

# III. Solving and Analyzing Least-Squares Problems

## Solving systems of symmetric equations

In the previous chapter, we saw a few examples of how to set up least-squares problems in machine learning. We saw that finding the least-squares solution always amounted to solving a symmetric system of linear equations. In this section, we will discuss how to solve, and analyze the solution to, systems which are square and symmetric.

For all of this set of notes, we will consider  $\mathbf{A}$  which are  $N \times N$  (square) and symmetric (or Hermitian for complex-valued  $\mathbf{A}$ ).

**Definition:** If  $\mathbf{A}$  is real-valued, then we call it **symmetric** if  $\mathbf{A}^T = \mathbf{A}$ . Example:

$$\mathbf{A} = \begin{bmatrix} 1 & 3 & 7 \\ 3 & -5 & -2 \\ 7 & -2 & 6 \end{bmatrix}$$

In other words,  $\mathbf{A}^T = \mathbf{A}$  means that  $A[m, n] = A[n, m]$  for all  $m, n = 1, \dots, N$ .

**Definition:** If  $\mathbf{A}$  is complex-valued, then we call it **Hermitian** if  $\mathbf{A}^H = \mathbf{A}$ . Recall that  $\mathbf{A}^H$  is the conjugate-transpose of  $\mathbf{A}$ : we exchange rows and columns and then conjugate all the entries. Example:

$$\mathbf{A} = \begin{bmatrix} 1 & 3 + j2 & 1 - j3 \\ 3 - j2 & -5 & 2 + j4 \\ 1 + j3 & 2 - j4 & 6 \end{bmatrix}$$

In other words,  $\mathbf{A} = \mathbf{A}^H$  means that  $A[m, n] = \overline{A[n, m]}$ . Of course, if  $\mathbf{A}$  is real-valued, then  $\mathbf{A}^H = \mathbf{A}^T$ .

In the end, we will be interested in the general case of non-symmetric and even non-square matrices, but symmetric/Hermitian matrices will be the starting point of our study — they actually play a fundamental role in solving general “least-squares” problems and our analysis of them here will prove very useful later on.

The main mathematical construct we will use to understand the solution to symmetric systems of equations is the **eigenvalue decomposition**.

## Eigenvalue decompositions of symmetric matrices

**Definition:** An **eigenvector** of an  $N \times N$  matrix is a vector  $\mathbf{v}$  such that

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$$

for some  $\lambda \in \mathbb{C}$ . The scalar  $\lambda$  is called the **eigenvalue** associated with  $\mathbf{v}$ .

A matrix  $\mathbf{A}$  is called **diagonalizable** if it has  $N$  linearly independent eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_N$ . In this case, we can write the eigenvalue relations

$$\begin{aligned}\mathbf{A}\mathbf{v}_1 &= \lambda_1\mathbf{v}_1 \\ \mathbf{A}\mathbf{v}_2 &= \lambda_2\mathbf{v}_2 \\ &\vdots \\ \mathbf{A}\mathbf{v}_N &= \lambda_N\mathbf{v}_N\end{aligned}$$

as

$$\mathbf{A} \underbrace{\begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_N \end{bmatrix}}_{\mathbf{V}} = \underbrace{\begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_N \end{bmatrix}}_{\mathbf{V}} \underbrace{\begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_N \end{bmatrix}}_{\mathbf{\Lambda}}$$

or more compactly

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

and since the  $\mathbf{v}_n$  are linearly independent,  $\mathbf{V}^{-1}$  exists, and we can write

$$\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^{-1}.$$

Not all matrices are diagonalizable, but it happens that all symmetric matrices are. This fact is a consequence of the Schur Triangularity Lemma, which states that **any**  $N \times N$  matrix is “unitarily similar” to an upper-triangular matrix. That is, for a given  $\mathbf{A}$ , there is an orthonormal matrix<sup>1</sup>  $\mathbf{V}$  such that

$$\mathbf{A} = \mathbf{V}\mathbf{\Delta}\mathbf{V}^H,$$

where

$$\mathbf{\Delta} = \begin{bmatrix} \Delta[1,1] & \Delta[1,2] & \cdots & \Delta[1,N] \\ 0 & \Delta[2,2] & \cdots & \Delta[2,N] \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \Delta[N,N] \end{bmatrix}.$$

(For the interested reader, we will prove the Schur Triangularity Lemma in the “Technical Details” section at the end of these notes.)

---

<sup>1</sup>We call a  $N \times N$  matrix  $\mathbf{V}$  *orthonormal* if its columns are orthogonal to one another and have unit norm, so  $\mathbf{V}^H\mathbf{V} = \mathbf{I}$ . This means, of course, that  $\mathbf{V}^{-1} = \mathbf{V}^H$ .

If  $\mathbf{A}$  is Hermitian, then  $\mathbf{A} = \mathbf{A}^H$  implies

$$\mathbf{V}\mathbf{\Delta}\mathbf{V}^H = \left(\mathbf{V}\mathbf{\Delta}\mathbf{V}^H\right)^H = \mathbf{V}\mathbf{\Delta}^H\mathbf{V}^H$$

and so  $\mathbf{\Delta} = \mathbf{\Delta}^H$ . Since  $\mathbf{\Delta}$  is upper-triangular, this means it must also be diagonal and real:

$$\mathbf{\Delta} = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & \cdots & \cdots & \lambda_N \end{bmatrix}, \quad \lambda_n \in \mathbb{R}.$$

Thus every Hermitian matrix is **diagonalized** by an orthonormal matrix; it is unitarily similar to a real-valued diagonal matrix. It is also not hard to see that if  $\mathbf{v}_n$  is a column of  $\mathbf{V}$ , then

$$\mathbf{A}\mathbf{v}_n = \mathbf{V}\mathbf{\Delta}\mathbf{V}^H\mathbf{v}_n = \mathbf{V}\mathbf{\Delta}\mathbf{e}_n = \lambda_n\mathbf{V}\mathbf{e}_n = \lambda_n\mathbf{v}_n,$$

where  $\mathbf{e}_n$  is the  $n$ th standard basis vector ( $e_n[k] = 1$  for  $k = n$  and is zero elsewhere). This means we can interpret  $\mathbf{A} = \mathbf{V}\mathbf{\Delta}\mathbf{V}^H$  as an *eigen-decomposition* of  $\mathbf{A}$ . Hermitian  $N \times N$  matrices have  $N$  orthogonal eigenvectors and real eigenvalues.

We will sometimes find it convenient to write the eigenvalue decomposition as a weighted sum of outer products between each eigenvector with itself:

$$\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^H = \sum_{n=1}^N \lambda_n \mathbf{v}_n \mathbf{v}_n^H.$$

An  $N \times N$  symmetric/Hermitian matrix  $\mathbf{A}$  has:

- Real eigenvalues (even if  $\mathbf{A}$  is itself complex valued),  $\lambda_1, \dots, \lambda_N$ .
- $N$  orthogonal eigenvectors,  $\mathbf{v}_1, \dots, \mathbf{v}_N$ .
- If  $\mathbf{A}$  is real-valued, then the  $\mathbf{v}_n$  can also be chosen to be real-valued.

We can decompose real-valued  $\mathbf{A}$  as

$$\mathbf{A} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T = \sum_{n=1}^N \lambda_n \mathbf{v}_n \mathbf{v}_n^T,$$

where  $\mathbf{V}$  contains the eigenvectors as columns, so  $\mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}$ , and  $\mathbf{\Lambda}$  contains the eigenvalues along its diagonal. If  $\mathbf{A}$  is complex-valued and Hermitian, simply replace  $\mathbf{V}^T$  and  $\mathbf{v}^T$  with  $\mathbf{V}^H$  and  $\mathbf{v}^H$  above.

For the remainder of these notes, we will concentrate below on the real-valued/symmetric case. But there is a straightforward extension to everything we say to the complex-valued/Hermitian case.

One way to think about  $\mathbf{V}$  is as a **transform** which greatly simplifies the action of  $\mathbf{A}$ :

$$\begin{aligned} \mathbf{A} \mathbf{x} &= \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \mathbf{x} \\ &= (\text{inverse } \mathbf{V} \text{ transform})(\text{pointwise multiply})(\mathbf{V} \text{ transform}) [\mathbf{x}]. \end{aligned}$$

## Example

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

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

$$\mathbf{V} = [\mathbf{v}_1 \quad \mathbf{v}_2] = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \quad \text{and} \quad \mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

Sketch the action of  $\mathbf{A}$  on  $\mathbf{x} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ .

Calculating  $\boldsymbol{\alpha} = \mathbf{V}^T \mathbf{x} = \begin{bmatrix} \langle \mathbf{x}, \mathbf{v}_1 \rangle \\ \langle \mathbf{x}, \mathbf{v}_2 \rangle \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$  tells us that

$$\mathbf{x} = \frac{1}{\sqrt{2}} \mathbf{v}_1 + \frac{1}{\sqrt{2}} \mathbf{v}_2.$$

The eigenvalue decomposition show us that applying  $\mathbf{A}$  stretches the  $\mathbf{v}_1$  component by 2 and the  $\mathbf{v}_2$  component by 1 to get

$$\begin{aligned} \mathbf{Ax} &= \lambda_1 \langle \mathbf{x}, \mathbf{v}_1 \rangle \mathbf{v}_1 + \lambda_2 \langle \mathbf{x}, \mathbf{v}_2 \rangle \mathbf{v}_2 \\ &= \sqrt{2} \mathbf{v}_1 + \frac{1}{\sqrt{2}} \mathbf{v}_2 \\ &= \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} -1/2 \\ 1/2 \end{bmatrix} = \begin{bmatrix} 1/2 \\ 3/2 \end{bmatrix} \end{aligned}$$

## Symmetric positive definite matrices

**Definition:** A symmetric matrix  $\mathbf{A}$  is called **positive definite** if it has positive eigenvalues,

$$\lambda_n > 0 \quad \text{for } n = 1, \dots, N.$$

We will sometimes abbreviate this as **sym+def**. We will also call  $\mathbf{A}$  **positive semi-definite** if  $\lambda_n \geq 0$ .

sym+def  $\mathbf{A}$  are invertible, and obey

$$\mathbf{x}^T \mathbf{A} \mathbf{x} > 0 \quad \text{for all } \mathbf{x} \in \mathbb{R}^N.$$

We will use the typical convention for sym+def matrices of ordering the eigenvalues largest to smallest, so

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N > 0$$

## Variational form for extreme eigenvalues

For sym+def  $\mathbf{A}$ , there is a variational expression for the largest and smallest eigenvalues:

$$\max_{\mathbf{x} \in \mathbb{R}^N} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\|\mathbf{x}\|_2^2} = \max_{\substack{\mathbf{x} \in \mathbb{R}^N \\ \|\mathbf{x}\|_2=1}} \mathbf{x}^T \mathbf{A} \mathbf{x} = \lambda_1, \quad (\text{the maximizer is } \mathbf{v}_1)$$

$$\min_{\mathbf{x} \in \mathbb{R}^N} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\|\mathbf{x}\|_2^2} = \min_{\substack{\mathbf{x} \in \mathbb{R}^N \\ \|\mathbf{x}\|_2=1}} \mathbf{x}^T \mathbf{A} \mathbf{x} = \lambda_N, \quad (\text{the minimizer is } \mathbf{v}_N)$$



These identities can be verified using the fact that for any  $\mathbf{v} \in \mathbb{R}^N$ ,

$$\|\mathbf{V}^T \mathbf{x}\|_2^2 = \mathbf{x}^T \mathbf{V} \mathbf{V}^T \mathbf{x} = \mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|_2^2,$$

and so for example

$$\max_{\substack{\mathbf{x} \in \mathbb{R}^N \\ \|\mathbf{x}\|_2=1}} \mathbf{x}^T \mathbf{A} \mathbf{x} = \max_{\substack{\mathbf{x} \in \mathbb{R}^N \\ \|\mathbf{x}\|_2=1}} \mathbf{x}^T \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \mathbf{x} = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^N \\ \|\boldsymbol{\alpha}\|_2=1}} \boldsymbol{\alpha}^T \mathbf{\Lambda} \boldsymbol{\alpha} = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^N \\ \|\boldsymbol{\alpha}\|_2=1}} \sum_{n=1}^N |\alpha[n]|^2 \lambda_n,$$

where in the second equality above we have made the substitution  $\boldsymbol{\alpha} = \mathbf{V}^T \mathbf{x}$ . A moment's thought tells us that the sum above is largest, under the constraint that  $\sum_{n=1}^N |\alpha[n]|^2 = 1$ , for

$$\alpha[n] = \begin{cases} 1 & n = 1 \\ 0 & \text{otherwise} \end{cases}.$$

(Recall that  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N > 0$ .)

## Sylvester's Matrix Theorem

As we have seen above, another way to write  $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$  is as a weighted sum of *outer products* of each column of  $\mathbf{V}$  with itself:

$$\mathbf{A} = \sum_{n=1}^N \lambda_n \mathbf{v}_n \mathbf{v}_n^T.$$

This form, along with the orthogonality of the  $\mathbf{v}_n$ , makes it easy what happens to the eigenvalues when we square the matrix

$$\begin{aligned} \mathbf{A}^2 &= \left( \sum_{n=1}^N \lambda_n \mathbf{v}_n \mathbf{v}_n^T \right) \left( \sum_{\ell=1}^N \lambda_{\ell} \mathbf{v}_{\ell} \mathbf{v}_{\ell}^T \right) \\ &= \sum_{n=1}^N \sum_{\ell=1}^N \lambda_n \lambda_{\ell} \mathbf{v}_n (\mathbf{v}_n^T \mathbf{v}_{\ell}) \mathbf{v}_{\ell} \\ &= \sum_{n=1}^N \lambda_n^2 \mathbf{v}_n \mathbf{v}_n^T \quad (\text{Since } \mathbf{v}_n^T \mathbf{v}_{\ell} = 0 \text{ unless } n = \ell.) \\ &= \mathbf{V} \mathbf{\Lambda}^2 \mathbf{V}^T. \end{aligned}$$

It does not take too much imagination to see that this extends to any positive integer-valued power  $p$ ,

$$\mathbf{A}^p = \sum_{n=1}^N \lambda_n^p \mathbf{v}_n \mathbf{v}_n^T,$$

and thus to any polynomial function of  $\mathbf{A}$ :

$$c_p \mathbf{A}^p + c_{p-1} \mathbf{A}^{p-1} + \cdots + c_1 \mathbf{A} + c_0 \mathbf{I} = \sum_{n=1}^N (c_p \lambda_n^p + c_{p-1} \lambda_n^{p-1} + \cdots + c_1 \lambda_n + c_0) \mathbf{v}_n \mathbf{v}_n^T.$$

From polynomials, we can move to any analytic function using a *Taylor expansion*,

$$f(\mathbf{A}) = \sum_{n=1}^N f(\lambda_n) \cdot \mathbf{v}_n \mathbf{v}_n^T.$$

This is known as **Sylvester's matrix theorem**.

Here are several important examples of Sylvester's matrix theorem in action:

1. Inverting a matrix. Take  $f(x) = x^{-1}$ ,

$$\mathbf{A}^{-1} = \sum_{n=1}^N \frac{1}{\lambda_n} \cdot \mathbf{v}_n \mathbf{v}_n^T = \mathbf{V} \mathbf{\Lambda}^{-1} \mathbf{V}^T.$$

Of course, for this expression to make sense, all of the eigenvalues  $\lambda_n$  must be non-zero.

2. Taking the square root of a positive matrix. If  $\mathbf{A}$  is non-negative in that all of the eigenvalues  $\lambda_n$  are greater than or equal to zero, then we can write

$$\mathbf{A}^{1/2} = \sum_{n=1}^N \sqrt{\lambda_n} \cdot \mathbf{v}_n \mathbf{v}_n^T = \mathbf{V} \mathbf{\Lambda}^{1/2} \mathbf{V}^T.$$

3. Matrix exponential,

$$e^{\mathbf{A}} = \sum_{n=1}^N e^{\lambda_n} \cdot \mathbf{v}_n \mathbf{v}_n^T.$$

(As an aside, the matrix exponential is of paramount importance when studying *dynamical systems* — if we have a system of linear homogenous first-order differential equations with

constant coefficients

$$\begin{aligned}\mathbf{x}'(t) &= \mathbf{A}\mathbf{x}(t) \\ \mathbf{x}(0) &= \mathbf{x}_0,\end{aligned}$$

the solution is given by  $\mathbf{x}(t) = e^{\mathbf{A}t}\mathbf{x}_0$ .)

## Solving systems of sym+def equations

Given  $\mathbf{y}$ , we are interested in finding  $\mathbf{x}$  such that

$$\mathbf{y} = \mathbf{A}\mathbf{x}.$$

When  $\mathbf{A}$  is sym+def, it is invertible, so this system has a unique solution:

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{y}$$

which we can also write as

$$\begin{aligned}\mathbf{x} &= \mathbf{V}\mathbf{\Lambda}^{-1}\mathbf{V}^T\mathbf{y} \\ &= (\text{inverse } \mathbf{V} \text{ transform})(\text{pointwise multiply})(\mathbf{V} \text{ transform})[\mathbf{y}].\end{aligned}$$

We can also write

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{y} = \sum_{n=1}^N \frac{1}{\lambda_n} \langle \mathbf{y}, \mathbf{v}_n \rangle \mathbf{v}_n.$$

Notice that once we have the eigenvalue decomposition in hand, solving  $\mathbf{y} = \mathbf{A}\mathbf{x}$  simply amounts to a matrix-vector multiply.

Now suppose that there is some observation error:

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e},$$

where  $\mathbf{e}$  is an unknown error vector in  $\mathbb{R}^N$ . We reconstruct as before, by applying  $\mathbf{A}^{-1}$  to  $\mathbf{y}$ :

$$\begin{aligned}\tilde{\mathbf{x}} &= \mathbf{A}^{-1}\mathbf{y} = \mathbf{A}^{-1}(\mathbf{A}\mathbf{x} + \mathbf{e}) \\ &= \mathbf{x} + \mathbf{A}^{-1}\mathbf{e}.\end{aligned}$$

The **reconstruction error** is

$$\tilde{\mathbf{x}} - \mathbf{x} = \mathbf{A}^{-1}\mathbf{e}.$$

### Questions:

1. What is the largest eigenvalue of  $\mathbf{A}^{-1}$ ?
2. What is an upper bound on the reconstruction error energy?

$$\begin{aligned}\|\tilde{\mathbf{x}} - \mathbf{x}\|_2^2 &= \|\mathbf{A}^{-1}\mathbf{e}\|_2^2 \\ &= \mathbf{e}^T (\mathbf{A}^{-1})^T \mathbf{A}^{-1} \mathbf{e}\end{aligned}$$

$$\leq \underline{\hspace{2cm}}$$

3. What is the smallest eigenvalue of  $\mathbf{A}^{-1}$ ?
4. What is a lower bound on the construction error?

In the end, we have

$$\frac{1}{\lambda_1^2} \|\mathbf{e}\|_2^2 \leq \|\tilde{\mathbf{x}} - \mathbf{x}\|_2^2 \leq \frac{1}{\lambda_N^2} \|\mathbf{e}\|_2^2.$$

## Average reconstruction error

The maximum and minimum eigenvalues give us the “best case” and “worst case” errors,

$$\tilde{\mathbf{x}} - \mathbf{x} = \mathbf{V} \mathbf{\Lambda}^{-1} \mathbf{V}^T \mathbf{e},$$

and so

$$\frac{1}{\lambda_1^2} \|\mathbf{e}\|_2^2 \leq \|\tilde{\mathbf{x}} - \mathbf{x}\|_2^2$$

where the lower bound is achieved for  $\mathbf{e} \in \text{span}(\mathbf{v}_1)$ . Similarly,

$$\|\tilde{\mathbf{x}} - \mathbf{x}\|_2^2 \leq \frac{1}{\lambda_N^2} \|\mathbf{e}\|_2^2$$

where the upper bound is achieved for  $\mathbf{e} \in \text{span}(\mathbf{v}_N)$ .

We can also get an “average case” reconstruction error when  $\mathbf{e}$  is **generic** (i.e. **random**). Our model for random noise is that the entries of  $\mathbf{e}$  are iid Gaussian:

$$e[n] \sim \text{Normal}(0, \sigma^2), \quad n = 1, 2, \dots, N$$

$$\mathbb{E}[e[n]e[\ell]] = \begin{cases} \sigma^2 & n = \ell \\ 0 & n \neq \ell \end{cases}.$$

What is the mean-energy of  $\mathbf{e}$ ? Answer:

$$\begin{aligned} \mathbb{E}[\|\mathbf{e}\|_2^2] &= \mathbb{E}\left[\sum_{n=1}^N |e[n]|^2\right] \\ &= \sum_{n=1}^N \mathbb{E}[|e[n]|^2] \\ &= \sum_{n=1}^N \sigma^2 \\ &= N\sigma^2 \end{aligned}$$

Since  $\mathbf{e}$  is random,  $\tilde{\mathbf{x}} - \mathbf{x}$  is a random vector. The expected reconstruction energy is

$$\mathbb{E}[\|\tilde{\mathbf{x}} - \mathbf{x}\|_2^2] = \mathbb{E}[\|\mathbf{A}^{-1}\mathbf{e}\|_2^2] .$$

We can get a nice expression for this by using the eigenvalue description. Just plug and chug<sup>2</sup>:

---

<sup>2</sup>There is a fact that we are using here time and time again: if  $\mathbf{u}$  and  $\mathbf{v}$  are vectors and  $\mathbf{Z}$  is a matrix, then  $\langle \mathbf{v}, \mathbf{Z}\mathbf{u} \rangle = \langle \mathbf{Z}^T \mathbf{v}, \mathbf{u} \rangle$ .



$$\begin{aligned}
\mathbb{E} [\|\mathbf{A}^{-1}\mathbf{e}\|_2^2] &= \mathbb{E} [\|\mathbf{V}\mathbf{\Lambda}^{-1}\mathbf{V}^T\mathbf{e}\|_2^2] \\
&= \mathbb{E} [\langle \mathbf{V}\mathbf{\Lambda}^{-1}\mathbf{V}^T\mathbf{e}, \mathbf{V}\mathbf{\Lambda}^{-1}\mathbf{V}^T\mathbf{e} \rangle] \\
&= \mathbb{E} [\langle \mathbf{\Lambda}^{-1}\mathbf{V}^T\mathbf{e}, \mathbf{V}^T\mathbf{V}\mathbf{\Lambda}^{-1}\mathbf{V}^T\mathbf{e} \rangle] \\
&= \mathbb{E} [\langle \mathbf{\Lambda}^{-1}\mathbf{V}^T\mathbf{e}, \mathbf{\Lambda}^{-1}\mathbf{V}^T\mathbf{e} \rangle] \quad (\text{since } \mathbf{V}^T\mathbf{V} = \mathbf{I}) \\
&= \mathbb{E} \left[ \sum_{n=1}^N \left| \frac{1}{\lambda_n} \langle \mathbf{e}, \mathbf{v}_n \rangle \right|^2 \right] \\
&= \sum_{n=1}^N \frac{1}{\lambda_n^2} \mathbb{E} [|\langle \mathbf{e}, \mathbf{v}_n \rangle|^2] .
\end{aligned}$$

Note that

$$\begin{aligned}
|\langle \mathbf{e}, \mathbf{v}_n \rangle|^2 &= \left( \sum_{m=1}^N v_n[m] e[m] \right)^2 \\
&= \sum_{m=1}^N \sum_{\ell=1}^N v_n[m] v_n[\ell] e[m] e[\ell],
\end{aligned}$$

so

$$\begin{aligned}
\mathbb{E} [|\langle \mathbf{e}, \mathbf{v}_n \rangle|^2] &= \sum_{m=1}^N \sum_{\ell=1}^N v_n[m] v_n[\ell] \underbrace{\mathbb{E} [e[m] e[\ell]]}_{=0 \text{ unless } m=\ell} \\
&= \sum_{m=1}^N |v_n[m]|^2 \sigma^2 \\
&= \sigma^2, \quad \text{since } \|\mathbf{v}_n\|_2^2 = 1.
\end{aligned}$$

Thus

$$\begin{aligned} \mathbb{E} [\|\mathbf{A}^{-1}\mathbf{e}\|_2^2] &= \sigma^2 \sum_{n=1}^N \frac{1}{\lambda_n^2} \\ &= \underbrace{N\sigma^2}_{\mathbb{E} \|\mathbf{e}\|_2^2} \underbrace{\left( \frac{1}{N} \sum_{n=1}^N \frac{1}{\lambda_n^2} \right)}_{\text{average of } \lambda_n^{-2}}. \end{aligned}$$

The intuition here is that a random vector is spread out more or less equally over the basis  $\{\mathbf{v}_n\}$  rather than being concentrated in  $\text{span}(\mathbf{v}_1)$  or  $\text{span}(\mathbf{v}_n)$ .

## Summary of sym+def reconstruction

**Observe**

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e},$$

where  $\mathbf{A}$  is sym+def.

**Reconstruct**

$$\tilde{\mathbf{x}} = \mathbf{A}^{-1}\mathbf{y} = \mathbf{x} + \mathbf{A}^{-1}\mathbf{e}.$$

**Best/Worst case reconstruction errors**

$$\frac{1}{\lambda_1^2} \|\mathbf{e}\|_2^2 \leq \|\tilde{\mathbf{x}} - \mathbf{x}\|_2^2 \leq \frac{1}{\lambda_N^2} \|\mathbf{e}\|_2^2$$

The max (min) is obtained for  $\mathbf{e}$  pointing in the same direction as  $\mathbf{v}_N$  ( $\mathbf{v}_1$ ).

**Average reconstruction error**

$$\begin{aligned} e[n] &\sim \text{Normal}(0, \sigma^2), \quad e[n] \text{ iid} \\ \text{E} [\|\tilde{\mathbf{x}} - \mathbf{x}\|_2^2] &= \sigma^2 \sum_{n=1}^N \frac{1}{\lambda_n^2} \\ &= \left( \frac{1}{N} \sum_{n=1}^N \frac{1}{\lambda_n^2} \right) \cdot \text{E} [\|\mathbf{e}\|_2^2] \end{aligned}$$

## Technical details: Schur decomposition

In this section we prove one of the fundamental results in linear algebra: that any  $N \times N$  matrix is *unitarily similar* to an upper-triangular matrix. That is, given an  $N \times N$  matrix  $\mathbf{A}$ , there is an orthonormal matrix  $\mathbf{V}$  (meaning  $\mathbf{V}^H \mathbf{V} = \mathbf{I}$ ) such that

$$\mathbf{A} = \mathbf{V} \mathbf{\Delta} \mathbf{V}^H,$$

where

$$\mathbf{\Delta} = \begin{bmatrix} \Delta[1, 1] & \Delta[1, 2] & \cdots & \Delta[1, N] \\ 0 & \Delta[2, 2] & \cdots & \Delta[2, N] \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \Delta[N, N] \end{bmatrix}.$$

This is known as the **Schur Decomposition** or the **Schur Triangulation**. It is also possible to choose  $\mathbf{V}$  so that  $\mathbf{\Delta}$  is lower-triangular.

The proof works by induction. First, we use the fact that every matrix has at least one eigenvector. Let  $\mathbf{v}_1$  be an eigenvector of  $\mathbf{A}$ ; we may assume that  $\mathbf{v}_1$  is normalized, since all scalar multiples of eigenvectors are also eigenvectors. Then we take  $\mathbf{V}_1$  to be any orthogonal matrix with  $\mathbf{v}_1$  as one of its columns:

$$\mathbf{V}_1 = [\mathbf{v}_1 \ \mathbf{U}_1], \quad \mathbf{U}_1 \in \mathbb{R}^{N \times N-1}, \quad \mathbf{U}_1^H \mathbf{U}_1 = \mathbf{I}, \quad \mathbf{U}_1^H \mathbf{v}_1 = \mathbf{0}.$$

This is equivalent to finding an orthobasis for  $\mathbb{R}^N$  where  $\mathbf{v}_1$  is one of the basis vectors and the  $N - 1$  columns of  $\mathbf{U}_1$  are the others. There are many such choices for  $\mathbf{U}_1$ ; one can be found using the Gram-Schmidt algorithm.

Since  $\mathbf{v}_1$  is an eigenvector of  $\mathbf{A}$  (call the corresponding eigenvalue  $\lambda_1$ ),

$$\mathbf{A} \mathbf{V}_1 = [\lambda_1 \mathbf{v}_1 \ \mathbf{A} \mathbf{U}_1],$$

and

$$\mathbf{V}_1^H \mathbf{A} \mathbf{V}_1 = \begin{bmatrix} \lambda_1 & & \\ 0 & & \\ \vdots & & \\ \vdots & & \\ 0 & & \end{bmatrix} \mathbf{V}_1^H \mathbf{A} \mathbf{U}_1.$$

Now suppose we have an  $N \times N$  matrix of the form

$$\mathbf{A}_p = \begin{bmatrix} \mathbf{\Delta}_p & \mathbf{W}_p \\ \mathbf{0} & \mathbf{M}_p \end{bmatrix}, \quad (1)$$

where  $\mathbf{\Delta}_p$  is a  $p \times p$  upper-triangular matrix,  $\mathbf{W}_p$  is an arbitrary  $p \times (N-p)$  matrix, and  $\mathbf{M}_p$  is an arbitrary  $(N-p) \times (N-p)$  square matrix. Now let  $\mathbf{v}_{p+1}$  be an eigenvector of  $\mathbf{M}_p$  with corresponding eigenvalue  $\lambda_{p+1}$ , and let  $\mathbf{U}_{p+1}$  be an  $(N-p) \times (N-p-1)$  matrix such that

$$\mathbf{Z}_{p+1} = [\mathbf{v}_{p+1} \quad \mathbf{U}_{p+1}]$$

is a  $(N-p) \times (N-p)$  orthonormal matrix. Set

$$\mathbf{V}_{p+1} = \begin{bmatrix} \mathbf{I}_p & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}_{p+1} \end{bmatrix},$$

where  $\mathbf{I}_p$  is the  $p \times p$  identity matrix. It should be clear that  $\mathbf{V}_{p+1}$  is an orthonormal matrix. Applying  $\mathbf{V}_{p+1}$  to the right of  $\mathbf{A}_p$  yields

$$\mathbf{A}_p \mathbf{V}_{p+1} = \begin{bmatrix} \mathbf{\Delta}_p & \mathbf{W}_p \mathbf{Z}_{p+1} \\ \mathbf{0} & [\lambda_{p+1} \mathbf{v}_{p+1} \quad \mathbf{M}_p \mathbf{U}_{p+1}] \end{bmatrix},$$

and so

$$\mathbf{V}_{p+1}^H \mathbf{A}_p \mathbf{V}_{p+1} = \begin{bmatrix} \mathbf{\Delta}_p & \mathbf{W}_p \mathbf{Z}_{p+1} \\ \mathbf{0} & \begin{bmatrix} \lambda_{p+1} & \\ 0 & \\ \vdots & \\ \vdots & \\ 0 & \end{bmatrix} \mathbf{Z}_{p+1}^H \mathbf{M}_p \mathbf{U}_{p+1} \end{bmatrix} = \begin{bmatrix} \mathbf{\Delta}_{p+1} & \mathbf{W}_{p+1} \\ \mathbf{0} & \mathbf{M}_{p+1} \end{bmatrix},$$

where  $\Delta_{p+1}$  is a  $(p+1) \times (p+1)$  upper-triangular matrix, and  $\mathbf{W}_{p+1}$  and  $\mathbf{M}_{p+1}$  are arbitrary  $(p+1) \times (N-p-1)$  and  $(N-p-1) \times (N-p-1)$  matrices, respectively.

Given an arbitrary  $\mathbf{A}$ ,

$$\mathbf{A}_p = \mathbf{V}_{p-1}^H \cdots \mathbf{V}_2^H \mathbf{V}_1^H \mathbf{A} \mathbf{V}_1 \mathbf{V}_2 \cdots \mathbf{V}_{p-1}$$

will have the form (1). Applying the construction over  $N$  iterations gives

$$\Delta = \mathbf{V}_N^H \cdots \mathbf{V}_2^H \mathbf{V}_1^H \mathbf{A} \mathbf{V}_1 \mathbf{V}_2 \cdots \mathbf{V}_N,$$

which will be upper-triangular. Since each of the  $\mathbf{V}_p$  are orthonormal,  $\mathbf{V} := \mathbf{V}_1 \mathbf{V}_2 \cdots \mathbf{V}_N$  will also be orthonormal. Thus

$$\Delta = \mathbf{V}^H \mathbf{A} \mathbf{V} \quad \Leftrightarrow \quad \mathbf{A} = \mathbf{V} \Delta \mathbf{V}^H,$$

where  $\Delta$  is upper-triangular and  $\mathbf{V}^H \mathbf{V} = \mathbf{I}$ .

## Eigenvalues of $\mathbf{A}$

The diagonal entries of the matrix  $\Delta$  will contain the  $\lambda_p$  used in the construction above (which we might recall are the eigenvalues of the submatrices  $\mathbf{M}_p$ ):

$$\Delta[p, p] = \lambda_p.$$

We can see now that the  $\lambda_p$  are also eigenvalues of  $\mathbf{A}$ . Since  $\Delta$  is triangular, its diagonal entries  $\lambda_1, \dots, \lambda_N$  are its eigenvalues. If  $\mathbf{x}_p$  is the eigenvector of  $\Delta$  corresponding to  $\lambda_p$ , then taking  $\mathbf{y}_p = \mathbf{V} \mathbf{x}_p$  we have

$$\mathbf{A} \mathbf{y}_p = \mathbf{V} \Delta \mathbf{V}^H \mathbf{V} \mathbf{x}_p = \mathbf{V} \Delta \mathbf{x}_p = \lambda_p \mathbf{V} \mathbf{x}_p = \lambda_p \mathbf{y}_p,$$

and so the  $\lambda_1, \dots, \lambda_N$  are eigenvalues of  $\mathbf{A}$  as well.

## Real-valued decompositions

If  $\mathbf{A}$  is real-valued but non-symmetric, then both  $\mathbf{V}$  and  $\mathbf{\Delta}$  can be complex-valued. However, there does exist real-valued  $\mathbf{U}$  and  $\mathbf{\Upsilon}$  such that

$$\mathbf{A} = \mathbf{U}\mathbf{\Upsilon}\mathbf{U}^T,$$

where  $\mathbf{U}$  is orthonormal,  $\mathbf{U}^T\mathbf{U} = \mathbf{I}$ , and  $\mathbf{\Upsilon}$  is almost upper-triangular:

$$\mathbf{\Upsilon} = \begin{bmatrix} \mathbf{\Lambda}_1 & * & \cdots & * \\ 0 & \mathbf{\Lambda}_2 & \cdots & * \\ \vdots & & \ddots & \vdots \\ 0 & \cdots & \cdots & \mathbf{\Lambda}_K \end{bmatrix}.$$

The  $\mathbf{\Lambda}_p$  above are either  $2 \times 2$  matrices or scalars; there is a  $2 \times 2$  block for every pair of complex-conjugate eigenvalues of  $\mathbf{A}$ , and a scalar for every real eigenvalue. Although this decomposition is not strictly upper-triangular, it carries many of the same advantages. For example, with  $\mathbf{U}$  pre-computed and given a  $\mathbf{b} \in \mathbb{R}^N$ , we can still compute the solution to  $\mathbf{A}\mathbf{x} = \mathbf{b}$  with  $O(N^2)$  operations.