

Computing the solution to least-squares problems

Here are the least-squares problems we have talked about so far:

Pseudo-inverse when \mathbf{A} has full column rank $\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$

Pseudo-inverse when \mathbf{A} has full row rank $\hat{\mathbf{x}} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1} \mathbf{y}$

Ridge Regression / Tikhonov regularization $\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A} + \delta \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}$

Kernel Regression $\hat{\boldsymbol{\alpha}} = (\mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{y}$

Here are two that we did not talk about, but are easy extensions:

Constrained least-squares:

minimize $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ subject to $\mathbf{x} = \mathbf{G}\boldsymbol{\alpha}$ for some $\boldsymbol{\alpha}$,

has solution

$$\hat{\mathbf{x}} = \mathbf{G}(\mathbf{G}^T \mathbf{A}^T \mathbf{A} \mathbf{G})^{-1} \mathbf{G}^T \mathbf{A}^T \mathbf{y}.$$

Generalized Tikhonov regularization:

minimize $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \delta \|\mathbf{D}\mathbf{x}\|_2^2$,

has solution

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A} + \delta \mathbf{D}^T \mathbf{D})^{-1} \mathbf{A}^T \mathbf{y}.$$

Each of the problems on the previous page involves solving a system of **symmetric positive definite** system of equations.

There are many ways to solve general systems of equations. Most general methods revolve around factoring the matrix \mathbf{A} into a series of systems that are much easier to solve — doing this allows us to re-use our work if we have multiple right-hand sides.

We will start with a (very) brief overview of how to solve general systems of equations using explicit matrix computations. If you want to read more about this, the classic reference is:

G. H. Golub and C. F. van Loan, *Matrix Computations*, now in its 4th edition (2012).

First, to set some context, let's look at some particular types of systems which are “easy” to solve. In all of the examples below, \mathbf{A} is an invertible $N \times N$ matrix.

Diagonal systems. If

$$\mathbf{A} = \begin{bmatrix} a_{11} & 0 & \cdots & \\ 0 & a_{22} & & \\ & & \ddots & \\ & & & a_{NN} \end{bmatrix}$$

with $a_{nn} \neq 0$, then solving $\mathbf{A}\mathbf{x} = \mathbf{b}$ for a given \mathbf{b} is easy; simply take

$$\hat{x}[n] = b[n]/a_{nn}.$$

This is of course very efficient computationally; we can compute $\hat{\mathbf{x}}$ in $O(N)$ time.

Example. Solve

$$\begin{bmatrix} 3 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 4 \end{bmatrix} \begin{bmatrix} x[1] \\ x[2] \\ x[3] \end{bmatrix} = \begin{bmatrix} 0.5 \\ 14 \\ 7 \end{bmatrix}$$

Orthogonal systems. If the columns (or equivalently the rows) of \mathbf{A} are orthonormal, $\mathbf{A}^T \mathbf{A} = \mathbf{I}$, then $\mathbf{A}^{-1} = \mathbf{A}^T$, and we can solve $\mathbf{A}\mathbf{x} = \mathbf{b}$ with a single matrix-vector multiply:

$$\hat{\mathbf{x}} = \mathbf{A}^T \mathbf{b}.$$

In general, the cost of this matrix-vector multiply is $O(N^2)$.

Example: Solve

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1/\sqrt{2} & 1/\sqrt{2} \\ 0 & -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \begin{bmatrix} x[1] \\ x[2] \\ x[3] \end{bmatrix} = \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix}$$

Triangular systems. If \mathbf{A} is lower triangular in that all of the terms above its main diagonal are zero,

$$\mathbf{A} = \begin{bmatrix} a_{11} & 0 & 0 & & \\ a_{21} & a_{22} & 0 & & \\ a_{31} & a_{32} & a_{33} & & \\ \vdots & & & \ddots & \\ a_{N1} & a_{N2} & \cdots & & a_{NN} \end{bmatrix},$$

with $a_{nn} \neq 0$, then given \mathbf{b} , we can solve $\mathbf{Ax} = \mathbf{b}$ using *forward substitution*:

$$\begin{aligned} \hat{x}[1] &= b[1]/a_{11} \\ \hat{x}[2] &= (b[2] - a_{21}\hat{x}[1])/a_{22} \\ &\vdots \\ \hat{x}[N] &= \left(b[N] - \sum_{n=1}^{N-1} a_{Nn}\hat{x}[n] \right) / a_{NN}. \end{aligned}$$

The total cost of this is about the same as a vector-matrix multiply, $O(N^2)$. If \mathbf{A} is upper triangular,

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1N} \\ 0 & a_{22} & a_{23} & & a_{2N} \\ 0 & 0 & a_{33} & & a_{3N} \\ \vdots & & & \ddots & \\ 0 & & & & a_{NN} \end{bmatrix},$$

then given \mathbf{y} we can solve $\mathbf{Ax} = \mathbf{y}$ using *backward substitution*. (Write it down at home!)

Example: Solve

$$\begin{bmatrix} 2 & 0 & 0 \\ -1 & 3 & 0 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} x[1] \\ x[2] \\ x[3] \end{bmatrix} = \begin{bmatrix} 14 \\ -2 \\ 1 \end{bmatrix}$$

Matrix factorization

The general strategy for solving a system of equations is to **factor** the matrix \mathbf{A} into multiple components, each of which has one of the structures above. These factorizations are not significantly more expensive than solving a system using Gaussian elimination, and once they are performed, solving another system $\mathbf{A}\mathbf{x} = \mathbf{b}$ with the same \mathbf{A} (but different \mathbf{b}) is fast.

LU factorization

Every $N \times N$ matrix \mathbf{A} can be written as a product $\mathbf{A} = \mathbf{L}\mathbf{U}$, where \mathbf{L} is lower diagonal, and \mathbf{U} is upper diagonal. The decomposition is in general not unique, but if we restrict the diagonal entries of \mathbf{L} (or \mathbf{U}) to be 1 (or anything $\neq 0$), then it becomes unique (when \mathbf{A} is invertible).

Example:

$$\begin{bmatrix} 2 & 1 & -1 \\ -3 & -1 & 2 \\ -2 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ -1.5 & 1 & 0 \\ -1 & 4 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & -1 \\ 0 & 0.5 & -0.5 \\ 0 & 0 & -1 \end{bmatrix}$$
$$\mathbf{A} = \mathbf{L} \mathbf{U}$$

If \mathbf{A} is invertible, then both \mathbf{L} and \mathbf{U} are invertible. So given $\mathbf{A}\mathbf{x} = \mathbf{b}$, we can solve for \mathbf{x} as follows:

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b} = \mathbf{U}^{-1}\mathbf{L}^{-1}\mathbf{b}.$$

More explicitly, we solve for \mathbf{w} in $\mathbf{L}\mathbf{w} = \mathbf{b}$, then solve for \mathbf{x} in $\mathbf{U}\mathbf{x} = \mathbf{w}$. As we have argued above, the cost of solving each of

these systems of equations with the factorization in place is $O(N^2)$ (as opposed to $O(N^3)$).

Computing the **LU** factorization is basically the same as recording all of your work while you are doing Gaussian elimination. We start with the original matrix, perform row operations on it until it is in upper triangular form. Each of the row operations corresponds to a lower-diagonal matrix with one off diagonal term, and their product will also be lower diagonal.

Example: Start with the original matrix

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

Eliminate the lower-left term. In this case we can add the first row to the third row, i.e.

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & -1 \\ -3 & -1 & 2 \\ -2 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 2 & 1 & -1 \\ -3 & -1 & 2 \\ 0 & 2 & 1 \end{bmatrix}$$

$$\mathbf{L}_1 \quad \mathbf{A} \quad = \quad \mathbf{A}_1$$

Eliminate the term in (row,column) = (2, 1) by adding 3/2 the top row to the second row:

$$\begin{bmatrix} 1 & 0 & 0 \\ 1.5 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & -1 \\ -3 & -1 & 2 \\ 0 & 2 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 1 & -1 \\ 0 & 0.5 & 0.5 \\ 0 & 2 & 1 \end{bmatrix}$$

$$\mathbf{L}_2 \quad \mathbf{A}_1 \quad = \quad \mathbf{A}_2$$

Finally, we eliminate (3, 2) by subtracting 4 times the second row

from the third row:

$$\begin{array}{ccc} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -4 & 1 \end{bmatrix} & \begin{bmatrix} 2 & 1 & -1 \\ 0 & 0.5 & 0.5 \\ 0 & 2 & 1 \end{bmatrix} & = \begin{bmatrix} 2 & 1 & -1 \\ 0 & 0.5 & 0.5 \\ 0 & 0 & -1 \end{bmatrix} \\ \mathbf{L}_3 & \mathbf{A}_2 & = \mathbf{U}. \end{array}$$

So we have:

$$\mathbf{L}_3 \mathbf{L}_2 \mathbf{L}_1 \mathbf{A} = \mathbf{U},$$

where \mathbf{U} is upper triangular and the \mathbf{L}_i are all lower triangular with 1 along the diagonal and exactly 1 non-zero off-diagonal term. Using the facts that

1. the inverses of \mathbf{L}_i are also lower triangular with 1 down the diagonal and exactly one non-zero off-diagonal term (show this at home!), and
2. the product of two lower triangular matrices with 1 down the diagonal is again lower-triangular with 1 down the diagonal (show this at home!),

we have

$$\mathbf{A} = (\mathbf{L}_1^{-1} \mathbf{L}_2^{-1} \mathbf{L}_3^{-1}) \mathbf{U} = \mathbf{L} \mathbf{U},$$

where \mathbf{L} is lower diagonal with 1 down the diagonal. In the example above, we have

$$\begin{aligned} \mathbf{L} = \mathbf{L}_1^{-1} \mathbf{L}_2^{-1} \mathbf{L}_3^{-1} &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ -1.5 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 4 & 1 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & 0 \\ -1.5 & 1 & 0 \\ -1 & 4 & 1 \end{bmatrix} \end{aligned}$$

If we count up the operations above, we have to eliminate $(N-1)N/2$ terms, each at a cost of $O(N)$, so computing the \mathbf{LU} factorization is $O(N^3)$ (same as using Gaussian elimination to solve $\mathbf{Ax} = \mathbf{b}_0$ for a particular \mathbf{b}_0).

Cholesky factorization

When \mathbf{A} is symmetric positive definite, we can take the lower- and upper-triangular factors to be transposes of one another:

$$\mathbf{A} = \mathbf{LL}^T.$$

In this case, \mathbf{L} will not have 1s along the diagonal. Algorithms used to compute Cholesky factorization are similar to those that compute \mathbf{LU} factorizations. The basic idea is that you can use elimination to find a lower-triangular matrix \mathbf{R}_1 that eliminates all but the first entry in the first column of \mathbf{A} :

$$\mathbf{R}_1 \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1N} \\ a_{21} & a_{22} & \cdots & a_{2N} \\ \vdots & & \ddots & \\ a_{N1} & a_{N2} & \cdots & a_{NN} \end{bmatrix} = \begin{bmatrix} \sqrt{a_{11}} & a'_{12} & \cdots & a'_{1N} \\ 0 & a'_{22} & \cdots & a'_{2N} \\ \vdots & & \ddots & \\ 0 & a'_{N2} & \cdots & a'_{NN} \end{bmatrix}.$$

Since \mathbf{A} is symmetric, we can do the same elimination on the columns to get

$$\mathbf{R}_1 \mathbf{A} \mathbf{R}_1^T = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & a''_{22} & \cdots & a''_{2N} \\ \vdots & & \ddots & \\ 0 & a''_{N2} & \cdots & a''_{NN} \end{bmatrix} = \begin{bmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{A}_1 \end{bmatrix}$$

It is easy to see that \mathbf{A}_1 must also be symmetric (and positive definite), so we continue by eliminating all but the first entry in its first

column using another lower-triangular matrix \mathbf{R}_2 , then doing the same to the columns to get

$$\mathbf{R}_2 \mathbf{R}_1 \mathbf{A} \mathbf{R}_1^T \mathbf{R}_2^T = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix}.$$

After N such iterations, we have

$$\mathbf{R}_N \cdots \mathbf{R}_1 \mathbf{A} \mathbf{R}_1^T \cdots \mathbf{R}_N^T = \mathbf{R} \mathbf{A} \mathbf{R}^T = \mathbf{I}.$$

Since the product of two lower-triangular matrices is again lower triangular (show this at home!), \mathbf{R} is again lower-triangular. Since the inverse of a lower-triangular matrix is again lower triangular (show this at home!), we can take

$$\mathbf{A} = \mathbf{L} \mathbf{L}^T, \quad \text{with } \mathbf{L} = \mathbf{R}^{-1}.$$

The Cholesky decomposition is unique. It can actually be computed slightly faster than a general \mathbf{LU} decomposition, and is easier to stabilize.

Example.

$$\begin{bmatrix} 6 & 3 & 0 \\ 3 & 4 & 1 \\ 0 & 1 & 3 \end{bmatrix} = \begin{bmatrix} 2.4495 & 0 & 0 \\ 1.2247 & 1.5811 & 0 \\ 0 & 0.6325 & 1.6125 \end{bmatrix} \begin{bmatrix} 2.4495 & 1.2247 & 0 \\ 0 & 1.5811 & 0.6325 \\ 0 & 0 & 1.6125 \end{bmatrix}$$

QR factorization

The \mathbf{QR} decomposition factors \mathbf{A} as

$$\mathbf{A} = \mathbf{Q} \mathbf{R},$$

where \mathbf{Q} is orthogonal, and \mathbf{R} is upper triangular. It can be computed by running (a stabilized version of) Gram-Schmidt on the columns of \mathbf{A} ; its computational complexity is again $O(N^3)$. Once we have it in hand, solving $\mathbf{A}\mathbf{x} = \mathbf{b}$ has the cost of solving an orthogonal system ($O(N^2)$) and a triangular system ($O(N^2)$).

We will explore this connection more on the next homework.

Like all of the other decompositions in this section, Computing a \mathbf{QR} decomposition costs $O(N^3)$, but once it is in place, we can solve $\mathbf{A}\mathbf{x} = \mathbf{b}$ using $\mathbf{x} = \mathbf{R}^{-1}\mathbf{Q}^T\mathbf{b}$ in $O(N^2)$ time.

Symmetric \mathbf{QR}

When \mathbf{A} is symmetric, we can write

$$\mathbf{A} = \mathbf{Q}\mathbf{T}\mathbf{Q}^T,$$

where \mathbf{Q} is orthonormal, and \mathbf{T} is symmetric and **tri-diagonal**:

$$\mathbf{T} = \begin{bmatrix} t_{11} & t_{12} & 0 & 0 & \cdots & 0 \\ t_{21} & t_{22} & t_{23} & 0 & \cdots & 0 \\ 0 & t_{32} & t_{33} & t_{34} & \cdots & 0 \\ \vdots & & & \ddots & & \vdots \\ 0 & \cdots & & 0 & t_{NN-1} & t_{NN} \end{bmatrix}.$$

This is a handy fact, since in general, tridiagonal matrices are easier to manipulate (invert, compute eigenvalues/eigenvectors of, etc) than general symmetric matrices.

We call this “symmetric QR” since algorithms to compute this decomposition are very similar to those used to compute $\mathbf{A} = \mathbf{QR}$. We

describe one such algorithm, which is very efficient and fairly easy to understand, called *Household tri-diagonalization*.

Let \mathbf{A} be a symmetric matrix

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1N} \\ a_{21} & a_{22} & \cdots & a_{2N} \\ \vdots & & \ddots & \vdots \\ a_{N1} & a_{N2} & \cdots & a_{NN} \end{bmatrix},$$

and set

$$\mathbf{x}_1 = \begin{bmatrix} 0 \\ a_{21} \\ a_{31} \\ \vdots \\ a_{N1} \end{bmatrix}, \quad \mathbf{y}_1 = \begin{bmatrix} 0 \\ r_1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \text{with } r_1 = \|\mathbf{x}_1\|_2.$$

From \mathbf{x}_1 and \mathbf{y}_1 we create the *Householder matrix*

$$\mathbf{H}_1 = \mathbf{I} - 2\mathbf{u}_1\mathbf{u}_1^T, \quad \mathbf{u}_1 = \frac{\mathbf{x}_1 - \mathbf{y}_1}{\|\mathbf{x}_1 - \mathbf{y}_1\|_2}.$$

Then by direct calculation, you can check that

$$\mathbf{A}_2 = \mathbf{H}_1\mathbf{A}\mathbf{H}_1 = \begin{bmatrix} a_{11} & r_1 & 0 & \cdots & 0 \\ r_1 & \tilde{a}_{22} & \tilde{a}_{23} & \cdots & \tilde{a}_{2N} \\ 0 & \tilde{a}_{32} & \tilde{a}_{33} & \cdots & \tilde{a}_{3N} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \tilde{a}_{N2} & \tilde{a}_{N3} & \cdots & \tilde{a}_{NN} \end{bmatrix}.$$

So applying \mathbf{H}_1 on either side eliminates all but two entries in both the first row and first column of \mathbf{A} .

Notice also that the $N \times N$ matrix \mathbf{H}_1 is orthonormal, as

$$\begin{aligned}\mathbf{H}_1^T \mathbf{H}_1 &= (\mathbf{I} - 2\mathbf{u}_1 \mathbf{u}_1^T)(\mathbf{I} - 2\mathbf{u}_1 \mathbf{u}_1^T) \\ &= \mathbf{I} - 4\mathbf{u}_1 \mathbf{u}_1^T + 4\mathbf{u}_1 \mathbf{u}_1^T \mathbf{u}_1 \mathbf{u}_1^T \\ &= \mathbf{I},\end{aligned}$$

where the last equality above follows from the fact that $\mathbf{u}_1^T \mathbf{u}_1 = \|\mathbf{u}_1\|_2^2 = 1$.

We can continue this process on the $(N-1) \times (N-1)$ submatrix on the lower right. We take

$$\mathbf{x}_2 = \begin{bmatrix} 0 \\ 0 \\ \tilde{a}_{32} \\ \tilde{a}_{42} \\ \vdots \\ \tilde{a}_{N2} \end{bmatrix}, \quad \mathbf{y}_2 = \begin{bmatrix} 0 \\ 0 \\ r_2 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \text{with } r_2 = \|\mathbf{x}_2\|_2,$$

and

$$\mathbf{H}_2 = \mathbf{I} - 2\mathbf{u}_2 \mathbf{u}_2^T, \quad \mathbf{u}_2 = \frac{\mathbf{x}_2 - \mathbf{y}_2}{\|\mathbf{x}_2 - \mathbf{y}_2\|_2}.$$

Then

$$\mathbf{A}_3 = \mathbf{H}_2 \mathbf{A}_2 \mathbf{H}_2 = \begin{bmatrix} a_{11} & r_1 & 0 & 0 & \cdots & 0 \\ r_1 & \tilde{a}_{22} & r_2 & 0 & \cdots & 0 \\ 0 & r_2 & \hat{a}_{33} & \hat{a}_{34} & \cdots & \hat{a}_{3N} \\ 0 & 0 & \hat{a}_{43} & \hat{a}_{44} & \cdots & \hat{a}_{4N} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \hat{a}_{N3} & \hat{a}_{N4} & \cdots & \hat{a}_{NN} \end{bmatrix}$$

Repeating this procedure $N-1$ times gives us a tri-diagonal matrix

$$\mathbf{T} = \mathbf{A}_N = \mathbf{H}_N \mathbf{H}_{N-1} \cdots \mathbf{H}_1 \mathbf{A} \mathbf{H}_1 \mathbf{H}_2 \cdots \mathbf{H}_N = \mathbf{Q}^T \mathbf{A} \mathbf{Q}.$$

where

$$\mathbf{Q} = \mathbf{H}_1 \mathbf{H}_2 \cdots \mathbf{H}_N.$$

Note that since all of the \mathbf{H}_n are orthonormal, so is \mathbf{Q} .

SVD and Eigenvalue decompositions

We have already become familiar with the singular value decomposition for general matrices:

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T.$$

When \mathbf{A} is square and invertible ($\text{rank}(\mathbf{A}) = N$), then all of \mathbf{U} , $\mathbf{\Sigma}$, and \mathbf{V} are $N \times N$ and $\mathbf{U} \mathbf{U}^T = \mathbf{U}^T \mathbf{U} = \mathbf{V} \mathbf{V}^T = \mathbf{V}^T \mathbf{V} = \mathbf{I}$. As we have seen, we can solve $\mathbf{A} \mathbf{x} = \mathbf{b}$ with

$$\mathbf{x} = \mathbf{V} \mathbf{\Sigma}^{-1} \mathbf{U}^T \mathbf{b}.$$

We can see that with the SVD in place, the cost of solving a system is $O(N^2)$.

For symmetric \mathbf{A} , we can write

$$\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T,$$

and then $\mathbf{A} \mathbf{x} = \mathbf{b}$ is solved with $\mathbf{x} = \mathbf{V} \mathbf{\Lambda}^{-1} \mathbf{V}^T \mathbf{b}$.

Both of these decompositions represent a matrix as orthogonal-diagonal-orthogonal. Computing either costs $O(N^3)$, and they are slightly more expensive than the \mathbf{QR} and \mathbf{LU} decompositions above. In fact, computing a \mathbf{QR} decomposition is often used as a stepping stone to computing the SVD or eigenvalue decomposition.

Computing eigenvalue decompositions of symmetric matrices

For $N \times N$ symmetric positive semi-definite \mathbf{A} , there are many ways to compute the eigenvalue decomposition $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$. We discuss here one particular technique, popular for its stability, flexibility, and speed, based on *power iterations*.

Power iterations for computing \mathbf{v}_1

To start, let's consider the simpler problem of computing the largest eigenvalue λ_1 and corresponding eigenvector \mathbf{v}_1 of \mathbf{A} . We do this with the following iteration.

Let \mathbf{q}_0 be an arbitrary vector in \mathbb{R}^N with unit norm, $\|\mathbf{q}_0\|_2 = 1$. Then for $k = 1, 2, \dots$ compute

$$\begin{aligned} \mathbf{z}_k &= \mathbf{A}\mathbf{q}_{k-1} \\ \mathbf{q}_k &= \frac{\mathbf{z}_k}{\|\mathbf{z}_k\|_2} \\ \gamma_k &= \mathbf{q}_k^T \mathbf{A}\mathbf{q}_k \end{aligned}$$

Then, as long as \mathbf{q}_0 is not orthogonal to \mathbf{v}_1 , $\langle \mathbf{q}_0, \mathbf{v}_1 \rangle \neq 0$, and $\lambda_1 > \lambda_2$, as k gets large

$$\mathbf{q}_k \rightarrow \mathbf{v}_1 \quad \text{and} \quad \gamma_k \rightarrow \lambda_1.$$

A detailed proof of this, including rates of convergence, can be found in Section 8.2 of the Golub/van Loan book. But we can see roughly why it works with a simple calculation.

The vector \mathbf{q}_k above can be written as

$$\mathbf{q}_k = \frac{\mathbf{A}^k \mathbf{q}_0}{\|\mathbf{A}^k \mathbf{q}_0\|_2}.$$

Since the eigenvectors of \mathbf{A} , $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N$ form an orthobasis for \mathbb{R}^N , we can write

$$\mathbf{q}_0 = \alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2 + \dots + \alpha_N \mathbf{v}_N,$$

for $\alpha_n = \langle \mathbf{q}_0, \mathbf{v}_n \rangle$. Then the expression for \mathbf{q}_k above becomes

$$\mathbf{q}_k = \frac{\alpha_1 \lambda_1^k \mathbf{v}_1 + \alpha_2 \lambda_2^k \mathbf{v}_2 + \dots + \alpha_N \lambda_N^k \mathbf{v}_N}{\sqrt{\alpha_1^2 \lambda_1^{2k} + \alpha_2^2 \lambda_2^{2k} + \dots + \alpha_N^2 \lambda_N^{2k}}}.$$

If $\alpha_1 \neq 0$ and $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_N$, then as k gets large, the first term in each of the sums above will dominate. Thus for large k ,

$$\mathbf{q}_k \approx \frac{\alpha_1 \lambda_1^k \mathbf{v}_1}{\sqrt{\alpha_1^2 \lambda_1^{2k}}} = \mathbf{v}_1.$$

Since $\mathbf{v}_1^T \mathbf{A} \mathbf{v}_1 = \lambda_1$ and $\mathbf{q}_k \approx \mathbf{v}_1$, we also have that $\mathbf{q}_k^T \mathbf{A} \mathbf{q}_k \approx \lambda_1$.

QR iterations

Now consider the problem of computing all the eigenvectors and eigenvalues of \mathbf{A} . We might be tempted to extend the the power method by starting with an entire orthobasis¹ \mathbf{Q}_0 , then take

$$\mathbf{Z}_k = \mathbf{A} \mathbf{Q}_{k-1},$$

¹So \mathbf{Q}_0 is $N \times N$ and satisfies $\mathbf{Q}_0^T \mathbf{Q}_0 = \mathbf{I}$.

and then renormalize the columns of \mathbf{Z}_k to get \mathbf{Q}_k . The problem with this is that all of the columns of \mathbf{Z}_k will converge to \mathbf{v}_1 — this is just running the power method with N different starting points.

What we do instead is “orthonormalize” the columns of \mathbf{Z}_k , we make them orthogonal to each other at every iteration as well as unit norm. This gives us the following iteration:

Let \mathbf{Q}_0 be any orthonormal matrix. For $k = 1, 2, \dots$, take

$$\mathbf{Z}_k = \mathbf{A}\mathbf{Q}_{k-1} \tag{1}$$

$$[\mathbf{Q}_k, \mathbf{R}_k] = \text{qr}(\mathbf{Z}_k) \quad (\text{so } \mathbf{Z}_k = \mathbf{Q}_k\mathbf{R}_k). \tag{2}$$

Using arguments not too different than for the power method, you can show (again, see Golub/van Loan Section 8.2 for details) that

$$\mathbf{Q}_k \rightarrow \mathbf{V}, \quad \text{and} \quad \mathbf{\Gamma}_k = \mathbf{Q}_k^T \mathbf{A} \mathbf{Q}_k \rightarrow \mathbf{\Lambda}.$$

A more popular way to state the iteration (1)–(2) above, and the one you will see in almost every textbook on numerical linear algebra, is the following.

Set $\mathbf{\Gamma}_0 = \mathbf{A}$. Then for $k = 1, 2, \dots$, take

$$[\mathbf{U}_k, \mathbf{R}_k] = \text{qr}(\mathbf{\Gamma}_{k-1}) \quad (\text{so } \mathbf{\Gamma}_{k-1} = \mathbf{U}_k\mathbf{R}_k) \tag{3}$$

$$\mathbf{\Gamma}_k = \mathbf{R}_k\mathbf{U}_k \tag{4}$$

So at each iteration, we are computing a QR factorization, then reversing it (cute!). This gives us the relation

$$\mathbf{R}_{k-1}\mathbf{U}_{k-1} = \mathbf{U}_k\mathbf{R}_k.$$

To see the relationship between version 1 in (1)–(2) and version 2 in

(3)–(4), notice that if we initialize version 1 with $\mathbf{Q}_0 = \mathbf{I}$, then

$$\begin{aligned} \mathbf{A} &= \mathbf{Q}_1 \mathbf{R}_1 \Rightarrow \mathbf{Q}_1 = \mathbf{A} \mathbf{R}_1^{-1} \\ \mathbf{A} \mathbf{Q}_1 &= \mathbf{Q}_2 \mathbf{R}_2 \Rightarrow \mathbf{Q}_2 = \mathbf{A}^2 \mathbf{R}_1^{-1} \mathbf{R}_2^{-1} \\ &\vdots \\ \mathbf{A} \mathbf{Q}_{k-1} &= \mathbf{Q}_k \mathbf{R}_k \Rightarrow \mathbf{Q}_k = \mathbf{A}^k \mathbf{R}_1^{-1} \mathbf{R}_2^{-1} \cdots \mathbf{R}_k^{-1}. \end{aligned}$$

In version 2, we have

$$\begin{aligned} \mathbf{A} &= \mathbf{U}_1 \mathbf{R}_1 \\ \mathbf{A}^2 &= \mathbf{U}_1 \mathbf{R}_1 \mathbf{U}_1 \mathbf{R}_1 = \mathbf{U}_1 \mathbf{U}_2 \mathbf{R}_2 \mathbf{R}_1 \\ \mathbf{A}^3 &= \mathbf{U}_1 \mathbf{R}_1 \mathbf{U}_1 \mathbf{U}_2 \mathbf{R}_2 \mathbf{R}_1 = \mathbf{U}_1 \mathbf{U}_2 \mathbf{R}_2 \mathbf{U}_2 \mathbf{R}_2 \mathbf{R}_1 = \mathbf{U}_1 \mathbf{U}_2 \mathbf{U}_3 \mathbf{R}_3 \mathbf{R}_2 \mathbf{R}_1 \\ &\vdots \\ \mathbf{A}^k &= \mathbf{U}_1 \mathbf{U}_2 \cdots \mathbf{U}_k \mathbf{R}_k \mathbf{R}_{k-1} \cdots \mathbf{R}_1, \end{aligned}$$

which is the same thing as version 1 with

$$\mathbf{Q}_k = \mathbf{U}_1 \mathbf{U}_2 \cdots \mathbf{U}_k.$$

To keep track of the eigenvectors, we can restate the QR algorithm as follows.

Let $\mathbf{\Gamma}_0 = \mathbf{A}$, $\mathbf{Q}_0 = \mathbf{I}$. For $k = 1, 2, \dots$, do

$$\begin{aligned} [\mathbf{U}_k, \mathbf{R}_k] &= \text{qr}(\mathbf{\Gamma}_{k-1}) \quad (\text{so } \mathbf{\Gamma}_{k-1} = \mathbf{U}_k \mathbf{R}_k) \\ \mathbf{\Gamma}_k &= \mathbf{R}_k \mathbf{U}_k \\ \mathbf{Q}_k &= \mathbf{U}_k \mathbf{Q}_{k-1} \end{aligned}$$

Then $\mathbf{\Gamma}_k \rightarrow \mathbf{\Lambda}$ and $\mathbf{Q}_k \rightarrow \mathbf{V}$.

Comments on computational complexity

The methods above have been the object of intense study over the past 50-60 years, and their cost and stability is very well understood. Notice that all of the methods cost $O(N^3)$ in the general case — this is the essential cost of solving a system of linear equations.

When N is small, $O(N^3)$ is OK. For example, I have a nice desktop system (two 2.93 GHz Intel Xeons with 6 cores each, 32 GB memory), and this about how long it takes MATLAB to solve $\mathbf{Ax} = \mathbf{b}$ for different N :

$N = 100$, 0.0007 seconds ($700\mu\text{s}$)

$N = 1000$, 0.009 seconds (9ms)

$N = 5000$, 0.550 seconds

$N = 10,000$, 3.35 seconds

$N = 20,000$, 24.53 seconds

There are many applications where N is in the millions (or even billions). In these situations, solving $\mathbf{Ax} = \mathbf{b}$ directly is infeasible. You can roughly divide problems into three categories:

Small scale. $N \lesssim 10^3$. Here $O(N^3)$ algorithms are OK, and exact algorithms are appropriate.

Medium scale. $N \sim 10^4$. Here $O(N^3)$ is not OK, $O(N^2)$ is OK. It is getting hard to even store the matrix in memory at this point.

Large scale. $N \gtrsim 10^5$. Here $O(N^2)$ is not OK, $O(N^3)$ is unthinkable. We need algorithms that are $O(N)$ or $O(N \log N)$, possibly at the cost of finding an exact solution.

We will start by looking at certain types of **structured** symmetric matrices. This structure allows us to solve $\mathbf{Ax} = \mathbf{b}$ significantly faster than $O(N^3)$.

After this, we will look at two different **iterative algorithms** for finding an appropriate solution to $\mathbf{Ax} = \mathbf{b}$ when \mathbf{A} is sym+def. These algorithms have the nice feature that the matrix \mathbf{A} does not need to be held in memory — all we need is a “black box” that computes \mathbf{Ax} given \mathbf{x} as input. This is especially nice if you have a fast implicit method for computing \mathbf{Ax} (e.g \mathbf{A} involved FFTs, or is sparse).

Structured matrices

We will look at four types of structured matrices, although plenty of other types exist:

- identity+low rank
- circulant
- Toeplitz (we will actually skip this)²
- banded (homework)

In each of these cases, the structure of the system allows us to do a solve in much better than $O(N^3)$ operations.

Identity + low rank

Consider a system of the form

$$(\gamma \mathbf{I} + \mathbf{B}\mathbf{B}^T)\mathbf{x} = \mathbf{b},$$

where $\gamma > 0$ is some scalar, and \mathbf{B} is a $N \times R$ matrix with $R < N$. These types of systems are prevalent in array signal processing and machine learning. We will see that if $R \ll N$, this system can be solved in (much) faster than $O(N^3)$ time.

Note that while $\mathbf{B}\mathbf{B}^T$ is not at all invertible (since it is rank deficient), $\gamma \mathbf{I} + \mathbf{B}\mathbf{B}^T$ will be. To see this, set

$$\mathbf{z} = \mathbf{B}^T \mathbf{x}, \quad \mathbf{z} \in \mathbb{R}^R.$$

²See <https://goo.gl/uCTcLQ> if you are interested

Then we can solve the system by jointly solving for \mathbf{x} and \mathbf{z} :

$$\gamma \mathbf{x} + \mathbf{B} \mathbf{z} = \mathbf{b} \tag{5}$$

$$\mathbf{B}^T \mathbf{x} - \mathbf{z} = \mathbf{0}. \tag{6}$$

Solving the first equations (5) yields

$$\mathbf{x} = \gamma^{-1}(\mathbf{b} - \mathbf{B} \mathbf{z}),$$

and then plugging this into (6) gives us

$$\begin{aligned} \gamma^{-1} \mathbf{B}^T (\mathbf{b} - \mathbf{B} \mathbf{z}) - \mathbf{z} &= \mathbf{0} \\ \Rightarrow (\gamma \mathbf{I} + \mathbf{B}^T \mathbf{B}) \mathbf{z} &= \mathbf{B}^T \mathbf{b} \end{aligned}$$

and so

$$\mathbf{z} = (\gamma \mathbf{I} + \mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}^T \mathbf{b}.$$

But notice that this is an $R \times R$ system of equations.

So it takes $O(NR^2)$ to construct $\gamma \mathbf{I} + \mathbf{B}^T \mathbf{B}$,
then $O(R^3)$ to solve for \mathbf{z} ,
then $O(NR)$ to calculate $\mathbf{B} \mathbf{z}$ (and hence find \mathbf{x}).

The dominant cost in all of this is $O(NR^2)$, which is much less than $O(N^3)$ if $R \ll N$.

Circulant systems.

A circulant matrix has the form

$$\mathbf{H} = \begin{bmatrix} h_0 & h_{N-1} & h_{N-2} & \cdots & h_1 \\ h_1 & h_0 & h_{N-1} & \cdots & h_2 \\ h_2 & h_1 & h_0 & & \vdots \\ \vdots & \vdots & & \ddots & h_{N-1} \\ h_{N-1} & & & h_1 & h_0 \end{bmatrix}$$

For \mathbf{H} symmetric, we have $h_k = h_{N-k}$ for $k = 1, \dots, N-1$, although symmetry does not play too big a role in exploiting this structure.

Circulant matrices have two very nice properties:

- We know their eigenvectors already — they are the discrete harmonic sinusoids (i.e. the columns of the $N \times N$ DFT matrix).
- Transforming into the eigenbasis is **fast** thanks to the FFT (which is $O(N \log N)$).

We can write

$$\mathbf{H} = \mathbf{F} \mathbf{\Lambda} \mathbf{F}^H, \quad F[m, n] = \frac{1}{\sqrt{N}} e^{j2\pi mn/N}$$

and

$$\mathbf{H}^{-1} = \mathbf{F} \mathbf{\Lambda}^{-1} \mathbf{F}^H,$$

so

$$\mathbf{H}^{-1} \mathbf{b} = \underbrace{\mathbf{F}}_{\text{FFT, } O(N \log N)} \underbrace{\mathbf{\Lambda}^{-1}}_{\text{diagonal weighting, } O(N)} \underbrace{\mathbf{F}^H \mathbf{b}}_{\text{FFT, } O(N \log N)}$$

\Rightarrow solving an $N \times N$ system of equations can be done in $O(N \log N)$ time!

This is **fast** compared to $O(N^3)$ — on my computer, I can solve a system like this in $N = 20000$ in $800\mu\text{s}$ (compare to 24 seconds for the general case).