- The fill-in is typically generated irregularly throughout the factorization: initially few updates are needed but later steps involve many updates, leading to large amounts of dropping.
- The amount of fill-in can grow quickly with increasing  $\ell$  and, as a result, l is typically small.
- Level-based dropping is often combined with threshold-based dropping or with sparsifying *A* before the factorization commences.

### Incomplete factorizations based on shortest fill-paths



IC(0)



### Incomplete factorizations based on shortest fill-paths

- The level-based strategy comes from observing that in practical examples the absolute values of the entries in the factors in positions for which *level* is large are often small. This is the case for model problems arising from discretized PDEs.
- Theoretical understanding follows.

#### Theorem

Consider the  $ILU(\ell)$  factorization of A. level(i, j) = k for some  $k \leq \ell$  if and only if there is a shortest fill path  $i \Longrightarrow j$  of length k + 1 in the adjacency graph  $\mathcal{G}(A)$ .

Incomplete factorizations based on shortest fill-paths

Algorithm (Find the sparsity pattern of row *i* of the ILU( $\ell$ ) factor  $\widetilde{U}$  of *A*: breadth first search)

1:	$\mathcal{S}{\widetilde{U}_{i,i:n}} = {i}, \mathcal{Q} = {i}$	Queue holds i initially
2:	length(i) = 0	
3:	visited(i) = i	
4:	while $Q$ is not empty <b>do</b>	
5:	$pop(\mathcal{Q},k)$	$\triangleright$ Take $k$ from the queue
6:	for $j \in adj_{\mathcal{G}(A)}(k)$ with $visited(j)  eq i$ do	
7:	visited(j) = i	
8:	if $j < i$ and $length(k) < \ell$ then	
9:	$append(\mathcal{Q},j)$	$\triangleright$ Add $j$ to the queue
10.	length(j) = length(k) + 1	
11.	else if $j > i$ then	
12.	$\mathcal{S}\{\widetilde{U}_{i,i:n}\} = \mathcal{S}\{\widetilde{U}_{i,i:n}\} \cup \{j\}$	$ ho$ Add $j$ to the sparsity pattern of row $i$ of $\widetilde{U}$
13.	end if	
14.	end for	
15: end while		

- Assume that the target sparsity pattern  $S{\widetilde{L} + \widetilde{U}}$  contains  $S{A}$ .
- Modified incomplete factorizations (MILU or MIC in the SPD case) seek to maintain equality between the row sums of A and  $\widetilde{L}\widetilde{U}$ , that is,  $\widetilde{L}\widetilde{U}e = Ae$  (e is the vector of all ones).

# Algorithm (Modified incomplete factorization (MILU))

1: Initialise 
$$\tilde{l}_{ij} = (I + L_A)_{ij}$$
 for all  $(i, j) \in S(\tilde{L})$   
2: Initialise  $\tilde{u}_{ij} = (D_A + U_A)_{ij}$  for all  $(i, j) \in S(\tilde{U})$   
3: for  $k = 1 : n - 1$  do  
4: for  $i = k + 1 : n$  such that  $(i, k) \in S\{\tilde{L}\}$  do  
5:  $\tilde{l}_{ik} = \tilde{l}_{ik}/\tilde{u}_{kk}$   $\triangleright$  Check that  $\tilde{u}_{kk}$  is nonzero  
6: for  $j = i : n$  such that  $(k, j) \in S\{\tilde{U}\}$  do  
7: if  $(i, j) \in S\{\tilde{U}\}$  then  
8:  $\tilde{u}_{ij} = \tilde{u}_{ij} - \tilde{l}_{ik} \tilde{u}_{kj}$   
9: else  
10:  $\tilde{u}_{ii} = \tilde{u}_{ii} - \tilde{l}_{ik} \tilde{u}_{kj}$   $\triangleright$  Modify diagonal instead of creating fill-in  
11: end if  
12: end for  
13: for  $j = k + 1 : i - 1$  such that  $(k, j) \in S\{\tilde{U}\}$  do  
14: if  $(i, j) \in S\{\tilde{L}\}$  then  
15:  $\tilde{l}_{ij} = \tilde{l}_{ij} - \tilde{l}_{ik} \tilde{u}_{kj}$   
16: else  
17:  $\tilde{u}_{ii} = \tilde{u}_{ii} - \tilde{l}_{ik} \tilde{u}_{kj}$   $\triangleright$  Modify diagonal instead of creating fill-in  
18: end if  
19: end for  
20: end for

- Equality of row sums: If all the filled entries are retained (that is,  $S{\{\tilde{L}+\tilde{U}\}} = S{\{L+U\}}$ ) then the claim holds trivially.
- Otherwise, if an entry in column j of row i of A belongs to the target sparsity pattern then its value is modified in Step 8 if  $i \leq j$  or in Step 15 if i > j. Otherwise, the *i*-th diagonal entry of  $\tilde{U}$  is modified (Step 10 or Step 17). In each case,  $\tilde{l}_{ik} \tilde{u}_{kj}$  is subtracted from entries of the *i*-th row of the incomplete factors.
- Consider row *i* of  $\widetilde{L}\widetilde{U}$ . This product is given by

$$\begin{split} \sum_{j=1}^{i-1} \tilde{l}_{ij} \sum_{k=j}^{n} \tilde{u}_{jk} &= \sum_{j=1}^{i-1} \tilde{l}_{ij} \tilde{u}_{jj} + \sum_{j=1}^{i-1} \tilde{l}_{ij} \sum_{k=j+1}^{n} \tilde{u}_{jk} + \sum_{k=i}^{n} \tilde{u}_{ik} = \\ &= \sum_{j=1}^{i-1} \left( a_{ij} - \sum_{k=1}^{j-1} \tilde{l}_{ik} \tilde{u}_{kj} \right) + \sum_{j=1}^{i-1} \tilde{l}_{ij} \sum_{k=j+1}^{n} \tilde{u}_{jk} + \sum_{k=i}^{n} \left( a_{ik} - \sum_{j=1}^{i-1} \tilde{l}_{ij} \tilde{u}_{jk} \right) \\ &= \sum_{j=1}^{n} a_{ij} + \sum_{j=1}^{i-1} \tilde{l}_{ij} \sum_{k=j+1}^{n} \tilde{u}_{jk} - \left( \sum_{j=1}^{i-1} \sum_{k=1}^{j-1} \tilde{l}_{ik} \tilde{u}_{kj} + \sum_{k=i}^{n} \sum_{j=1}^{i-1} \tilde{l}_{ij} \tilde{u}_{jk} \right). \end{split}$$

- Rearranging the indices in the double summations, the last three sums cancel out.
- Moreover, the added double summation is the sum of all the modification terms l
   *i*<sub>ik</sub> u
   *i*<sub>kj</sub> in the MILU Algorithm, and the sum of the two subtracted double summations also comprises all the modification terms.
- Consequently, the row sums of *A* are preserved in the product of the incomplete factors.

- Theorem above provides motivation for maintaining constant row sums in the case of a model PDE problem.
- The result is also valid for Neumann or mixed boundary conditions, and there are extensions to three-dimensional problems and MIC(ℓ) with ℓ > 0. However, although Theorem holds for MILU factorizations, the approach may not be useful for general *A*.

#### Theorem

Let *A* come from a discretized Poisson problem on a uniform two-dimensional rectangular grid with Dirichlet boundary conditions and discretization parameter *h*. Then the condition number  $\kappa((\widetilde{L}\widetilde{U})^{-1}A)$  for the level-based MIC(0) preconditioner is  $O(h^{-1})$ .

- RILU/RIC: the update term *l*<sub>ik</sub> *ũ*<sub>kj</sub> may be multiplied by a parameter θ (0 < θ < 1) before it is subtracted from the diagonal entry *ũ*<sub>ii</sub>.
- This is a practical way to extend MILU to linear systems not coming from discretized PDEs. Clearly, using  $\theta < 1$  reduces the amount by which the diagonal entries are modified.

#### Dynamic compensation

- Instead of accepting a filled entry in position (i, j), the idea is to add a (weighted) multiple of its absolute value to the corresponding diagonal entries  $\tilde{u}_{ii}$  and  $\tilde{u}_{jj}$ .
- Provided the number of modifications is small, this can be useful if A is diagonally dominant and scaled so that its diagonal entries are nonnegative.
- The parameter  $\omega$  controls the amount by which the diagonal entries of  $\widetilde{U}$  are modified but if  $\omega < 1$  then breakdown can still occur.
- Dynamic compensation can be successful when incorporated into an IC factorization of a SPD matrix A because the resulting local modifications correspond to adding positive semidefinite matrices to A.
- In practice, the behaviour of the resulting preconditioner can be very different from that computed using the MIC approach.

# Algorithm (ILU factorization with dynamic compensation)

$$\begin{array}{ll} 1: \ \tilde{l}_{ij} = (I + L_A)_{ij} \ \text{for all } (i, j) \in \mathcal{S}(\tilde{L}) \\ 2: \ \tilde{u}_{ij} = (D_A + U_A)_{ij} \ \text{for all } (i, j) \in \mathcal{S}(\tilde{U}) \\ 3: \ \text{for } k = 1: n - 1 \ \text{do} \\ 4: \quad \text{for } i = k + 1: n \ \text{such that } (i, k) \in \mathcal{S}\{\tilde{L}\} \ \text{do} \\ 5: \quad \tilde{l}_{ik} = \tilde{l}_{ik}/\tilde{u}_{kk} \\ 6: \quad \text{for } j = i: n \ \text{such that } (k, j) \in \mathcal{S}\{\tilde{U}\} \ \text{do} \\ 7: \quad \text{if } (i, j) \in \mathcal{S}\{\tilde{U}\} \ \text{then} \\ 8: \quad \tilde{u}_{ij} = \tilde{u}_{ij} - \tilde{l}_{ik} \ \tilde{u}_{kj} \\ 9: \quad \text{else} \\ 10: \quad \rho_{ij} = (\tilde{u}_{ii}/\tilde{u}_{jj})^{1/2} \\ 11: \quad \tilde{u}_{ii} = \tilde{u}_{ii} + \omega \rho_{ij} \ | \tilde{l}_{ik} \ \tilde{u}_{kj} |, \ \tilde{u}_{jj} = \tilde{u}_{jj} + \omega | \tilde{l}_{ik} \ \tilde{u}_{kj} | /\rho_{ij}, \ \tilde{u}_{ij} = 0. \\ 12: \quad \text{end if} \\ 13: \quad \text{end for} \\ 14: \quad \text{for } j = k + 1: i - 1 \ \text{such that } (k, j) \in \mathcal{S}\{\tilde{U}\} \ \text{do} \\ 15: \quad \text{if } (i, j) \in \mathcal{S}\{\tilde{L}\} \ \text{then} \\ 16: \quad \tilde{l}_{ij} = \tilde{l}_{ij} - \tilde{l}_{ik} \ \tilde{u}_{kj} \\ 17: \quad \text{else} \\ 18: \qquad \rho_{ij} = (\tilde{u}_{ii}/\tilde{u}_{jj})^{1/2} \\ 19: \quad \tilde{u}_{ii} = \tilde{u}_{ii} + \omega \rho_{ij} \ | \tilde{l}_{ik} \ \tilde{u}_{kj} |, \ \tilde{u}_{jj} = \tilde{u}_{jj} + \omega | \tilde{l}_{ik} \ \tilde{u}_{kj} | /\rho_{ij}, \ \tilde{l}_{ij} = 0. \\ 20: \quad \text{end if} \end{array}$$

550/609

Dynamic compensation: getting closer to special matrices

- A related scheme, called diagonally compensated reduction, modifies *A* before the factorization begins by adding the values of all of its positive off-diagonal entries to the corresponding diagonal entries and then setting these off-diagonal entries to zero.
- If A is SPD then the resulting matrix is a symmetric M-matrix and the incomplete factorization of an M-matrix is breakdown-free. However, the modified matrix may be too far from A for its incomplete factors to be useful.

### Memory-limited incomplete factorizations

• *A* be SPD, consider the decomposition

$$A = (\widetilde{L} + \widetilde{R}) \, (\widetilde{L} + \widetilde{R})^T - E.$$

- The error matrix E is  $E = \widetilde{R}\widetilde{R}^T$ .
- On step j of the incomplete factorization, the first column of the Schur complement  $S^{(j)}$  is split into the sum

$$\widetilde{L}_{j:n,j} + \widetilde{R}_{j:n,j},$$

where  $\widetilde{L}_{j:n,j}$  contains the entries that are retained in column j of the final incomplete factorization and  $\widetilde{R}_{jj} = 0$  and  $\widetilde{R}_{j+1:n,j}$  contains the entries that are discarded.

#### Memory-limited incomplete factorizations

• If a complete factorization was being computed then the Schur complement would be updated by subtracting

$$(\widetilde{L}_{j+1:n,j} + \widetilde{R}_{j+1:n,j}) (\widetilde{L}_{j+1:n,j} + \widetilde{R}_{j+1:n,j})^T$$

However, the incomplete factorization discards the term

$$E^{(j)} = \widetilde{R}_{j+1:n,j} \,\widetilde{R}_{j+1:n,j}^T.$$

Thus,  $E^{(j)}$  is implicitly added to A and because  $E^{(j)}$  is positive semidefinite, the approach is naturally breakdown-free.

• The obvious choice for  $\widetilde{R}_{j+1:n,j}$  is the smallest off-diagonal entries in the column.

### Memory-limited incomplete factorizations

- Figure depicts the first step of this approach. In the first row and column, \* and  $\delta$  denote the entries of  $\widetilde{L}_{1:n,1}$  and  $\widetilde{R}_{1:n,1}$ , respectively.
- Standard sparsification scheme: no fill (left)
- Using intermediate memory: right.

$$\begin{pmatrix} * & * & * & \delta & \delta \\ * & f & f & & \\ * & f & f & & \\ \delta & & & & & \\ \delta & & & & & \end{pmatrix} \begin{pmatrix} * & * & * & \delta & \delta \\ * & f & f & f & f \\ * & f & f & f & f \\ \delta & f & f & & \\ \delta & f & f & & \\ \end{pmatrix}$$

Figure: An illustration of the fill-in in a standard sparsification-based IC factorization (left) and in the approach that uses intermediate memory (right) after one step of the factorization. Entries with small absolute value in row and column 1 are denoted by  $\delta$ . The filled entries are denoted by *f*.

### Memory-limited incomplete factorizations

• Enables the structure of the complete factorization to be followed more closely than is possible using a standard approach. If the small entries at positions (1,3) and (3,1) are not discarded then there is a filled entry in position (3,2) and this allows the incomplete factorization using intermediate memory to involve the (large) off-diagonal entries in positions (5,2) and (6,2) in the second step of the IC factorization.



## Memory-limited incomplete factorizations

- Unfortunately, because the column  $\widetilde{R}_{j+1:n,j}$  must be retained to perform the updates on the next step, the total memory requirements are essentially as for a complete factorization.
- Relaxations are needed: e.g., introducing two drop tolerances so that only entries of absolute value at least  $\tau_1$  are kept in  $\tilde{L}$  and entries smaller than  $\tau_2$  are dropped from  $\tilde{R}$ .
- Or, limiting the fill-in.
- But, then, no longer breakdown-free approach.

## Algorithm (Crout memory-limited IC factorization)

1: w(1:n) = 0*2:* for j = 1 : n do 3: for i = j : n such that  $a_{ij} \neq 0$  do 4:  $w_i = a_{ij}$ 5: end for 6: for k < j such that  $\tilde{l}_{jk} \neq 0$  do 7: for i = j : n such that  $\tilde{l}_{ik} \neq 0$  do 8:  $w_i = w_i - \tilde{l}_{ik} \, \tilde{l}_{ik}$ 9: end for 10: for i = j : n such that  $\tilde{r}_{ik} \neq 0$  do 11:  $w_i = w_i - \tilde{r}_{ik} \, \tilde{l}_{ik}$ 12: end for 13: end for 14: for k < j such that  $\tilde{r}_{jk} \neq 0$  do 15: for i = j : n such that  $\tilde{l}_{ik} \neq 0$  do 16:  $w_i = w_i - \tilde{l}_{ik} \, \tilde{r}_{ik}$ 17: end for 18: end for 19: Copy into  $L_{j:n,j}$  the  $lsize + nz(A_{j:n,j})$  entries of w of largest absolute value 00. 

 $\triangleright w$  is a vector of length n

557/609

#### Fixed-point iterations for computing ILU factorizations

• Given the target sparsity pattern  $S{\{\tilde{L} + \tilde{U}\}}$ , the goal is to iteratively generate incomplete factors fulfilling the ILU property

$$(\widetilde{L}\widetilde{U})_{ij} = a_{ij}, \quad (i,j) \in \mathcal{S}\{\widetilde{L} + \widetilde{U}\}$$

Parallel computation using the constraints

$$\sum_{\substack{k=1\\i,k),(k,j)\in\mathcal{S}\{\widetilde{L}+\widetilde{U}\}}}^{\min(i,j)} \tilde{l}_{ik}\,\tilde{u}_{kj} = a_{ij}, \quad (i,j)\in\mathcal{S}\{\widetilde{L}+\widetilde{U}\},$$

and the normalization  $\tilde{l}_{ii} = 1$ .

#### Fixed-point iterations for computing ILU factorizations

Using the relations

$$\tilde{l}_{ij} = \left(a_{ij} - \sum_{k=1}^{j-1} \tilde{l}_{ik} \, \tilde{u}_{kj}\right) / \, \tilde{u}_{jj}, \quad i > j,$$

$$\tilde{u}_{ij} = a_{ij} - \sum_{k=1}^{i-1} \tilde{l}_{ik} \, \tilde{u}_{kj}, \quad i \le j,$$
(92)
(93)

the approach can be formulated as a fixed-point iteration method of the form  $w^{k+1} = f(w^k)$ , k = 0, 1, ..., where w is a vector containing the unknowns  $\tilde{l}_{ij}$  and  $\tilde{u}_{ij}$ . Each fixed-point iteration is called a sweep.

## Fixed-point iterations for computing ILU factorizations

# Algorithm (Fixed-point ILU factorization)

Input: Matrix A, the target sparsity pattern  $S{\widetilde{L} + \widetilde{U}}$ , and initial incomplete factors  $\widetilde{L}$  and  $\widetilde{U}$ . Output: Updated incomplete factors.

1: Set  $\tilde{l}_{ij}$  and  $\tilde{u}_{ij}$  to initial values 2: for  $sweep = 1, 2, \dots$  do for  $(i, j) \in \mathcal{S}{\{\widetilde{L} + \widetilde{U}\}}$  do 3: 4: if i > j then 5: Compute  $\tilde{l}_{ij}$  using (92) 6: else 7: Compute  $\tilde{u}_{ij}$  using (93) 8: end if 9: end for 10: end for

# Ordering in incomplete factorizations

- Can have a positive effect on the robustness and performance of preconditioned Krylov subspace methods.
- The best choice of ordering for an incomplete factorization preconditioner may not be the same as for a complete factorization.
- When the natural (lexicographic) ordering is used, the incomplete triangular factors resulting from a no-fill ILU factorization can be highly ill-conditioned, even if the matrix *A* is well conditioned.
- Allowing more fill-in in the factors, for example, using ILU(1) instead of ILU(0), may solve the problem but it is not guaranteed.
- Minimum degree orderings: the rows (and columns) of the permuted matrix can have significantly different counts.
- A strategy is to specify that the permitted fill-in is proportional to the row/column counts of the complete factorization.

# Ordering in incomplete factorizations

- Global orderings cut local connections within the graph of A.
- When used with incomplete factorizations, can lead to poor quality preconditioners.
- A global ordering that specifically targets incomplete factorizations is a red-black (or checker board) ordering.

## Ordering in incomplete factorizations



Figure: A model problem to illustrate a red-black ordering.
## Exploiting block structure

- Blocking methods for complete factorizations can be adapted to incomplete factorizations. The aim is to speed up the computation of the factors and to obtain more effective preconditioners.
- In a block factorization, scalar operations of the form

$$\tilde{l}_{ik} = a_{ik} / \tilde{u}_{kk}$$

are replaced by matrix operations

$$\widetilde{L}_{ib,kb} = A_{ib,kb} \widetilde{U}_{kb,kb}^{-1},$$

and scalar multiplications of entries of the factors are replaced by matrix-matrix products. When dropping entries, instead of considering the absolute values, simple norms of the block entries (such as the one-norm, max-norm or Frobenius norm) are used.

• An incomplete factorization can start with the supernodal structure of the complete factors.

# Outline

- Sparse approximate inverses

# Approximate inverse preconditioners

- Standard solves by substitution steps can present a computational bottleneck. In particular, in parallel computational environment.
- But it is  $M^{-1}$ , which represents an approximation of  $A^{-1}$ , that is applied by performing forward and back substitution steps
- Therefore, an alternative strategy to standard (incomplete) factorizations is to directly approximate  $A^{-1}$  by explicitly computing  $M^{-1}$ .

Approximate inverse preconditioners

- But, there is a problem: The sparsity pattern of the inverse of an irreducible matrix *A* is dense, even when *A* is sparse.
- But, perhaps there is a way ...: although  $A^{-1}$  is fully dense, the following result shows this is not the case for the factors of factorized inverses.

# Theorem

Assume the matrix A is SPD and let L be its Cholesky factor. Then  $S\{L^{-1}\}$  is the union of all entries (i, j) such that i is an ancestor of j in the elimination tree T(A).

- A consequence of this result is that  $L^{-1}$  need not be fully dense.
- Algorithmically, if A is SPD it may be advantageous to preorder A to limit the number of ancestors of vertices in  $\mathcal{T}(A)$ .
- For example, by ND applied to  $S{A}$  or to  $S{A + A^T}$ .

## Basic approaches

- An obvious way: to compute an incomplete LU factorization of *A* and then perform an approximate inversion of the incomplete factors.
- But, two levels of approximation.

#### Basic approaches

- Another straightforward approach is based on bordering.
- Let  $A_j = A_{1:j,1:j}$  and its inverse factorization

$$A_j^{-1} = W_j D_j^{-1} Z_j^T$$

is known ( $W_j$  and  $Z_j$  are unit upper triangular matrices and  $D_j$  is a diagonal matrix).

$$\begin{pmatrix} Z_j^T & 0 \\ z_{j+1}^T & 1 \end{pmatrix} \begin{pmatrix} A_j & A_{1:j,j+1} \\ A_{j+1,1:j} & a_{j+1,j+1} \end{pmatrix} \begin{pmatrix} W_j & w_{j+1} \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} D_j & 0 \\ 0 & d_{j+1,j+1} \end{pmatrix},$$
where for  $1 \le j < n$ 

$$w_{j+1} = -W_j D_j^{-1} Z_j^T A_{1:j,j+1}, z_{j+1} = -Z_j D_j^{-1} W_j^T A_{j+1,1:j}^T,$$
  
$$d_{j+1,j+1} = a_{j+1,j+1} + z_{j+1}^T A_j w_{j+1} + A_{j+1,1:j} w_{j+1} + z_{j+1}^T A_{1:j,j+1}.$$
  
Starting from  $i = 1$  this suggests a procedure for computing the

 Starting from j = 1, this suggests a procedure for computing the inverse factors of A. Sparsity can be preserved by dropping.

## Approximate inverse by bordering

# Algorithm (Nonsymmetric inverse bordering algorithm)

Input: Generally nonsymmetric A. Output: A unit upper triangular matrix  $\widetilde{Z}$  and diagonal matrix  $\widetilde{D}$  such that  $A^{-1} \approx \widetilde{Z}^T \widetilde{D}^{-1} \widetilde{W}^T$ .

 $\begin{array}{ll} 1: \ {\rm Set}\,\widetilde{Z}_1 = (1),\,\widetilde{W}_1 = (1),\,\widetilde{D}_1 = (a_{11}).\\ 2: \ {\rm for}\, j = 2:n \ {\rm do}\\ 3: & \ {\rm Set}\,\widetilde{z}_j = -\widetilde{Z}_{1:j-1,1:j-1}\widetilde{D}_{1:j-1,1:j-1}^{-1}\widetilde{W}_{1:j-1,1:j-1}^TA_{1:j-1,j}\\ 4: & \ {\rm Set}\,\widetilde{w}_j = -\widetilde{W}_{1:j-1,1:j-1}\widetilde{D}_{1:j-1,1:j-1}^{-1}\widetilde{Z}_{1:j-1,1:j-1}^TA_{j,1:j-1}^T\\ 5: & \ {\rm Set}\,\widetilde{d}_j = A_{jj} + A_{1:j-1,j}^T\widetilde{w}_j + A_{j,1:j-1}\widetilde{z}_j + \widetilde{w}_j^TA_{j-1}\widetilde{z}_j\\ 6: & \ {\rm Set}\,\widetilde{Z}_j = \begin{pmatrix} \widetilde{Z}_{1:j-1,1:j-1} & \widetilde{z}_j\\ 0 & 1 \end{pmatrix}\\ 7: & \ {\rm Set}\,\widetilde{W}_j = \begin{pmatrix} \widetilde{W}_{1:j-1,1:j-1} & \widetilde{w}_j\\ 0 & 1 \end{pmatrix}\\ 8: \ {\rm end}\ {\rm for}\\ 9: \ {\rm Set}\,\widetilde{Z} = \widetilde{Z}_n, \widetilde{W} = \widetilde{W}_n, \widetilde{D} = \widetilde{D}_n \end{array}$ 

## Inverse by bordering: notes

- If A is symmetric, W = Z and the required work is halved.
- Furthermore, if A is SPD then it can be shown that, in exact arithmetic,  $d_{jj} > 0$  for all j and the process does not break down.
- The computation of *Z* and *W* are tightly coupled restricting the potential to exploit parallelism.
- This implies a potential problem with efficient implementation.

# Frobenius norm minimization: SPAI

- Denote  $K = M^{-1}$ .
- Use minimization of

$$||I - AM^{-1}||_F^2 = ||I - AK||_F^2 = \sum_{j=1}^n ||e_j - Ak_j||_2^2,$$
(94)

over all K with pattern S.

- A left approximate inverse can be computed by solving a minimization problem for  $||I KA||_F = ||I A^T K^T||_F$ .
- The problem reduces to least squares problems for the columns of *K* that can be computed independently and, if required, in parallel.

### Frobenius norm minimization: SPAI

- These least squares (LS) problems are all of small dimension when S is chosen to ensure K is sparse.
- Let  $\mathcal{J} = \{i \mid k_j(i) \neq 0\}$  be the set of indices of the nonzero entries in column  $k_j$ . Further, denote  $\mathcal{I} = \{m \mid A_{m,\mathcal{J}} \neq 0\}$ .
- Let \$\heta\_j = e\_j(\mathcal{I})\$ be the vector of length \$|\mathcal{I}|\$ that is obtained by taking the entries of \$e\_j\$ with row indices belonging to \$\mathcal{I}\$.
- To solve the LS problem for  $k_j$ , construct the  $|\mathcal{I}| \times |\mathcal{J}|$  matrix  $\widehat{A} = A_{\mathcal{I},\mathcal{J}}$  and solve

$$\min_{\widehat{k}_j} \|\widehat{e}_j - \widehat{A}\,\widehat{k}_j\|_2^2. \tag{95}$$

This can be done using QR factorization of Â. Extending k<sub>j</sub> to have length n by zeros gives k<sub>j</sub>.

## SPAI algorithm

• Construction: starting with a chosen column sparsity pattern  $\mathcal{J}$  for  $k_j$ , construct  $\widehat{A}$ , solve (95) for  $\widehat{k}_j$ , set  $k_j(\mathcal{J}) = \widehat{k}_j$  and define the residual vector

$$r_j = e_j - A_{1:n,\mathcal{J}}\widehat{k}_j.$$

- If ||r<sub>j</sub>||<sub>2</sub> ≠ 0 then k<sub>j</sub> is not equal to the *j*-th column of A<sup>-1</sup> and a better approximation can be derived by augmenting J.
- Augmentation: let  $\mathcal{L} = \{l \mid r_j(l) \neq 0\}$  and define

$$\widetilde{\mathcal{J}} = \{i \,|\, A_{\mathcal{L},i} \neq 0\} \setminus \mathcal{J}.$$
(96)

 One or more candidate indices that can be added to *J* can be chosen. For example, such that most effectively reduce ||r<sub>j</sub>||<sub>2</sub>.

# SPAI algorithm

• A possible heuristic is to solve for each  $i \in \widetilde{\mathcal{J}}$  the minimization problem

$$\min_{\mu_i} ||r_j - \mu_i A e_i||_2^2.$$

- This has the solution  $\mu_i = r_j^T A e_i / ||Ae_i||_2^2$  with residual  $||r_j||^2 (r_j^T A e_i)^2 / ||Ae_i||_2^2$ .
- Indices  $i \in \widetilde{\mathcal{J}}$  for which this is small are appended to  $\mathcal{J}$ .
- The process can be repeated until either the required accuracy is attained or the maximum number of allowed entries in  $\mathcal{J}$  is reached.

#### SPAI algorithm

- Solving the unconstrained LS problem after extending  $\widehat{A}$  to  $A_{\mathcal{I}\cup\mathcal{I}',\mathcal{J}\cup\mathcal{J}'}$  is typically performed by updating the previous problems.
- Assume the QR factorization of  $\widehat{A}$  is

$$\widehat{A} = A_{\mathcal{I},\mathcal{J}} = Q \begin{pmatrix} R \\ 0 \end{pmatrix} = \begin{pmatrix} Q_1 & Q_2 \end{pmatrix} \begin{pmatrix} R \\ 0 \end{pmatrix},$$

where  $Q_1$  is  $|\mathcal{I}| \times |\mathcal{J}|$ .

The QR factorization of the extended matrix is

$$\begin{aligned} A_{\mathcal{I}\cup\mathcal{I}',\mathcal{J}\cup\mathcal{J}'} &= \begin{pmatrix} \widehat{A} & A_{\mathcal{I},\mathcal{J}'} \\ & A_{\mathcal{I}',\mathcal{J}'} \end{pmatrix} = \begin{pmatrix} Q \\ & I \end{pmatrix} \begin{pmatrix} R & Q_1^T A_{\mathcal{I},\mathcal{J}'} \\ & Q_2^T A_{\mathcal{I},\mathcal{J}'} \\ & A_{\mathcal{I}',\mathcal{J}'} \end{pmatrix} \\ &= \begin{pmatrix} Q \\ & I \end{pmatrix} \begin{pmatrix} I \\ & Q' \end{pmatrix} \begin{pmatrix} R & Q_1^T A_{\mathcal{I},\mathcal{J}'} \\ & R' \\ & 0 \end{pmatrix}. \end{aligned}$$

# SPAI algorithm

• Q' and R' are from the QR factorization of the  $(|\mathcal{I}'| + |\mathcal{I}| - |\mathcal{J}|) \times |\mathcal{J}'|$  matrix

$$\begin{pmatrix} Q_2^T A_{\mathcal{I},\mathcal{J}'} \\ A_{\mathcal{I}',\mathcal{J}'} \end{pmatrix}.$$

- Factorizing this matrix and updating the trailing QR factorization to get the new  $\hat{k}_j$  is much more efficient than computing the QR factorization of the extended matrix from scratch.
- Many variations of the basic approach.

## SPAI algorithm

# Algorithm (SPAI preconditioner (right-looking approach))

Input: Nonsymmetric matrix A, a convergence tolerance  $\eta > 0$ , an initial sparsity pattern  $\mathcal{J}_j$  and the maximum number  $nz_j$  of permitted entries for column j of K  $(1 \le j \le n)$ . Output:  $K \approx A^{-1}$  with columns  $k_j$   $(1 \le j \le n)$ .

1: for j = 1 : n do > The columns may be computed in parallel 2: Set  $\mathcal{J} = \mathcal{J}_i$  and  $\mathcal{I} = \{m \mid A(m, \mathcal{J}) \neq 0\}, \|r_i\|_2 = \infty$ 3: Construct  $\widehat{A} = A_{\mathcal{I},\mathcal{J}}$  and solve (95) for  $\widehat{k}_i$ 4:  $r_i = e_i - A_{1:n_{i}\mathcal{T}}\widehat{k}_i$ 5: while  $|\mathcal{J}| < nz_j$  and  $||r_j||_2 > \eta$  do 6:  $\triangleright \widetilde{\mathcal{T}}$  is the candidate set Construct  $\mathcal{J}$  given by (96) 7: Determine new indices  $\mathcal{J}' \subset \widetilde{\mathcal{J}}$  to add to  $\mathcal{J}$ 8:  $\mathcal{I}' = \{m \mid A_{m,\mathcal{I}'} \neq 0\} \setminus \mathcal{I}$ 9:  $\mathcal{I} = \mathcal{I} \cup \mathcal{I}'$  and  $\mathcal{J} = \mathcal{J} \cup \mathcal{J}'$ Augment the sparsity pattern 10: Construct new  $\widehat{A} = A_{\mathcal{I},\mathcal{I}}$  and new  $\widehat{k}_i$ Update the QR factorization 11:  $r_i = e_i - A_{1:n,\mathcal{T}} \widehat{k}_i$ 12: end while  $k_i(\mathcal{J}) = \widehat{k}_i$ 13:  $\triangleright$  Extend  $\hat{k}_i$  to  $k_i$  by setting entries not in  $\mathcal{J}$  to zero. 14: end for 578/609

#### SPAI algorithm

The example: The algorithm starts with  $\mathcal{J}_1 = \{1, 2\}$ .

$$A = \begin{pmatrix} 10 & -2 \\ -1 & 10 & -2 \\ & -1 & 10 & -2 \\ & & -1 & 10 & -2 \\ & & & -1 & 10 \end{pmatrix}, \ \widehat{A} = \begin{pmatrix} 10 & -2 \\ -1 & 10 \\ & & -1 \end{pmatrix}, \ \widehat{k}_1 = \begin{pmatrix} 0.1020 \\ 0.0101 \end{pmatrix}, \ r_1 = \begin{pmatrix} 1.00 \times 10^{-4} \\ 1.00 \times 10^{-3} \\ 1.01 \times 10^{-2} \\ 0 \\ 0 \end{pmatrix}$$

$$\widehat{A} = \begin{pmatrix} 10 & -2 \\ -1 & 10 & -2 \\ & -1 & 10 \\ & & -1 \end{pmatrix}, \ \widehat{k}_1 = \begin{pmatrix} 0.1021 \\ 0.0104 \\ 0.0010 \end{pmatrix}, \ r_1 = \begin{pmatrix} 1.0 \times 10^{-5} \\ 1.1 \times 10^{-4} \\ 1.1 \times 10^{-3} \\ 1.0 \times 10^{-2} \\ 0 \end{pmatrix}, \ k_1 = \begin{pmatrix} 0.1021 \\ 0.0104 \\ 0.0010 \\ 0 \\ 0 \end{pmatrix}.$$

Figure: An illustration of computing the first column of a sparse approximate inverse using the SPAI algorithm with  $nz_1 = 3$ . On the top line is the initial tridiagonal matrix A followed by the matrix  $\hat{A}$  and the vectors  $\hat{k}_1$  and  $r_1$  on the first loop of Algorithm. The bottom line presents the updated matrix  $\hat{A}$  that is obtained on the second loop by adding the third row and column of A and the corresponding vectors  $\hat{k}_1$  and  $r_1$  and, finally,  $k_1$ .

# SPAI algorithm

- When A is symmetric, there is no guarantee that the computed K will be symmetric. One possibility is to use  $(K + K^T)/2$  to approximate  $A^{-1}$ .
- The SPAI preconditioner is not sensitive to reorderings of *A*. This has the advantage that *A* can be partitioned and reordered in whatever way is convenient, for instance to better suit the needs of a distributed implementation.
- The disadvantage is that orderings cannot be used to reduce fill-in and/or improve the quality of the approximate inverse.
- For instance, if  $A^{-1}$  has no small entries, SPAI will not find a sparse *K*, and because the inverse of a permutation of *A* is just a permutation of  $A^{-1}$ , no permutation of *A* will change this.

• The factorized sparse approximate inverse (FSAI) preconditioner for an SPD matrix *A* is defined as the product

$$M^{-1} = G^T G,$$

where the sparse lower triangular matrix G is an approximation of the inverse of the (complete) Cholesky factor L of A.

• Theoretically, a FSAI preconditioner is computed by choosing a lower triangular sparsity pattern  $S_L$  and minimizing

$$||I - GL||_F^2 = tr \left[ (I - GL)^T (I - GL) \right]$$
(97)

over all G with sparsity pattern  $S_L$ .

 Differentiating the formula with respect to the entries of G and setting to zero yields

$$(GLL^T)_{ij} = (GA)_{ij} = (L^T)_{ij}$$
 for all  $(i,j) \in \mathcal{S}_L$ . (98)

• Because  $L^T$  is an upper triangular matrix while  $S_L$  is a lower triangular pattern, the matrix equation (98) can be rewritten as

$$(GA)_{ij} = \begin{cases} 0 & i \neq j, \quad (i,j) \in \mathcal{S}_L \\ l_{ii} & i = j. \end{cases}$$
(99)

• *G* is not available directly because *L* is unknown. Instead,  $\overline{G}$  is computed such that

$$(\overline{G}A)_{ij} = \delta_{i,j}$$
 for all  $(i,j) \in \mathcal{S}_L$ , (100)

where  $\delta_{i,j}$  is the Kronecker delta function ( $\delta_{i,j} = 1$  if i = j and is equal to 0, otherwise).

• The FSAI factor G is then obtained by setting

 $G = D\overline{G},$ 

where D is a diagonal scaling matrix.

• An appropriate choice for D is

$$D = [diag(\overline{G})]^{-1/2},$$
(101)

so that

$$(GAG^T)_{ii} = 1, \quad 1 \le i \le n.$$

#### Theorem

Assume A is SPD. If the lower triangular sparsity pattern  $S_L$  includes all diagonal positions then G exists and is unique.

#### Proof.

Set  $\mathcal{I}_i = \{j \mid (i, j) \in \mathcal{S}_L\}$  and let  $A_{\mathcal{I}_i, \mathcal{I}_i}$  denote the submatrix of order  $nz_i = |\mathcal{I}_i|$  of entries  $a_{kl}$  such that  $k, l \in \mathcal{I}_i$ . Let  $\bar{g}_i$  and  $g_i$  be dense vectors containing the nonzero coefficients in row *i* of  $\overline{G}$  and G, respectively. Using this notation, solving (100) decouples into solving *n* independent SPD linear systems

$$A_{\mathcal{I}_i, \mathcal{I}_i} \,\bar{g}_i = e_{nz_i}, \quad 1 \le i \le n,$$

where the unit vectors are of length  $nz_i$ . Moreover,

$$(\overline{G}A\overline{G}^T)_{ii} = \sum_{j \in \mathcal{I}_i} \delta_{i,j}\overline{G}_{ij} = \overline{G}_{ii} = (A_{\mathcal{I}_i,\mathcal{I}_i}^{-1})_{ii}.$$

This implies that the diagonal entries of D given by (101) are nonzero. Consequently, the computed rows of G exist and provide a unique solution.

# Algorithm (FSAI preconditioner)

Input: SPD matrix A and lower triangular sparsity pattern  $S_L$  that includes all diagonal positions. Output: Lower triangular matrix G such that  $A^{-1} \approx GG^T$ .

1: for i = 1 : n do 2: Construct  $\mathcal{I}_i = \{j \mid (i, j) \in S_L\}$ ,  $A_{\mathcal{I}_i, \mathcal{I}_i}$  and set  $nz_i = |\mathcal{I}_i|$ 3: Solve  $A_{\mathcal{I}_i, \mathcal{I}_i} \bar{g}_i = e_{nz_i}$ 4: Scale  $g_i = d_{ii}\bar{g}_i$  with  $d_{ii} = (\bar{g}_{i,nz_i})^{-1/2} \qquad \triangleright \bar{g}_{i,nz_i}$  is the last component of  $\bar{g}_i$ 5: Extend  $g_i$  to the row  $G_{i,1:i}$  by setting entries that are not in  $\mathcal{I}_i$  to zero 6: end for

Monotonicity property.

#### Theorem

Let *L* be the Cholesky factor of the SPD matrix *A*. Given the lower triangular sparsity pattern  $S_L$  that includes all diagonal positions, let *G* be the FSAI preconditioner computed using Algorithm above. Then any lower triangular matrix  $G_1$  with its sparsity pattern is contained in  $S_L$  a  $(G_1AG_1^T)_{ii} = 1$  ( $1 \le i \le n$ ) satisfies

 $||I - GL||_F \le ||I - G_1L||_F.$ 

• The performance is highly dependent on the choice of  $S_L$ .

#### Theorem

Let *L* be the Cholesky factor of the SPD matrix *A*. Given the lower triangular sparsity patterns  $S_{L1}$  and  $S_{L2}$  that include all diagonal positions, let the corresponding FSAI preconditioners computed using Algorithm 14.3 be  $G_1$  and  $G_2$ , respectively. If  $S_{L1} \subseteq S_{L2}$  then

 $||I - G_2L||_F \le ||I - G_1L||_F.$ 

# FSAI preconditioner: general case

- The FSAI algorithm can be extended to a general matrix *A*. Two input sparsity patterns are required.
- First, lower and upper triangular matrices  $\overline{G}_L$  and  $\overline{G}_U$  are computed such that

$$(\overline{G}_L A)_{ij} = \delta_{i,j}$$
 for all  $(i,j) \in \mathcal{S}_L$ ,

$$(A\overline{G}_U)_{ij} = \delta_{i,j}$$
 for all  $(i,j) \in \mathcal{S}_U$ .

- Then *D* is obtained as the inverse of the diagonal of the matrix  $\overline{G}_L A \overline{G}_U$ , and the final nonsymmetric FSAI factors are given by  $G_L = \overline{G}_L$  and  $G_U = \overline{G}_U D$ . The computation of the two approximate factors can be performed independently.
- This generalization is well defined if, for example, A is nonsymmetric positive definite. There is also theory that extends existence to special classes of matrices, including M- and H-matrices.

Determining a good sparsity pattern

- Input pattern is expected to filter out entries of A<sup>-1</sup> that contribute little to the quality of the preconditioner.
- For instance, it might be appropriate to ignore entries with a small absolute value, while retaining the largest ones. But, locations of large entries in A<sup>-1</sup> are generally unknown, and this makes the a priori sparsity choice difficult.
- A a banded SPD matrix: the entries of A<sup>-1</sup> are bounded in an exponentially decaying manner along each row or column: there exist 0 < ρ < 1 and a constant c such that for all i, j</li>

$$|(A^{-1})_{ij}| \le c\rho^{|i-j|}.$$

The scalars  $\rho$  and c depend on the bandwidth and  $\kappa(A)$ .

- A common choice for a general A is  $S_L + S_U = S\{A\}$ , .
- An alternative strategy uses the Neumann series expansion of  $A^{-1}$ : using the pattern of a small power of A, i.e.,  $S\{A^2\}$  or  $S\{A^3\}$ .

Factorized approximate inverses based on incomplete conjugation

- An alternative way: using incomplete conjugation (*A*-orthogonalization) in the SPD case and on incomplete *A*-biconjugation in the general case. For SPD matrices, the approach represents an approximate Gram-Schmidt orthogonalization that uses the *A*-inner product  $\langle ., . \rangle_A$ .
- Sparsity pattern not needed in advance.
- When *A* is a SPD matrix the AINV preconditioner is defined in the form

$$A^{-1} \approx M^{-1} = ZD^{-1}Z^T,$$

 ${\it Z}$  is unit upper triangular,  ${\it D}$  is a diagonal matrix with positive entries.

• Practical implementations need to employ sparse matrix techniques. The left-looking scheme computes the *j*-th column  $z_j$  of *Z* as a sparse linear combination of the previous columns  $z_1, \ldots, z_{j-1}$ . The key is determining which multipliers (the  $\alpha$ 's in Steps 4 and 5 of the two algorithms, respectively) are nonzero and

## Factorized approximate inverses based on incomplete conjugation

# Algorithm (AINV preconditioner (left-looking approach))

*Input:* SPD matrix A and sparsifying rule. *Output:*  $A^{-1} \approx ZD^{-1}Z^T$  with Z a unit upper triangular matrix and D a diagonal matrix with positive diagonal entries.



# Factorized approximate inverses based on incomplete conjugation

# Algorithm (AINV preconditioner (right-looking approach))

*Input:* SPD matrix A and sparsifying rule. *Output:*  $A^{-1} \approx ZD^{-1}Z^T$  with Z a unit upper triangular matrix and D a diagonal matrix with positive diagonal entries.



# AINV preconditioner: general case

 In the general case, the AINV preconditioner is given by an approximate inverse factorization of the form

$$A^{-1} \approx M^{-1} = W D^{-1} Z^T,$$

where Z and W are unit upper triangular matrix and D is a diagonal matrix.

• *Z* and *W* are sparse approximations of the inverses of the *L*<sup>*T*</sup> and *U* factors in the LDU factorization of *A*, respectively.

#### AINV preconditioner: general case

Algorithm (Nonsymmetric AINV preconditioner (right-looking approach))

1:  $[z_1^{(0)}, \ldots, z_n^{(0)}] = [e_1, \ldots, e_n]$  and  $[w_1^{(0)}, \ldots, w_n^{(0)}] = [e_1, \ldots, e_n]$ 2: for j = 1 : n do  $d_{ij} = (A_{1:n,j})^T z_i^{(j-1)}$  or  $d_{ij} = A_{i,1:n} w_i^{(j-1)}$ 3: 4: for k = j + 1 : n do  $\alpha = (A_{1:n,j})^T z_k^{(j-1)} / d_{jj}$ 5:  $z_{h}^{(j)} = z_{h}^{(j-1)} - \alpha z_{i}^{(j-1)}$ 6: Sparsify  $z_{L}^{(j)}$ 7:  $\triangleright$  Drop entries from  $z_{L}^{(j)}$  $\beta = A_{i,1:n} w_{i}^{(j-1)} / d_{ii}$ 8:  $w_{L}^{(j)} = w_{L}^{(j-1)} - \beta w_{i}^{(j-1)}$ 9: Sparsify  $w_{L}^{(j)}$  $\triangleright$  Drop entries from  $w_{L}^{(j)}$ 10: 11: end for 12: end for 13:  $Z = [z_1^{(0)}, \ldots, z_n^{(n-1)}]$  and  $W = [w_1^{(0)}, \ldots, w_n^{(n-1)}]$ 

## AINV preconditioner

• Matrix A, AINV preconditioner



# AINV preconditioner

# ILUT, inverse ILUT



## SAINV: stabilization of the AINV method

• The following result is analogous to the SPD case.

#### Theorem

If A is a nonsingular M- or H-matrix then the AINV factorization of A does not break down.

- For more general matrices breakdown can happen because of the occurrence of zero  $d_{jj}$  or, in the SPD case, negative  $d_{jj}$ .
- In practice, exact zeros are unlikely but very small d<sub>jj</sub> can occur (near breakdown), which may lead to uncontrolled growth in the size of entries in the incomplete factors and, because such entries are not dropped when using a threshold parameter, a large amount of fill-in.
- The next theorem indicates how breakdown can be prevented when *A* is SPD through reformulating the *A*-orthogonalization.

### SAINV: stabilization of the AINV method

#### Theorem

Consider AINV algorithm with no sparsification (Step 7 is removed). The following holds

$$A_{j,1:n} z_k^{(j-1)} \equiv e_j^T A z_k^{(j-1)} = \langle z_j^{(j-1)}, z_k^{(j-1)} \rangle_A, \ 1 \le j \le k \le n.$$

## Proof.

Because  $AZ = Z^{-T}D$  and  $Z^{-T}D$  is lower triangular, entries 1 to j - 1 of the vector  $Az_k^{(j-1)}$  are equal to zero. Z is unit upper triangular so entries j + 1 to n of its j-th column  $z_j^{(j-1)}$  are also equal to zero. Thus,  $z_j^{(j-1)}$  can be written as the sum  $z + e_j$ , where entries j to n of th vector z are zero. The result follows.  $\Box$ 

# SAINV: stabilization of the AINV method

# Algorithm (SAINV preconditioner (right-looking approach))

*Input:* SPD matrix A and sparsifying rule. *Output:*  $A^{-1} \approx ZD^{-1}Z^T$  with Z a unit upper triangular matrix and D a diagonal matrix with positive diagonal entries.

1:  $[z_1^{(0)}, \dots, z_n^{(0)}] = [e_1, \dots, e_n]$ 2: for j = 1 : n do 3:  $d_{jj} = \langle z_j^{(j-1)}, z_j^{(j-1)} \rangle_A$ 4: for k = j + 1 : n do 5:  $\alpha = \langle z_k^{(j-1)}, z_j^{(j-1)} \rangle_A / d_{jj}$ 6:  $z_k^{(j)} = z_k^{(j-1)} - \alpha z_j^{(j-1)}$ 7: Sparsify  $z_k^{(j)}$ 8: end for 9: end for 10:  $Z = [z_1^{(0)}, \dots, z_n^{(n-1)}]$ 

 $\triangleright$  Drop entries from  $z_k^{(j)}$
#### From AINV to Cholesky

- The factors *Z* and *D* obtained with no sparsification can be used to compute the square root-free Cholesky factorization of *A*.
- The *L* factor of *A* and the inverse factor *Z* computed using AINV Algorithm without sparsification satisfy

$$AZ = LD$$
 or  $L = AZD^{-1}$ .

• Using  $d_{jj} = \langle z_j^{(j-1)}, z_j^{(j-1)} \rangle_A$ , and equating corresponding entries of  $AZD^{-1}$  and L gives

$$l_{ij} = \frac{\langle z_j^{(j-1)}, \, z_i^{(j-1)} \rangle_A}{\langle z_j^{(j-1)}, \, z_j^{(j-1)} \rangle_A}, \quad 1 \le j \le i \le n.$$

Thus, the SAINV algorithm generates the L factor of the square root-free Cholesky factorization of A as a by-product of orthogonalization in the inner product (.,.)<sub>A</sub> at no extra cost and without breakdown. Stabilization-like for general A

• The stabilization strategy can be extended to the nonsymmetric AINV algorithm using the following result.

#### Theorem

Consider nonsymmetric AINV Algorithm with no sparsification (Steps 7 and 10 removed). The following identities hold:

$$A_{j,1:n} z_k^{(j-1)} \equiv e_j^T A z_k^{(j-1)} = \langle w_j^{(j-1)}, \, z_k^{(j-1)} \rangle_A,$$

$$(A_{1:n,j})^T w_k^{(j-1)} \equiv e_j^T A^T w_k^{(j-1)} = \langle z_j^{(j-1)}, w_k^{(j-1)} \rangle_A, \quad 1 \le j \le k \le n.$$

 The nonsymmetric SAINV algorithm obtained using this reformulation can improve the preconditioner quality but it is not guaranteed to be breakdown free. Approximate inverse by global iterations

Consider one-dimensional Newton-Raphson iterations to find a scalar value p which is the root of a given function f, that is

$$f(p) = 0.$$

The method approaches p by a sequence of approximations  $p_0, p_1, \ldots$ . Consider a tangent of f at  $p_k$  for some integer  $k \ge 0$  in the following form

$$y = f'(p_k)p_k + b.$$
 (102)

The tangent crosses  $(p_k, f(p_k))$  and this can be put down as

$$f(p_k) = f'(p_k)p_k + b.$$
 (103)

This implies

$$b = f(p_k) - f'(p_k)p_k$$
 (104)

and we get a function of x given by

$$y = f'(x)x + f(p_k) - f'(p_k)p_k.$$
 (105)

602/609

# Approximate factorizations, splitting and preconditioning

Assume that the root is achieved at  $p_{k+1}$ . Then

$$0 = f'(p_{k+1})p_{k+1} + f(p_k) - f'(p_k)p_k$$
(106)

and therefore

$$p_{k+1} = p_k - \frac{f(p_k)}{f'(p_k)}.$$
 (107)

For f beeing the function of the inverse given by

$$f(x) = 1/x - a$$

we have

$$p_{k+1} = p_k - \frac{1/p_k - a}{-1/p_k^2} = p_k(2 - ap_k).$$
 (108)

Finding the matrix inverse in case it is well-defined:

$$G_{i+1} = G_i(2I - AG_i), i = 1, \dots$$

for the sequence of non-factorized approximate inverses  $G_0, \ldots$ . The main problem: *G* is for irreducible *A* fully dense.

## Approximate factorizations, splitting and preconditioning

Consider the computation of the j-th diagonal entry and assume the exactly computed quantities. The computation from the formula

$$d_j = A_{jj} + A_{1:j-1,j}^T w_j + A_{j,1:j-1} z_j + w_j^T A_{j-1} z_j$$
(109)

can be easily replaced by the mathematically equivalent formula which we used in the algorithms using biconjugation computation:

$$d_j = A_{j,1:j-1} z_j \text{ or } d_j = A_{1:j-1,j}^T w_j.$$
 (110)

#### 1. Shared memory computers

- 1st level of parallelism: tree structure of the decomposition.
- 2nd level of parallelism: local node parallel enhancements.



- Both may/should be coordinated.
- **Tree parallelism** potential decreases towards its root.
- Potential for the local parallelism (larger dense matrices) increases towards the root.

### Decomposition and computer architectures: 1st level of parallelism

#### Two basic possibilities for the tree parallelism

- Dynamic task scheduling on shared memory computers
- Direct static mapping: subtree to subcube
  - 1. Dynamic task scheduling on shared memory computers
- Dynamic scheduling of the tasks
- Each processor selects a task
- Again, problem of elimination tree reordering
- Not easy to optimize memory, e.g., in the multifrontal method



## Decomposition and computer architectures: 1st level of parallelism: II

#### 2. Direct static mapping: subtree to subcube

- Recursively map processors to the tree parts from the top
- Various ways of mapping.
- Note: In the SPD (non-pivoting) case the arithmetic work can be computed and considered
- Localized communication
- More difficult to share the work among processors in more complex models



### Decomposition and computer architectures: 2nd level of parallelism

- Block Cholesky/LU factorization
- BLAS / parallel BLAS operations



1D and 2D block cyclic distribution

(Only illustrative figures for the talk!)

### Decomposition and computer architectures: Distributed memory parallelism

Basic classical parallelization approaches (consider Cholesky)

- Fan-in approach
  - Demand-driven column-based algorithm
  - Required data are aggregated updates asked from previous columns
- o bf Fan-out approach
  - Data-driven column-based algorithm
  - Updates are broadcasted once computed and aggregated
  - Historically the first approach; greater interprocessor communication than fan-in
- Multifrontal approach
  - Example: MUMPS