

Chap 7. Computation with Matrices

7.1 Introduction

In numerical analysis, there is a survival of the fittest, and we describe some ideas that have survived.

① Techniques for Solving $A\mathbf{x} = \mathbf{b}$.

Elimination is a perfect algorithm, except when the particular problem has special properties – as almost every problem has. For sparse matrices, we develop iterative methods for solving $A\mathbf{x} = \mathbf{b}$.

The object is to get close more quickly than elimination. In some problems, that can be done; in many others, elimination is safer and faster if it takes advantage of the zeros. The competition is far from over.

② Techniques for solving $A\mathbf{x} = \lambda\mathbf{x}$.

Until recently no one knew how to solve the eigenvalue problem. Dozens of algorithms have been suggested.

Two or three ideas have superseded others: the QR algorithms, the power methods, and the preprocessing of a symmetric matrix to make it tridiagonal. The last is a direct method, but it produces a much simpler matrix to use in the iterative steps.

③ The Condition Number of a Matrix.

The condition number measures the sensitivity of A by multiplying the norms of A and A^{-1} .

The matrices in this chapter are square.

7.3 Computation of Eigenvalues

The ordinary power method starts with u_0 and each step $u_{k+1} = A u_k$ is a matrix-vector multiplication. $\Rightarrow u_k = A^k u_0$. The multiplication by A should be easy, i.e., a large matrix A should be sparse. Assuming A has a full set of eigenvectors x_1, \dots, x_n , the vector u_k will be

$$u_k = c_1 \lambda_1^{k_1} x_1 + \dots + c_n \lambda_n^{k_n} x_n.$$

Suppose $|\lambda_1| \leq |\lambda_2| \leq \dots \leq |\lambda_{n-1}| < |\lambda_n|$. When u_0 contained some component of x_n , i.e., $c_n \neq 0$, the component $c_n x_n$ will dominate in u_k :

$$\frac{u_k}{(\lambda_n)^{k_n}} = c_1 \left(\frac{\lambda_1}{\lambda_n} \right)^{k_n} x_1 + \dots + c_{n-1} \left(\frac{\lambda_{n-1}}{\lambda_n} \right)^{k_n} x_{n-1} + c_n x_n$$

The convergence factor is the ratio $r = |\lambda_{n-1}|/\lambda_n$. Often we divide each u_k by its first component a_k before taking the next step: $u_{k+1} = A u_k / a_k$, where a_k will approach λ_n .

Ex 1:

$$u_k \rightarrow \begin{bmatrix} 2/3 \\ 1/3 \end{bmatrix} = \begin{bmatrix} 0.667 \\ 0.333 \end{bmatrix} \text{ when } A = \begin{bmatrix} 0.9 & 0.2 \\ 0.1 & 0.8 \end{bmatrix}$$
$$u_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, u_1 = \begin{bmatrix} 0.9 \\ 0.1 \end{bmatrix}, u_2 = \begin{bmatrix} 0.83 \\ 0.17 \end{bmatrix}, u_3 = \begin{bmatrix} 0.781 \\ 0.219 \end{bmatrix}, u_4 = \begin{bmatrix} 0.740 \\ 0.253 \end{bmatrix}$$

If $r = |\lambda_{n-1}|/\lambda_n$ is close to 1, then convergence is slow.

If $r=1$ (i.e., $|\lambda_{n-1}|=\lambda_n$), then convergence will not occur.

There are several ways to get around this limitation.

① The block power method works with several vectors at once.

If we multiply p orthonormal vectors by A , and then apply Gram-Schmidt to orthogonalize them again, the convergence ratio becomes $r' = |\lambda_{n-p}|/\lambda_n$. We obtain approximations to p different eigenvalues and their eigenvectors.

② The inverse power method operates with A^T instead of A .

A single step is $\mathbf{v}_{k+1} = A^T \mathbf{v}_k$ by solving $A\mathbf{v}_{k+1} = \mathbf{v}_k$.
(Save the factors L and U !) Now we converge to
the smallest eigenvalue λ_1 and its eigenvector \mathbf{x} if $|\lambda_1| < |\lambda_2|$.
Often it is λ_1 that is wanted in the applications.

③ The shifted inverse power method is best of all.

Replace A by $A - \alpha I$. Each eigenvalue is shifted by α , and
the convergence factor for the inverse method will change to
 $r'' = |\lambda_1 - \alpha| / |\lambda_2 - \alpha|$. If α is a good approximation to λ_1 ,
 r'' will be very small and the convergence is enormously
accelerated. Each solves $(A - \alpha I) \mathbf{w}_{k+1} = \mathbf{w}_k$

$$\mathbf{w}_k = \frac{c_1 \mathbf{x}_1}{(\lambda_1 - \alpha)^k} + \frac{c_2 \mathbf{x}_2}{(\lambda_2 - \alpha)^k} + \dots + \frac{c_n \mathbf{x}_n}{(\lambda_n - \alpha)^k}$$

When α is close to λ_1 , the first term dominates. If a good
approximation to any eigenvalue has been computed by another algorithm
(such as QR), this is a good way to find the eigenvector.

If λ_1 is not already approximated, the shifted inverse power
method has to generate its own choice of α . We can vary $\alpha = \alpha_k$
at every step : $(A - \alpha_k I) \mathbf{w}_{k+1} = \mathbf{w}_k$. When A is symmetric,
a very accurate choice is the Rayleigh quotient :

$$\text{shift by } \alpha_k = R(\mathbf{w}_k) = \frac{\mathbf{w}_k^T A \mathbf{w}_k}{\mathbf{w}_k^T \mathbf{w}_k}$$

This quotient $R(\mathbf{x})$ has a minimum at the true eigenvector \mathbf{x}_1 .

The error $\lambda_1 - \alpha_k$ is roughly the square of the error in the eigenvector.
The convergence factors $|\lambda_1 - \alpha_k| / |\lambda_2 - \alpha_k|$ converge to zero. Then
these Rayleigh quotient shifts give cubic convergence of α_k to λ_1 .

Tridiagonal and Hessenberg Forms

o. The power method is reasonable only for a large sparse matrix. Our goal is to find a simple way to create zeros. After computing a similar matrix $Q^T A Q$ with more zeros than A , we don't go back to the power method. The QR algorithm is a much more powerful variant. (The shifted inverse power method has its place at the very end, in finding the eigenvector.) The first step is to produce quickly many zeros, using Q .

If A is symmetric, then so is $Q^T A Q$. Q is an orthogonal matrix.

o. The Hessenberg form accepts one nonzero diagonal below the main diagonal. If a Hessenberg matrix is symmetric, it only has three nonzero diagonals. A Householder matrix is a reflection matrix determined by one vector v :

$$\text{Householder matrix: } H = I - 2 \frac{v \cdot v^T}{\|v\|^2} = I - 2 u u^T,$$

where $u = v/\|v\|$. H is both symmetric and orthogonal.

$$H^T H = (I - 2 u u^T)(I - 2 u u^T) = I - 4 u u^T + 4 u u^T u u^T = I$$

Thus $H = H^T = H^{-1}$.

o. Householder's plan was to produce zeros with these matrices.

Its success depends on the following identity $Hx = -\alpha z$:

Suppose $z = (1, 0, \dots, 0)^T$, $\alpha = \|x\|$, and $v = x + \alpha z$.

Then $Hx = -\alpha z = (-\alpha, 0, \dots, 0)^T$.

$$\begin{aligned} \text{Pf: } Hx &= x - \frac{2 v \cdot v^T}{\|v\|^2} x = x - (x + \alpha z) \frac{2(x + \alpha z)^T x}{(x + \alpha z)^T (x + \alpha z)} \\ &= x - (x + \alpha z) \quad (\because x^T x = \alpha^2) \\ &= -\alpha z \end{aligned}$$

This identity can be used right away, on the first column of A . The final $Q^{-1}AQ$ is allowed one nonzero diagonal below the main diagonal (Hessenberg form). Therefore *only the entries strictly below the diagonal will be involved*:

$$x = \begin{bmatrix} a_{21} \\ a_{31} \\ \vdots \\ a_{n1} \end{bmatrix}, \quad z = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad Hx = \begin{bmatrix} -\sigma \\ 0 \\ \vdots \\ 0 \end{bmatrix}. \quad (3)$$

At this point Householder's matrix H is only of order $n - 1$, so it is embedded into the lower right-hand corner of a full-size matrix U_1 :

$$U_1 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & & & & \\ 0 & H & & & \\ 0 & & & & \\ 0 & & & & \end{bmatrix} = U_1^{-1}, \quad \text{and} \quad U_1^{-1}AU_1 = \begin{bmatrix} a_{11} & * & * & * & * \\ -\sigma & * & * & * & * \\ \mathbf{0} & * & * & * & * \\ \mathbf{0} & * & * & * & * \\ \mathbf{0} & * & * & * & * \end{bmatrix}.$$

The first stage is complete, and $U_1^{-1}AU_1$ has the required first column. At the second stage, x consists of the last $n - 2$ entries in the second column (three bold stars). Then H_2 is of order $n - 2$. When it is embedded in U_2 , it produces

$$U_2 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & & & \\ 0 & 0 & H_2 & & \\ 0 & 0 & & & \end{bmatrix} = U_2^{-1}, \quad U_2^{-1}(U_1^{-1}AU_1)U_2 = \begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ \mathbf{0} & * & * & * & * \\ \mathbf{0} & \mathbf{0} & * & * & * \\ \mathbf{0} & \mathbf{0} & * & * & * \end{bmatrix}.$$

U_3 will take care of the third column. For a 5 by 5 matrix, the Hessenberg form is achieved (it has six zeros). In general Q is the product of all the matrices $U_1 U_2 \cdots U_{n-2}$, and the number of operations required to compute it is of order n^3 .

Example 2

(to change $a_{13} = a_{31}$ to zero)

$$A = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 0 \end{bmatrix}, \quad x = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad v = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad H = \begin{bmatrix} 0 & -1 \\ -1 & 0 \end{bmatrix}.$$

Embedding H into Q , the result $Q^{-1}AQ$ is tridiagonal:

$$Q = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & -1 & 0 \end{bmatrix}, \quad Q^{-1}AQ = \begin{bmatrix} 1 & -1 & \mathbf{0} \\ -1 & 0 & 1 \\ \mathbf{0} & 1 & 1 \end{bmatrix}.$$

Two other applications of the Householder matrix H

① The Gram-Schmidt factorization $A = QR$.

The first step is to work with the first column of A:

$$x = \begin{bmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{n1} \end{bmatrix}, z = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow H_1 = I - 2 \frac{v v^T}{\|v\|^2}$$

The first column of $H_1 A$ equals $- \|x\| z$, which is also the first column of R. The second step works with the second column of $H_1 A$, from the pivot on down, and produces an $H_2 H_1 A$ which is zero below that pivot. The result of $(n-1)$ steps is an upper triangular R.

$$H_m H_{m-1} \cdots H_1 A = R \Rightarrow A = (H_1 \cdots H_{m-1}) R = QR$$

The product $Q = H_1 \cdots H_{m-1}$ can be stored in the factored form (keep only the v's) and never computed explicitly.

② The Singular Value Decomposition $U^T A V = \Sigma$.

The first step is as above. The next step is to multiply on the right by an $H^{(1)}$, which produces zeros as indicated along the first row:

$$A \rightarrow H_1 A = \begin{bmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \end{bmatrix} \rightarrow H_1 A H^{(1)} = \begin{bmatrix} * & * & 0 & 0 \\ 0 & * & * & * \\ 0 & * & * & * \end{bmatrix}.$$

Repeating similar steps produces the **bidagonal form**

$$H_2 H_1 A H^{(1)} = \begin{bmatrix} * & * & 0 & 0 \\ 0 & * & * & * \\ 0 & 0 & * & * \end{bmatrix} \quad \text{and}$$

$$H_2 H_1 A H^{(1)} H^{(2)} = \begin{bmatrix} * & * & 0 & 0 \\ 0 & * & * & 0 \\ 0 & 0 & * & * \end{bmatrix}.$$

The QR Algorithm for Computing Eigenvalues

It starts with $A_0 = Q_0 R_0$ and reverses the factors: $A_0 = R_0 Q_0$.
 A_1 is similar to A_0 ($\because Q_0^{-1} A_0 Q_0 = Q_0^{-1} (Q_0 R_0) Q_0 = A_1$)

The process continues with no change in the eigenvalues.

All A_K are similar: $A_K = Q_K R_K$ then $A_{K+1} = R_K Q_K$

Almost always A_K approaches a triangular form. Its diagonal entries approach its eigenvalues. If there was some processing to a tridiagonal form A_0 , then A_0 is connected to the original A by $Q^T A Q = A_0$.

To speed up the QR algorithm, we must allow shifts to $A_K - \alpha_K I$, and the QR factorization should be very quick.

① The Shifted Algorithm. If α_K is close to an eigenvalue,

$$A_K - \alpha_K I = Q_K R_K \text{ and then } A_{K+1} = R_K Q_K + \alpha_K I$$

A_{K+1} is similar to A_K (always the same eigenvalues):

$$Q_K^{-1} A_K Q_K = Q_K^{-1} (Q_K R_K + \alpha_K I) Q_K = R_K Q_K + \alpha_K I = A_{K+1}.$$

The (m,n) entry of A_K is the first to approach an eigenvalue.

The entry is chosen for α_K . Normally it produces quadratic convergence, and in the symmetric case even cubic convergence to the smallest eigenvalue. After three or four steps,

$$A_K = \left[\begin{array}{ccc|c} * & * & * & * \\ * & * & * & * \\ 0 & * & * & * \\ \hline 0 & 0 & \epsilon & \lambda'_1 \end{array} \right], \quad \text{with } \epsilon \ll 1.$$

We accept $\lambda'_1 \approx \lambda_1$. To find λ_2 , we continue with the upper-left $(n-1) \times (n-1)$ submatrix. Another two steps are sufficient to find λ_2 To find the eigenvectors, a single inverse power step is usually good enough.

2. When A_0 is tridiagonal or Hessenberg, each QR step is very fast. The Gram–Schmidt process (factoring into QR) takes $O(n^3)$ operations for a full matrix A . For a Hessenberg matrix this becomes $O(n^2)$, and for a tridiagonal matrix it is $O(n)$. Fortunately, each new A_k is again in Hessenberg or tridiagonal form:

$$Q_0 \text{ is Hessenberg} \quad Q_0 = A_0 R_0^{-1} = \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ 0 & * & * & * \\ 0 & 0 & * & * \end{bmatrix} \begin{bmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & 0 & * & * \\ 0 & 0 & 0 & * \end{bmatrix}.$$

You can easily check that this multiplication leaves Q_0 with the same three zeros as A_0 . *Hessenberg times triangular is Hessenberg*. So is triangular times Hessenberg:

$$A_1 \text{ is Hessenberg} \quad A_1 = R_0 Q_0 = \begin{bmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & 0 & * & * \\ 0 & 0 & 0 & * \end{bmatrix} \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ 0 & * & * & * \\ 0 & 0 & * & * \end{bmatrix}.$$

The symmetric case is even better, since $A_1 = Q_0^{-1} A_0 Q_0 = Q_0^T A_0 Q_0$ stays symmetric. By the reasoning just completed, A_1 is also Hessenberg. So A_1 must be *tridiagonal*. The same applies to A_2, A_3, \dots , and *every QR step begins with a tridiagonal matrix*.

The last point is the factorization itself, producing the Q_k and R_k from each A_k (or really from $A_k - \alpha_k I$). We may use Householder again, but it is simpler to annihilate each subdiagonal element in turn by a “plane rotation” P_{ij} . The first is P_{21} :

$$\text{Rotation to kill } a_{21} \quad P_{21} A_k = \begin{bmatrix} \cos \theta & -\sin \theta & & \\ \sin \theta & \cos \theta & & \\ & & 1 & \\ & & & 1 \end{bmatrix} \begin{bmatrix} a_{11} & * & * & * \\ a_{21} & * & * & * \\ 0 & * & * & * \\ 0 & 0 & * & * \end{bmatrix} \quad (7)$$

The $(2, 1)$ entry in this product is $a_{11} \sin \theta + a_{21} \cos \theta$, and we choose the angle θ that makes this combination zero. The next rotation P_{32} is chosen in a similar way, to remove the $(3, 2)$ entry of $P_{32} P_{21} A_k$. After $n - 1$ rotations, we have R_0 :

$$\text{Triangular factor} \quad R_k = P_{nn-1} \cdots P_{32} P_{21} A_k. \quad (8)$$

7.4 Iterative Methods for $A\mathbf{x} = \mathbf{b}$

An iterative method is easy to invent, by splitting the matrix A . If $A = S - T$, then $A\mathbf{x} = \mathbf{b}$ is the same as $S\mathbf{x} = T\mathbf{x} + \mathbf{b}$.

$$\text{Iteration: } S\mathbf{x}_{k+1} = T\mathbf{x}_k + \mathbf{b}$$

- ① The new vector \mathbf{x}_{k+1} should be easy to compute. Thus S should be a simple (and invertible) matrix - diagonal or triangular.
- ② The sequence \mathbf{x}_k should converge to the true solution.

$$\text{Error equation: } S\mathbf{e}_{k+1} = T\mathbf{e}_k, \text{ where } \mathbf{e}_k = \mathbf{x}_k - \mathbf{x}.$$

It starts with the initial error \mathbf{e}_0 , and $\mathbf{e}_k = (S^{-1}T)^k \mathbf{e}_0$.

The convergence $\mathbf{x}_k \rightarrow \mathbf{x}$ exactly when $\mathbf{e}_k \rightarrow \mathbf{0}$.

$\boxed{\text{TF}}$

The iterative method is convergent if and only if every eigenvalue of $S^{-1}T$ satisfies $|\lambda| < 1$. Its rate of convergence depends on the maximum size of $|\lambda|$:

$$\text{Spectral radius "rho": } \rho(S^{-1}T) = \max_{\lambda \in \mathbb{C}} |\lambda|$$

A typical solution to $\mathbf{e}_{k+1} = S^{-1}T\mathbf{e}_k$ is a combination of eigenvectors:

$$\mathbf{e}_k = c_1 \lambda_1^{k_1} \mathbf{x}_1 + \dots + c_n \lambda_n^{k_n} \mathbf{x}_n.$$

The largest $|\lambda_i|$ will eventually be dominant, so the spectral radius $\rho = |\lambda_{\max}|$ will govern the rate at which \mathbf{e}_k converges to zero. We certainly need $\rho < 1$. We start with three possibilities

- ① $S = \text{diagonal part of } A$ (Jacobi's method)
- ② $S = \text{triangular part of } A$ (Gauss-Seidel method)
- ③ $S = \text{combination of ①&②}$ (Successive overrelaxation) or SOR

S is also called a preconditioner, and its choice is crucial in numerical analysis.

Ex 1: (Jacobi) Here S is the diagonal part of A

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}, S = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, T = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, S^+T = \begin{bmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{bmatrix}$$

For $\mathbf{x} = \begin{bmatrix} v \\ w \end{bmatrix}$, the Jacobi step is $S\mathbf{x}_{k+1} = T\mathbf{x}_k + \mathbf{b}$

$$\begin{aligned} 2v_{k+1} &= w_k + b_1 \quad \text{or} \quad \begin{bmatrix} v \\ w \end{bmatrix}_{k+1} = \begin{bmatrix} 0 & \frac{1}{2} \\ \frac{1}{2} & 0 \end{bmatrix} \begin{bmatrix} v \\ w \end{bmatrix}_k + \begin{bmatrix} b_1/2 \\ b_2/2 \end{bmatrix} \\ 2w_{k+1} &= v_k + b_2 \end{aligned}$$

The decisive matrix S^+T has eigenvalues $\pm \frac{1}{2}$, which means that the error is cut in half at every step.

The Gauss-Seidel method starts using each component of the new \mathbf{x}_{k+1} as soon as it is computed and thus requires only half as much space as the Jacobi.

$$\left\{ \begin{array}{l} \text{New } x_1 : a_{11}(x_1)_{k+1} = (-a_{12}x_2 - a_{13}x_3 - \dots - a_{1n}x_n)_k + b_1 \\ \text{New } x_2 : a_{22}(x_2)_{k+1} = -a_{21}(x_1)_{k+1} + (-a_{23}x_3 - \dots - a_{2n}x_n)_k + b_2 \\ \dots \\ \text{New } x_n : a_{nn}(x_n)_{k+1} = (-a_{n1}x_1 - a_{n2}x_2 - \dots - a_{n,n-1}x_{n-1})_{k+1} + b_n \end{array} \right.$$

Ex 2: (Gauss-Seidel) Here S^+T has smaller eigenvalues:

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}, S = \begin{bmatrix} 2 & 0 \\ -1 & 2 \end{bmatrix}, T = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}, S^+T = \begin{bmatrix} 0 & \frac{1}{4} \\ 0 & \frac{1}{4} \end{bmatrix}$$

$$\left\{ \begin{array}{l} 2v_{k+1} = w_k + b_1 \quad \text{or} \quad \begin{bmatrix} 2 & 0 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} v \\ w \end{bmatrix}_{k+1} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} v \\ w \end{bmatrix}_k + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \\ 2w_{k+1} = v_k + b_2 \end{array} \right.$$

The eigenvalues of S^+T are $\frac{1}{4}$ and 0 . The error is divided by 4 every time, so a single Gauss-Seidel step is worth two Jacobi steps. Both methods require the same number of operations. The Gauss-Seidel is better and saves storage.

There are examples in which Jacobi converges and Gauss-Seidel fails (or conversely). The symmetric case is straightforward: When all $a_{ii} > 0$, Gauss-Seidel converges $\Leftrightarrow A$ positive definite.

Convergence is faster if we go beyond the Gauss-Seidel correction $\hat{x}_{k+1} - \hat{x}_k$. Roughly speaking, these approximations stay on the same side of the solution x . An overrelaxation factor w moves us closer to the solution. When $w=1$, we recover Gauss-Seidel. With $w>1$, the method is known as SOR, Successive Overrelaxation. The optimal choice of w never exceeds 2. It is often ≈ 1.9 .

Let D, L, U be the parts of A on, below, and above the diagonal. (They are different from the $A=LDU$.) In fact, we now have $A=L+D+U$. The Jacobi method has $S=D$ and $T=-L-U$, and the Gauss-Seidel chose $S=D+L$ and $T=-U$. To accelerate the convergence,

$$\text{Overrelaxation: } [D+wL] \hat{x}_{k+1} = [(\mathbf{I}-w\mathbf{U})] \hat{x}_k + w\mathbf{b}$$

lower triangular upper triangular

\hat{x}_{k+1} can still replace \hat{x}_k , component by component, as soon as it is computed. A typical step is

$$a_{kk}(x_k)_{k+1} = a_{kk}(x_k)_k$$

$$+ w [(-a_{11}x_1 - \dots - a_{1,k-1}x_{k-1})_{k+1} + (-a_{22}x_2 - \dots - a_{2,k-1}x_{k-1})_{k+1} + \dots + (-a_{kk}x_k)_{k+1} + b_k]$$

Ex3: (SOR) For the same matrix A , each SOR step is

$$\begin{bmatrix} 2 & 0 \\ -w & 2 \end{bmatrix} \begin{bmatrix} w \\ w \end{bmatrix}_{k+1} = \begin{bmatrix} 2(1-w) & w \\ 0 & 2(1-w) \end{bmatrix} \begin{bmatrix} w \\ w \end{bmatrix}_k + w \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

The crucial matrix $L=S^{-1}T$ is

$$L = \begin{bmatrix} 2 & 0 \\ -w & 2 \end{bmatrix}^{-1} \begin{bmatrix} 2(1-w) & w \\ 0 & 2(1-w) \end{bmatrix} = \begin{bmatrix} 1-w & \frac{1}{2}w \\ \frac{1}{2}w(1-w) & 1-w+\frac{1}{4}w^2 \end{bmatrix}$$

The optimal w makes the largest eigenvalue of L (its spectral radius) as small as possible. The product of the eigenvalues of L equals $\det L = \det T / \det S$:

$$\lambda_1 \lambda_2 = \det L = (1-w)^2$$

Always $\det S = \det D$ and $\det T = \det(I - wD)$. Their product is $\det L = (I - wD)^n$. This explains why $w < 2$. At the optimal w , the two eigenvalues are equal to $w-1$. Their sum agrees to the sum of diagonal entries:

$$\text{Optimal } w: \lambda_1 + \lambda_2 = 2(w_{\text{opt}} - 1) = 2 - 2w_{\text{opt}} + \frac{1}{4}w_{\text{opt}}^2$$

$$\therefore w_{\text{opt}} = 4(2 - \sqrt{3}) \approx 1.07 \text{ and } w_{\text{opt}} - 1 \approx 0.07$$

This is a major reduction from the Gauss-Seidel: $\lambda = \frac{1}{4}$ at $w=1$. In this example, the choice of $w_{\text{opt}} = 1.07$ has again doubled the rate of convergence, because $(\frac{1}{4})^2 \approx 0.01$.

Young's 1950 thesis solved a simple formula for the optimal w . The key step is to connect the eigenvalue λ of L to the eigenvalue μ of the original Jacobi matrix $D^T(-L-U)$:

$$\text{Formula for } w: (\lambda + w - 1)^2 = \lambda w^2 \mu^2 \quad (7)$$

When $w=1$ (Gauss-Seidel), it yields $\lambda^2 = \lambda \mu^2 \Rightarrow \lambda=0$ and $\lambda=\mu^2$. (In Ex 2, $\mu = \pm \frac{1}{2}$ and $\lambda=0, \lambda=\frac{1}{4}$.) The important problem is to choose w so that λ_{\max} will be minimized. Young's Eq (7) is exactly our 2×2 example! The best w makes the two roots λ both equal to $w-1$

$$(w-1) + (w-1) = 2 - 2w + \mu^2 w^2 \text{ or } w = \frac{2(1 - \sqrt{1 - \mu^2})}{\mu^2}$$

For a large matrix, this pattern will be repeated for a number of different pairs $\pm \mu_i$. The largest μ gives the largest value of w and of $\lambda = w-1$. Since our goal is to make λ_{\max} as small as possible, that extremal pair specifies

$$\text{Optimal } w: w_{\text{opt}} = \frac{2(1 - \sqrt{1 - \mu_{\max}^2})}{\mu_{\max}^2} \text{ and } \lambda_{\max} = w_{\text{opt}} - 1.$$