

Gaussian Elimination

In its simplest form, Gaussian elimination proceeds much like a reduction to echelon form.

Example 1: Solve the system
$$\begin{cases} x_1 - 2x_2 + 2x_3 = 1 \\ -x_1 - x_2 + 3x_3 = 1 \\ 3x_1 + x_2 - 2x_3 = 2 \end{cases}$$

Writing down only the coefficients, and proceeding by removing the entries below the main diagonal in successive columns gives

$$\left[\begin{array}{ccc|c} 1 & -2 & 2 & 1 \\ -1 & -1 & 3 & 1 \\ 3 & 1 & -2 & 2 \end{array} \right] \Rightarrow \left\{ \begin{array}{l} \text{add multiples 1 and} \\ -3 \text{ resp. of top row} \\ \text{to two rows below} \end{array} \right\} \Rightarrow \left[\begin{array}{ccc|c} 1 & -2 & 2 & 1 \\ & -3 & 5 & 2 \\ & 7 & -8 & -1 \end{array} \right] \Rightarrow \left\{ \begin{array}{l} \text{multiply} \\ \text{second row} \\ \text{by } -\frac{1}{3} \end{array} \right\} \Rightarrow \left[\begin{array}{ccc|c} 1 & -2 & 2 & 1 \\ & 1 & -\frac{5}{3} & -\frac{2}{3} \\ & 7 & -8 & -1 \end{array} \right] \Rightarrow$$

$$\Rightarrow \left\{ \begin{array}{l} \text{add -7 times} \\ \text{second row to} \\ \text{third row} \end{array} \right\} \Rightarrow \left[\begin{array}{ccc|c} 1 & -2 & 2 & 1 \\ & 1 & -\frac{5}{3} & -\frac{2}{3} \\ & & \frac{11}{3} & \frac{11}{3} \end{array} \right] \Rightarrow \left\{ \begin{array}{l} \text{multiply} \\ \text{last row} \\ \text{by } \frac{3}{11} \end{array} \right\} \Rightarrow \left[\begin{array}{ccc|c} 1 & -2 & 2 & 1 \\ & 1 & -\frac{5}{3} & -\frac{2}{3} \\ & & 1 & 1 \end{array} \right].$$

Back substitution now proceeds by noting that the last equation tells us that $x_3 = 1$; knowing this, middle equation tells that $x_2 = 1$, and top equation that also $x_1 = 1$.



This basic elimination process is not quite satisfactory as it stands:

1. If the top left element (first *pivot*) had been zero, no multiple of that row could have eliminated entries lower down in the first column. Similar breakdowns could have occurred at any later stage. Or if a pivot is very small, a very large multiple of that row would have to be used - this loses accuracy in floating point,
2. It is quite common that one needs to solve a system with many right hand sides (RHSs). If they are all known initially, we can just place them side-by-side when we eliminate as above. But it is also common that they become known only one after another. We then do not want to have to repeat the work on the coefficient part.
3. The coefficient matrix has often a lot of zero elements, e.g. it may be tridiagonal or banded, maybe with some additional elements present as well. Such structures can often be exploited to greatly reduce both operation count and storage.

We now address these issues in turn:

1. The solution of a linear system is totally independent of in which order the equations have been written down. The usually satisfactory method of *partial pivoting* simply amounts to, at stage k , exchanging row k with whatever row below it that has the largest element in the k^{th} position (thus never needing to use multiples in magnitude greater than one in the elimination steps - this makes growth of intermediate quantities less likely). The elimination process leaves the

determinant of the coefficient matrix unchanged - its value can be read off as $(-1)^p$ times the product of their pivots (the diagonal elements of the resulting lower triangular matrix); p denotes how many times we needed to swap order of the equations. Breakdown of partial pivoting will occur if and only if the system is singular.

2. Although not obvious from how we presented the solution of Example 1 above, what we achieved was in fact a factorization of the original matrix

$$A = \begin{bmatrix} 1 & -2 & 2 \\ -1 & -1 & 3 \\ 3 & 1 & -2 \end{bmatrix}$$

into a product $A = LU$ where L is lower triangular, and U upper triangular. We were left with the U -matrix

$$U = \begin{bmatrix} 1 & -2 & 2 \\ & 1 & -\frac{5}{3} \\ & & 1 \end{bmatrix},$$

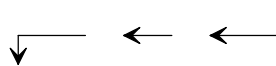
and just did not pay attention to the fact that

$$L = \begin{bmatrix} 1 & & \\ -1 & -3 & \\ 3 & 7 & \frac{11}{3} \end{bmatrix}$$

could have been written down by just taking note that these entries appeared at just these places during the elimination (Verify that the matrix product LU indeed recovers $A!$).

The reason this happened ($LU = A$) can be seen if we interpret every step we took as a matrix multiplication. Described in this way, the successive elimination steps were

$$\begin{bmatrix} 1 & & & \\ 1 & 1 & & \\ -3 & & 1 & \end{bmatrix} \begin{bmatrix} 1 & -2 & 2 & 1 \\ -1 & -1 & 3 & 1 \\ 3 & 1 & -2 & 2 \end{bmatrix} = \begin{bmatrix} 1 & -2 & 2 & 1 \\ -3 & 5 & 2 & \\ 7 & -8 & -1 & \end{bmatrix}$$



$$\begin{bmatrix} 1 & & & \\ -\frac{1}{3} & & & \\ & & 1 & \end{bmatrix} \begin{bmatrix} 1 & -2 & 2 & 1 \\ -3 & 5 & 2 & \\ 7 & -8 & -1 & \end{bmatrix} = \begin{bmatrix} 1 & -2 & 2 & 1 \\ 1 & -\frac{5}{3} & -\frac{2}{3} & \\ 7 & -8 & -1 & \end{bmatrix}$$



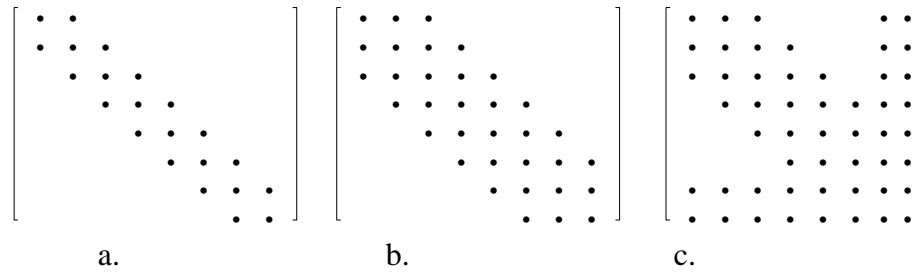


Figure 1. Three common matrix structures for linear systems.

Although much more complex sparsity patterns arise in many applications, strategies for them fall well outside the scope of this summary. For the three cases above:

- a. A matrix is said to be *diagonally dominant* if it holds that

$$|a_{i,i}| \geq |a_{i,i-1}| + |a_{i,i+1}| .$$

If this is not the case, we need to pivot. Tridiagonal then becomes a special case of type b. - general banded matrix. If it happens to be diagonally dominant, pivoting can be shown to be unnecessary, and the very fast *Thomas algorithm* can be used. This is equivalent to unpivoted Gaussian elimination, but is usually described differently: We aim for a standard $A = L U$ - factorization, which will be of the form

$$\begin{bmatrix} a_{11} & a_{12} & & & \\ a_{21} & a_{22} & a_{23} & & \\ & a_{32} & a_{33} & a_{34} & \\ & & \ddots & \ddots & \ddots \\ & & & \ddots & \ddots \end{bmatrix} = \begin{bmatrix} l_{11} & & & & \\ l_{21} & l_{22} & & & \\ & l_{32} & l_{33} & & \\ & & \ddots & \ddots & \\ & & & \ddots & \ddots \end{bmatrix} \begin{bmatrix} 1 & u_{12} & & & \\ & 1 & u_{23} & & \\ & & 1 & u_{34} & \\ & & & \ddots & 1 \\ & & & & \ddots \end{bmatrix}$$

For this $L U$ - product to become equal to A , we need

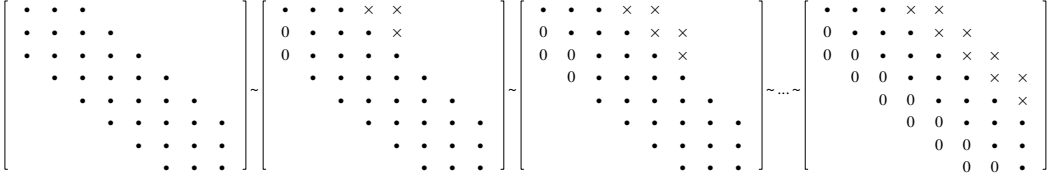
$$\begin{array}{llll} \text{row 1:} & & a_{11} = l_{11} & , & a_{12} = l_{11}u_{12} & , \\ \text{row 2:} & a_{21} = l_{21} & , & a_{22} = l_{21}u_{12} + l_{22} & , & a_{23} = l_{22}u_{23} & , \\ \text{row 3:} & a_{32} = l_{32} & , & a_{33} = l_{32}u_{23} + l_{33} & , & a_{34} = l_{33}u_{34} & , \\ \dots & \dots & & \dots & & \dots \\ \text{etc.} & & & & & \end{array}$$

We could simply have copied the subdiagonal of A over into the subdiagonal of L . The remaining equations give us recursively $l_{11}, u_{12}, l_{22}, u_{23}, l_{33}, u_{34}$ etc. The two back substitutions becomes equally quick recursions.

In a computer program, one saves a lot of storage by just storing the diagonal elements in a few vectors - there obviously is no need to store all the zero entries.

- b. In the banded case, we basically have to use the same strategy as for full matrices. The only difference is that L and U (even when partial pivoting is employed) will have non-zero entries only in some central diagonals. So again, it is much better to only store

only these diagonals (in the columns of a matrix), and confine all the arithmetic to within it. The following schematic illustrations show how the reduction and fill-in would proceed in a case where the band is equally wide on both sides of the main diagonal, and partial (row-) pivoting is used. Removed entries are marked as "0" and locations where new entries are introduced by "×":



The U -part will not only fit back into the matrix that used to hold the diagonals - it can directly overwrite in the space that the diagonals of A used to occupy. Had the RHS been available initially - and been manipulated along with the matrix - we would now be ready for back substitution. No additional storage would be needed.

If the RHS was not available at the time of the LU factorization, we also need to store one integer vector of pivot information - telling at each stage which row is moved into pivot position. Furthermore, we need to store the multiples used in the elimination. Using the same diagonal-based storage as above, the matrix structures would be

Input: Output:

$$A = \begin{bmatrix} \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \end{bmatrix}, \quad L = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}, \quad U = \begin{bmatrix} 1 & \cdot & \cdot & \cdot & \cdot \\ 1 & \cdot & \cdot & \cdot & \cdot \\ 1 & \cdot & \cdot & \cdot & \cdot \\ 1 & \cdot & \cdot & \cdot & \cdot \\ 1 & \cdot & \cdot & \cdot & \cdot \\ 1 & \cdot & \cdot & \cdot & \cdot \\ 1 & \cdot & \cdot & \cdot & \cdot \end{bmatrix}, \quad P = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

With the information in L , U and P , we can rapidly - and repeatedly - solve the linear system $A \underline{x} = \underline{b}$ whenever a RHS vector \underline{b} becomes available.

- c. If one happens to have available solvers for general full and banded linear systems, one can solve this case of banded with extra rows/columns quickly. Schematically, we write the system to solve as (assuming for simplicity we have the RHS available from the start)

$$\left[\begin{array}{c|c} A & B \\ \hline C & D \end{array} \right] \begin{bmatrix} \underline{x}_1 \\ \underline{x}_2 \end{bmatrix} = \begin{bmatrix} \underline{b}_1 \\ \underline{b}_2 \end{bmatrix} .$$

First, we solve (by a band solver) the many-RHS-system

$$\left[\begin{array}{c} A \end{array} \right] \left[\begin{array}{c|c} E & \underline{b}_3 \end{array} \right] = \left[\begin{array}{c|c} B & \underline{b}_1 \end{array} \right] .$$

The original system is now equivalent to

$$\left[\begin{array}{c|c} I & E \\ \hline C & D \end{array} \right] \left[\begin{array}{c} \underline{x}_1 \\ \underline{x}_2 \end{array} \right] = \left[\begin{array}{c} \underline{b}_3 \\ \underline{b}_2 \end{array} \right] ,$$

where I is the *identity matrix*, "1"s in the diagonal, and "0"s for the rest. This hold because we could have achieved the same thing (more costly) by multiplying the original system from the left by

$$\left[\begin{array}{c|c} A^{-1} & 0 \\ \hline 0 & I \end{array} \right] .$$

Now, we can add multiples of the top rows to clear out C , in the process turning \underline{b}_2 into \underline{b}_4 and D into F . As can be seen by multiplying the system with

$$\left[\begin{array}{c|c} I & 0 \\ \hline -C & I \end{array} \right] .$$

from the left, this step can be carried out by a simple matrix multiply-subtraction:

$$\left[\begin{array}{c|c} F & \underline{b}_4 \end{array} \right] = \left[\begin{array}{c|c} D & \underline{b}_2 \end{array} \right] - \left[\begin{array}{c} C \end{array} \right] \left[\begin{array}{c|c} E & \underline{b}_3 \end{array} \right]$$

The system is logically now in the form

$$\left[\begin{array}{c|c} I & E \\ \hline 0 & F \end{array} \right] \left[\begin{array}{c} \underline{x}_1 \\ \underline{x}_2 \end{array} \right] = \left[\begin{array}{c} \underline{b}_3 \\ \underline{b}_4 \end{array} \right] .$$

Next we solve the square system in the bottom right corner

$$\begin{bmatrix} F \end{bmatrix} \begin{bmatrix} \underline{x}_2 \end{bmatrix} = \begin{bmatrix} \underline{b}_4 \end{bmatrix} ,$$

which has become independent of the other equations. Finally, knowing \underline{x}_2 , we find the remaining unknowns \underline{x}_1 by means of

$$\begin{bmatrix} \underline{x}_1 \end{bmatrix} = \begin{bmatrix} \underline{b}_3 \end{bmatrix} - \begin{bmatrix} E \end{bmatrix} \begin{bmatrix} \underline{x}_2 \end{bmatrix} .$$