

Linear Algebra Review

Akhil Vasvani

February 2019

1 Questions

Exercise 1. What is broadcasting in connection to Linear Algebra?

Proof. Typically for matrix addition, each matrix must be of the same dimension. However, in the context of deep learning, there is a shorthand notation for adding a matrix and a vector (the vector in this case is represented as a scalar). This vector (\mathbf{b} for instance) is added to each row of the matrix, which eliminates the need to define a matrix with \mathbf{b} copied into each row before doing the addition. This implicit copying of \mathbf{b} to many locations is called **broadcasting**.

Example.

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

In this notation, the \mathbf{b} is portrayed as a scalar 1, but in actuality it is a matrix of of 1's. This shorthand is broadcasting. \square

Exercise 2. What are scalars, vectors, matrices, and tensors?

Proof. A **scalar** is just a single number or a matrix with a single entry.

Example. Let $s \in \mathbb{R}$ be the slope of the line, 4, or $[3]$

A **vector** is a 1D array of numbers arranged in a particular order (a matrix with only one column). Another way to think of vectors is identifying points in space with each element giving the coordinate along a different axis.

Example. $\mathbf{x} = [x_1 \ x_2 \ x_3 \ x_4]^\top$

A matrix is a 2D array of numbers where each element is identified by two indices (ROW then COLUMN).

Example. $\mathbf{A} = \begin{bmatrix} 3 & 1 \\ 1 & 4 \end{bmatrix}$, where A has a height of two rows and a width of two columns. Hence, it is a shape of 2×2 . The element 3 is found at $A_{1,1}$ and element 4 is found at $A_{2,2}$.

A tensor is a multi-dimensional array with more than two axes denoted as \mathbf{A} where elements are at coordinates $A_{i,j,k}$. \square

Exercise 3. What is the Hadamard product of two matrices?

Proof. The Hadamard product (or element-wise product) is a matrix containing the product of the individual elements. It has a special symbol \odot . NOT TO BE CONFUSED WITH MATRIX MULTIPLICATION.

Example. $\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ and $\mathbf{B} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}$, so

$$\mathbf{AB} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} 1 \cdot (2) + 1 \cdot (2) & 1 \cdot (2) + 1 \cdot (2) \\ 1 \cdot (2) + 1 \cdot (2) & 1 \cdot (2) + 1 \cdot (2) \end{bmatrix} = \begin{bmatrix} 4 & 4 \\ 4 & 4 \end{bmatrix}, \text{ while}$$

$$\mathbf{A} \odot \mathbf{B} = \begin{bmatrix} 1 \cdot 2 & 1 \cdot 2 \\ 1 \cdot 2 & 1 \cdot 2 \end{bmatrix} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}$$

Note the difference? □

Exercise 4. What is an inverse matrix?

Proof. The inverse matrix of \mathbf{A} (denoted as \mathbf{A}^{-1}) is defined such that:

$$\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}_n,$$

where \mathbf{I}_n is the identity matrix. □

Exercise 5. If the inverse of a matrix exists, how to calculate it?

Proof. There are a couple ways to calculate the inverse matrix (if it exists).

Option 1: Create an augmented matrix with \mathbf{A} and the identity matrix \mathbf{I}_n (represented as $\mathbf{A}|\mathbf{I}_n$). The goal would be to turn \mathbf{A} into \mathbf{I}_n using elementary row operations. In turn, whatever we do to \mathbf{A} , we do the same to \mathbf{I}_n thereby turning \mathbf{I}_n into \mathbf{A}^{-1} .

Example. Let's say $\mathbf{A} = \begin{bmatrix} 2 & -1 & 0 \\ 1 & 3 & -1 \\ -3 & 0 & 1 \end{bmatrix}$, so

$$\mathbf{A}|\mathbf{I}_n = \left[\begin{array}{ccc|ccc} 2 & -1 & 0 & 1 & 0 & 0 \\ 1 & 3 & -1 & 0 & 1 & 0 \\ -3 & 0 & 1 & 0 & 0 & 1 \end{array} \right]$$

Row Operation: $R_2 + R_3 \rightarrow R_3$

$$\Rightarrow \left[\begin{array}{ccc|ccc} 2 & -1 & 0 & 1 & 0 & 0 \\ 1 & 3 & -1 & 0 & 1 & 0 \\ -2 & 3 & 0 & 0 & 1 & 1 \end{array} \right]$$

Row Operation: $R_3 + R_1 \rightarrow R_1$

$$\Rightarrow \left[\begin{array}{ccc|ccc} 0 & 2 & 0 & 1 & 1 & 1 \\ 1 & 3 & -1 & 0 & 1 & 0 \\ -2 & 3 & 0 & 0 & 1 & 1 \end{array} \right]$$

Row Operations: $\frac{3}{2}R_1 - R_2 \rightarrow R_2, \frac{3}{2}R_1 - R_3 \rightarrow R_3$

$$\Rightarrow \left[\begin{array}{ccc|ccc} 0 & 3 & 0 & \frac{3}{2} & \frac{3}{2} & \frac{3}{2} \\ -1 & 0 & 1 & \frac{3}{2} & \frac{3}{2} & \frac{3}{2} \\ 2 & 0 & 0 & \frac{3}{2} & \frac{3}{2} & \frac{3}{2} \end{array} \right]$$

Row Operations: $\frac{1}{3}R_1 \rightarrow R_1, \frac{1}{2}R_3 \rightarrow R_3$

$$\Rightarrow \left[\begin{array}{ccc|ccc} 0 & 1 & 0 & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ -1 & 0 & 1 & \frac{3}{4} & \frac{3}{4} & \frac{3}{4} \\ 1 & 0 & 0 & \frac{3}{4} & \frac{3}{4} & \frac{3}{4} \end{array} \right]$$

Row Operation: $R_3 + R_2 \rightarrow R_2$

$$\Rightarrow \left[\begin{array}{ccc|ccc} 0 & 1 & 0 & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 & \frac{9}{4} & \frac{3}{4} & \frac{7}{4} \\ 1 & 0 & 0 & \frac{3}{4} & \frac{3}{4} & \frac{3}{4} \end{array} \right]$$

Row Operation: $R_3 \leftrightarrow R_2$

$$\Rightarrow \left[\begin{array}{ccc|ccc} 0 & 1 & 0 & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 & \frac{3}{4} & \frac{3}{4} & \frac{3}{4} \\ 0 & 0 & 1 & \frac{9}{4} & \frac{3}{4} & \frac{7}{4} \end{array} \right]$$

Row Operation: $R_1 \leftrightarrow R_2$

$$\Rightarrow \left[\begin{array}{ccc|ccc} 1 & 0 & 0 & \frac{3}{4} & \frac{1}{4} & \frac{1}{4} \\ 0 & 1 & 0 & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 & \frac{9}{4} & \frac{3}{4} & \frac{7}{4} \end{array} \right]$$

Thus,

$$\mathbf{A}^{-1} = \begin{bmatrix} \frac{3}{4} & \frac{1}{4} & \frac{1}{4} \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{9}{4} & \frac{3}{4} & \frac{7}{4} \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 3 & 1 & 1 \\ 2 & 2 & 2 \\ 9 & 3 & 7 \end{bmatrix}$$

Option 2: The next way involves using Minors, Cofactors, and Adjugate. First, calculate the matrix of minors: for each element of the matrix ignore the values of the current row and column and calculate the determinant of the remaining values. Second, apply a "checkerboard" of minuses to the "Matrix of Minors". In other words, we need to change the sign of alternate cells (+, -, +, etc.). Next, transpose all elements of the previous matrix (swap their positions over the diagonal) (This matrix will be called the Adjugate). Finally, find the determinant of the original matrix and multiply 1/determinant to the Adjugate matrix.

Example. Let's use the same \mathbf{A} matrix in the previous example. Now let's calculate the determinant, which will be helpful for later before calculating the matrix of minors.

$$\begin{aligned} \det(A) &= \begin{vmatrix} 2 & -1 & 0 \\ 1 & 3 & -1 \\ -3 & 0 & 1 \end{vmatrix} = 2 \begin{vmatrix} 3 & -1 \\ 0 & 1 \end{vmatrix} - (-1) \begin{vmatrix} 1 & -1 \\ -3 & 1 \end{vmatrix} + 0 \begin{vmatrix} 1 & 3 \\ -3 & 0 \end{vmatrix} \\ &= 2((3 \times 1) - (-1 \times 0)) + (1)((1 \times 1) - (-1 \times -3)) + 0 \end{aligned}$$

$$\begin{aligned}
&= 2(3) + (1 - 3) \\
&= 6 - 2 \\
&= 4
\end{aligned}$$

Now, we know the determinant let's solve for the matrix of minors.

Matrix of Minors =

$$\begin{bmatrix} ((3 \times 1) - (-1 \times 0)) & ((1 \times 1) - (-1 \times -3)) & ((1 \times 0) - (3 \times -3)) \\ ((-1 \times 1) - (0 \times 0)) & ((2 \times 1) - (-3 \times 0)) & ((2 \times 0) - (-1 \times -3)) \\ ((-1 \times -1) - (0 \times 3)) & ((2 \times -1) - (0 \times -3)) & ((2 \times 3) - (-1 \times 1)) \end{bmatrix} = \begin{bmatrix} 3 & -2 & 9 \\ -1 & 2 & -3 \\ 1 & -2 & 7 \end{bmatrix}$$

Matrix of Cofactors =

$$\begin{bmatrix} + & - & + \\ - & + & - \\ + & - & + \end{bmatrix} \Rightarrow \begin{bmatrix} +(3) & -(-2) & +(9) \\ -(-1) & +(2) & -(-3) \\ +(1) & -(-2) & +(7) \end{bmatrix} = \begin{bmatrix} 3 & 2 & 9 \\ 1 & 2 & 3 \\ 1 & 2 & 7 \end{bmatrix}$$

Adjugate Matrix is the transpose of the matrix of cofactors:

$$\begin{bmatrix} 3 & 2 & 9 \\ 1 & 2 & 3 \\ 1 & 2 & 7 \end{bmatrix}^T = \begin{bmatrix} 3 & 1 & 1 \\ 2 & 2 & 2 \\ 9 & 3 & 7 \end{bmatrix}$$

Finally, to get \mathbf{A}^{-1} , we multiply $\frac{1}{\det(A)}$ to the Adjugate matrix:

$$\frac{1}{\det(A)} \begin{bmatrix} 3 & 1 & 1 \\ 2 & 2 & 2 \\ 9 & 3 & 7 \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 3 & 1 & 1 \\ 2 & 2 & 2 \\ 9 & 3 & 7 \end{bmatrix} = \mathbf{A}^{-1}$$

□

Exercise 6. What is the determinant of a square matrix? How is it calculated? What is the connection of determinant to eigenvalues?

Proof. The determinant of a square matrix \mathbf{A} , denoted as $\det(A)$, is a function mapping matrices to real scalars. The absolute value of the determinant can be thought of as a measure of how much multiplication by the matrix expands or contracts space. If the determinant is 0, then space is contracted completely along at least one dimension, causing it to lose all of its volume. If the determinant is 1, then the transformation preserves volume. The determinant is equal to the product of all the eigenvalues of the matrix.

Example. You can calculate determinants either via Row or Column. Just remember the signs from the matrix of cofactors!

Using the matrix $\mathbf{A} = \begin{bmatrix} 2 & -1 & 0 \\ 1 & 3 & -1 \\ -3 & 0 & 1 \end{bmatrix}$, as our example let's start by solving for the $\det(A)$ via the first row.

$$\begin{aligned}
\det(A) &= \begin{vmatrix} 2 & -1 & 0 \\ 1 & 3 & -1 \\ -3 & 0 & 1 \end{vmatrix} = 2 \begin{vmatrix} 3 & -1 \\ 0 & 1 \end{vmatrix} - (-1) \begin{vmatrix} 1 & -1 \\ -3 & 1 \end{vmatrix} + 0 \begin{vmatrix} 1 & 3 \\ -3 & 0 \end{vmatrix} \\
&= 2((3 \times 1) - (-1 \times 0)) + (1)((1 \times 1) - (-1 \times -3)) + 0 \\
&= 2(3) + (1 - 3) \\
&= 6 - 2 \\
&= 4
\end{aligned}$$

Now, let's solve for $\det(A)$ via the second column:

$$\begin{aligned}
\det(A) &= \begin{vmatrix} 2 & -1 & 0 \\ 1 & 3 & -1 \\ -3 & 0 & 1 \end{vmatrix} = -(-1) \begin{vmatrix} 1 & -1 \\ -3 & 1 \end{vmatrix} + (3) \begin{vmatrix} 2 & 0 \\ -3 & 1 \end{vmatrix} - 0 \\
&= ((1 \times 1) - (-1 \times -3)) + 3((2 \times 1) - (0 \times -3)) \\
&= (1 - 3) + 3(2) \\
&= 6 - 2 \\
&= 4
\end{aligned}$$

□

Exercise 7. Discuss span and linear dependence.

Proof. Using $\mathbf{Ax} = \mathbf{b}$ as a reference, we can look at the columns of \mathbf{A} . Think of the columns of \mathbf{A} as specifying different directions we can travel from the origin (point specified to be 0 of all vectors), and determine how many ways there are to reach \mathbf{b} . Therefore, we can think of \mathbf{x} specifying how far we should travel in each of these directions like this:

$$\mathbf{Ax} = \sum_i x_i \mathbf{A}_{:,i}.$$

This is a linear combination in which each vector is multiplied by a scalar coefficient and adding the results. The span of a set of vectors is the set of all points obtainable by linear combination of the original vectors. (sounds sort of like a range?)

Example. Now consider $\mathbf{A} = \begin{bmatrix} 2 & 2 \\ 1 & 1 \end{bmatrix}$. Both columns $(\begin{bmatrix} 2 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 \\ 1 \end{bmatrix})$ are identical, so this matrix

has the same column space or more formally put: $\text{Span}(\mathbf{A}) = \text{Span}(\begin{bmatrix} 2 & 2 \\ 1 & 1 \end{bmatrix}) = \text{Span}(\begin{bmatrix} 2 \\ 1 \end{bmatrix})$.

The second column is a replica of the first column, so the span of \mathbf{A} is 2×1 . This type of redundancy is called linear dependence. A set of vectors is linearly independent if no vector in the set is a linear combination of the other vectors. Such as:

$$\mathbf{M} = \begin{bmatrix} 2 & 4 \\ 1 & 3 \end{bmatrix}$$

□

Exercise 8. What is $\mathbf{Ax} = \mathbf{b}$? When does $\mathbf{Ax} = \mathbf{b}$ has a unique solution?

Proof. $\mathbf{Ax} = \mathbf{b}$ is a system of linear equations where $\mathbf{A} \in \mathbb{R}^{m \times n}$ is a known matrix, $\mathbf{b} \in \mathbb{R}^m$ is a known vector, and $\mathbf{x} \in \mathbb{R}^n$ is a vector of unknown variables we would like to solve for. Observe matrix $\mathbf{A}^{m \times n}$. $\mathbf{Ax} = \mathbf{b}$ will have a unique solution if each column is linearly independent and \mathbf{b} lies in the column space of \mathbf{A} .

Example. Say $\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & -1 & 1 \\ 3 & 0 & -1 \end{bmatrix}$ and $\mathbf{b} = \begin{bmatrix} 9 \\ 8 \\ 3 \end{bmatrix}$, then the augmented matrix $\mathbf{A}|\mathbf{b}$ in row-reduced echelon form is:

$$\left[\begin{array}{ccc|c} 1 & 0 & 0 & 2 \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & 3 \end{array} \right]$$

Each of the columns is linearly independent and \mathbf{b} lies in the column space, so $\mathbf{Ax} = \mathbf{b}$ has a unique solution in which $\mathbf{x} = \begin{bmatrix} 2 \\ -1 \\ 3 \end{bmatrix}$.

□

Exercise 9. In $\mathbf{Ax} = \mathbf{b}$, what happens when \mathbf{A} is fat or tall?

Proof. Now depending on the **rank**—the number of first nonzero entry of each row (**pivots**) in a reduced row-echelon form matrix—affects the number of solutions for either a tall or fat matrix.

Note. A matrix is in reduced row-echelon form if (1) it is in row-echelon form, (2) all of the pivots are equal to 1, and (3) all entries in the pivot columns, except for the pivots themselves, are equal to zero.

$$A_{rref} = \begin{bmatrix} 1 & 0 & 4 & 0 \\ 0 & 1 & 2 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

Here the $\text{rank}[A_{rref}] = 3$.

Now the general rule of thumb is:

- 1) $\mathbf{Ax} = \mathbf{b}$ is inconsistent (i.e. no solution exists) if and only if $\text{rank}[\mathbf{A}] < \text{rank}[\mathbf{A}|\mathbf{b}]$.
- 2) $\mathbf{Ax} = \mathbf{b}$ has a unique solution if and only if $\text{rank}[\mathbf{A}] = \text{rank}[\mathbf{A}|\mathbf{b}] = n$.
- 3) $\mathbf{Ax} = \mathbf{b}$ has infinitely many solutions if and only if $\text{rank}[\mathbf{A}] = \text{rank}[\mathbf{A}|\mathbf{b}] < n$.

Since $\mathbf{A} \in \mathbb{R}^{m \times n}$, if \mathbf{A} is a tall matrix, then the numbers of rows $m >$ the number of columns n . If $\text{rank}(\mathbf{A}) < \text{rank}[\mathbf{A}|\mathbf{b}]$, then there is no solution—no linear combination will reach the desired \mathbf{b} —for all \mathbf{b} . If $\text{rank}(\mathbf{A}) = n$, then there is ONLY one unique solution for every \mathbf{b} . In laymen's terms, if $\text{rank}(\mathbf{A}) = n$, then there is one solution that is unique; otherwise, there does not exist a solution. Note if $\mathbf{b} = 0$, there there is unique (very obvious) solution if $\text{rank}(\mathbf{A}) = n$ —hint: it is 0—but if $\text{rank}(\mathbf{A}) < n$, there are infinitely many solutions.

TL;DR. For most tall matrices, there are usually no solutions.

However, if \mathbf{A} is a fat matrix then the number of columns $n >$ the number of rows m . If $\text{rank}[\mathbf{A}] < m$, then either there does not exist a solution or there are infinitely many solutions. However, if $\text{rank}[\mathbf{A}] = m$, there exists at least one solution for all \mathbf{b} —which suggests that there can be more than one solution (thereby infinitely many solutions). Note, if $\mathbf{b} = \mathbf{0}$ then there are infinite solutions.

TL;DR. For most fat matrices, there are usually infinitely many solutions that exist—there are an infinite number of linear combinations to reach the desired \mathbf{b} .

Example. Let's look for a solution to $\mathbf{Ax} = \mathbf{b}$, where

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 & 0 \\ 2 & -1 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \mathbf{x} = \begin{bmatrix} w \\ x \\ y \\ z \end{bmatrix}, \mathbf{b} = \begin{bmatrix} 9 \\ 8 \\ b_3 \end{bmatrix} \Rightarrow \begin{bmatrix} 1 & 2 & 3 & 0 \\ 2 & -1 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} w \\ x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 9 \\ 8 \\ b_3 \end{bmatrix}.$$

The augmented matrix $\mathbf{A}|\mathbf{b}$ in row-reduced echelon form:

$$\left[\begin{array}{cccc|c} 1 & 0 & -1 & 0 & -5 \\ 0 & 1 & 1 & 0 & 2 \\ 0 & 0 & 0 & 0 & b_3 \end{array} \right] \Rightarrow \begin{bmatrix} w \\ x \\ y \\ z \end{bmatrix} = \begin{bmatrix} y - 5 \\ 2 - y \\ y \\ z \end{bmatrix} = \begin{bmatrix} -5 + y \\ 2 - y \\ y \\ z \end{bmatrix} = \begin{bmatrix} -5 \\ 2 \\ 0 \\ 0 \end{bmatrix} + y \begin{bmatrix} 1 \\ -1 \\ 1 \\ 0 \end{bmatrix} + z \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}.$$

Both y, z can be any \mathbb{R} . Now, let's pay attention to the constant b_3 . If $b_3 = 0$, then $\text{rank}[\mathbf{A}] = \text{rank}[\mathbf{A}|\mathbf{b}]$, which $< n$ —it is also $< m$. So, there are an infinite number of linear combinations to reach the desired \mathbf{b} . On the other hand, if $b_3 \neq 0$, then $\text{rank}[\mathbf{A}] < \text{rank}[\mathbf{A}|\mathbf{b}]$, so there does not exist a solution. Therefore, matrix \mathbf{A} is a fat matrix.

$$\text{Note, if } \mathbf{A} = \begin{bmatrix} 1 & 2 & 3 & 0 \\ 2 & -1 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \Rightarrow \left[\begin{array}{cccc|c} 1 & 0 & -1 & 0 & -5 \\ 0 & 1 & 1 & 0 & 2 \\ 0 & 0 & 0 & 1 & b_3 \end{array} \right] \Rightarrow z = b_3.$$

So $b_3 \in \mathbb{R} \rightarrow z \in \mathbb{R}$, which means that any value of z can satisfy this equation. Hence, there are an infinite number of solutions.

Example. Let's look another example and find a solution to $\mathbf{Ax} = \mathbf{b}$, where

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & -1 & 1 \\ 3 & 0 & -1 \\ 0 & 0 & 0 \end{bmatrix}, \mathbf{x} = \begin{bmatrix} x \\ y \\ z \end{bmatrix}, \mathbf{b} = \begin{bmatrix} 9 \\ 8 \\ 3 \\ b_4 \end{bmatrix}.$$

The augmented matrix $\mathbf{A}|\mathbf{b}$ in row-reduced echelon form:

$$\left[\begin{array}{ccc|c} 1 & 0 & 0 & 2 \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & 3 \\ 0 & 0 & 0 & b_4 \end{array} \right] \Rightarrow \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

If the values for $b_4 = 0$ then $\text{rank}[\mathbf{A}] = \text{rank}[\mathbf{A}|\mathbf{b}]$, which $= n$. Thus, there is a unique solution. However, $b_4 \neq 0$, then $\text{rank}[\mathbf{A}] < \text{rank}[\mathbf{A}|\mathbf{b}]$, which means there does not exist a solution. Therefore, at most this system has one solution. Thus, \mathbf{A} is a tall matrix.

□

Exercise 10. When does an inverse of \mathbf{A} exist?

Proof. Inverse of \mathbf{A} (or \mathbf{A}^{-1}) exists if in $\mathbf{Ax} = \mathbf{b}$, there has to be n linearly independent columns and at most one solution for every value of \mathbf{b} . This will make the matrix a square matrix (number of m rows = number of n columns) and nonsingular (columns are linearly independent). Also, another check is that the determinant cannot be equal to 0. □

Exercise 11. What is a norm? What is L^1 , L^2 and L^∞ norm?

Proof. A norm is a function that measures the size of vectors. Simply speaking, the norm of a \mathbf{x} measures the distance from the origin to the point x .

L^2 norm (pronounced L-two) is known as the Euclidean Distance: $\|\mathbf{x}\|_2 = \sqrt{\sum_i |\mathbf{x}_i|^2}$.

It is so commonly used in machine learning that the 2 subscript is dropped out, so it's just $\|\mathbf{x}\|$. Remember, it's the Euclidean distance from the origin to the point x .

L^1 norm (pronounced L-one) is commonly referred to as Manhattan distance:

$$\|\mathbf{x}\|_1 = \sum_i |\mathbf{x}_i|.$$

L^∞ norm (also known as max norm) is the absolute value of the element with the largest magnitude in the vector: $\|\mathbf{x}\|_\infty = \max_i |x_i|$.

Example. Say $\mathbf{x} = \begin{bmatrix} -6 \\ 3 \\ 4 \end{bmatrix}$ then

L^2 norm =

$$\|\mathbf{x}\|_2 = \sqrt{(-6)^2 + (3)^2 + (4)^2} = \sqrt{36 + 9 + 16} = \sqrt{61} \approx 7.810.$$

L^1 norm =

$$\|\mathbf{x}\|_1 = |-6| + |3| + |4| = 13.$$

L^∞ norm =

$$\max_i |\mathbf{x}| = \max_{i=3} |x_3| = 6$$

□

Exercise 12. What are the conditions a norm has to satisfy?

Proof. More formally put, a norm is any function which satisfies the following properties:

- 1) $f(\mathbf{x}) = 0 \Rightarrow \mathbf{x} = \mathbf{0}$
- 2) $f(\mathbf{x}+\mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$ (the triangle inequality)
- 3) $\forall \alpha \in \mathbb{R}, f(\alpha\mathbf{x}) = |\alpha|f(\mathbf{x})$ □

Exercise 13. Why is squared of L^2 norm preferred in Machine Learning than just L^2 norm?

Proof. The squared L^2 norm is mathematically and computationally convenient to work with than the L^2 norm itself. For example, the derivatives of the squared L^2 norm with respect to each element of \mathbf{x} depend only on the corresponding element of \mathbf{x} , while all the derivatives of the L^2 norm may depend on the entire vector.

Example. From our definition,

$$L^2 \text{ norm} = \|\mathbf{x}\|_2 = \sqrt{\sum_i |\mathbf{x}_i|^2}, \text{ so squared of } L^2 \text{ norm} = \|\mathbf{x}\|_2^2 = \sum_i |\mathbf{x}_i|^2.$$

Now, if we take the derivative of the squared L^2 norm:

$$\frac{\partial \|\mathbf{x}\|_2^2}{\partial x_j} = \frac{\partial}{\partial x_j} \left[\sum_i |\mathbf{x}_i|^2 \right] = \sum_i \frac{\partial}{\partial x_j} |\mathbf{x}_i|^2 = 0 + \dots + 0 + 2|x_j| + 0 + \dots + 0 = 2|x_j|$$

when $i = j$.

Now, if we take the derivative of the L^2 norm:

$$\frac{\partial \|\mathbf{x}\|_2}{\partial x_j} = \frac{\partial}{\partial x_j} \left[\sqrt{\sum_i |\mathbf{x}_i|^2} \right] = \frac{1}{2} \left(\sum_i |\mathbf{x}_i|^2 \right)^{-\frac{1}{2}} \cdot \sum_i \frac{\partial}{\partial x_j} |\mathbf{x}_i|^2 = \frac{|x_j|}{\sqrt{\sum_i |\mathbf{x}_i|^2}}.$$

Also when $i = j$ for the numerator.

Therefore, it is easier to rely on the squared L^2 because you do not need the whole vector to compute its gradient. You just need the specific element. □

Exercise 14. When L^1 norm is preferred over L^2 norm?

Proof. However, in many contexts, the squared L^2 norm may be undesirable because it increases very slowly near the origin. And in several machine learning applications, it is important to discriminate between elements that are exactly 0 and elements that are small but nonzero. Hence, we turn to the L^1 norm which grows at the same rate in all locations. Every time an element of x moves away from 0 by ϵ , the L^1 norm increases by ϵ . □

Exercise 15. Can the number of nonzero elements in a vector be defined as L^0 norm? If not, why?

Proof. No, this is not correct. The number of nonzero entries in a vector is not a norm because scaling the vector by α (property 3) does not change the number of nonzero entries. The L^1 norm is often used as a substitute for the number of nonzero entries. □

Exercise 16. What is Frobenius norm?

Proof. The Frobenius norm measures the size of a matrix:

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i,j} \mathbf{A}_{i,j}^2},$$

which is analogous to L^2 norm of a vector.

Example. Say $\mathbf{A} = \begin{bmatrix} 2 & -1 & 5 \\ 0 & 2 & 1 \\ 3 & 1 & 1 \end{bmatrix}$, so $\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^3 \sum_{j=1}^3 \mathbf{A}_{i,j}^2}$

$$\begin{aligned}
&= \sqrt{(2)^2 + (-1)^2 + (5)^2 + (0)^2 + (2)^2 + (1)^2 + (3)^2 + (1)^2 + (1)^2} \\
&= \sqrt{4 + 1 + 25 + 4 + 1 + 9 + 1 + 1} \\
&= \sqrt{2 \cdot 4 + 4 \cdot 1 + 25 + 9} \\
&= \sqrt{3 \cdot 4 + 34} \\
&= \sqrt{46}
\end{aligned}$$

□

Exercise 17. What is a diagonal matrix?

Proof. Diagonal matrices consist mostly of zeros and have nonzero entries only along the main diagonal. Formally, a matrix \mathbf{D} is diagonal if and only if $D_{i,j} = 0$ for all $i \neq j$.

Notation: $\text{diag}(\mathbf{v})$ denotes a square diagonal matrix whose diagonal entries are given by the entries of the vector \mathbf{v} .

Example. Say $\mathbf{v} = \begin{bmatrix} 6 \\ 7 \\ 19 \end{bmatrix}$, so $\text{diag}(\mathbf{v}) = \begin{bmatrix} 6 & 0 & 0 \\ 0 & 7 & 0 \\ 0 & 0 & 19 \end{bmatrix}$. The identity matrix (I) is also a diagonal matrix.

□

Exercise 18. Why is multiplication by diagonal matrix computationally cheap? How is the multiplication different for square vs. non-square diagonal matrix?

Proof. Multiplication by diagonal matrix is computationally cheap because to compute $\text{diag}(\mathbf{v})\mathbf{x}$, we only need to scale each element x_i by v_i . Simply put: $\text{diag}(\mathbf{v})\mathbf{x} = \mathbf{v} \odot \mathbf{x}$.

Example. $\text{diag}(\mathbf{v}) = \begin{bmatrix} 6 & 0 & 0 \\ 0 & 7 & 0 \\ 0 & 0 & 19 \end{bmatrix}$ and $\mathbf{x} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$, so

$$\text{diag}(\mathbf{v})\mathbf{x} = \begin{bmatrix} 6 \cdot 1 & 0 \cdot 2 & 0 \cdot 3 \\ 0 \cdot 4 & 7 \cdot 5 & 0 \cdot 6 \\ 0 \cdot 7 & 0 \cdot 8 & 19 \cdot 9 \end{bmatrix} = \begin{bmatrix} 6 & 0 & 0 \\ 0 & 35 & 0 \\ 0 & 0 & 171 \end{bmatrix}$$

However, not diagonal matrices need be square. It is possible to construct a rectangular diagonal matrix. Non-square diagonal matrices do not have inverses but it is still possible to multiply by them cheaply. For a non-square diagonal matrix \mathbf{D} , the product $\mathbf{D}\mathbf{x}$ will involve scaling each element of \mathbf{x} , and either concatenating some zeros to the result if \mathbf{D} is taller than it is fat, or discarding some of the last elements of the vector if \mathbf{D} is fatter than it is tall.

Example. D is taller than it is fat. $D = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & -3 \\ 0 & 0 & 0 \end{bmatrix}$ and $\mathbf{x} = \begin{bmatrix} 3 & 2 & 1 \\ 6 & 5 & 4 \\ 9 & 8 & 7 \end{bmatrix}$, so

$$D\mathbf{x} = \begin{bmatrix} 1 \cdot 3 & 1 \cdot 2 & 1 \cdot 1 \\ 4 \cdot 6 & 4 \cdot 5 & 4 \cdot 4 \\ -3 \cdot 9 & -3 \cdot 8 & -3 \cdot 7 \\ 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 3 & 2 & 1 \\ 24 & 20 & 16 \\ -27 & -24 & -21 \\ 0 & 0 & 0 \end{bmatrix}.$$

Now, let's concatenate another column of 0's to make this a square matrix:

$$D\mathbf{x} + \mathbf{0} = \begin{bmatrix} 3 & 2 & 1 & 0 \\ 24 & 20 & 16 & 0 \\ -27 & -24 & -21 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Example. D is fat. $D = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 \\ 0 & 0 & -3 & 0 \end{bmatrix}$ and $\mathbf{x} = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 8 & 7 & 6 \\ 5 & 4 & 3 & 2 \end{bmatrix}$, so

$$D\mathbf{x} = \begin{bmatrix} 1 \cdot 1 & 1 \cdot 2 & 1 \cdot 3 & 1 \cdot 4 \\ 4 \cdot 5 & 4 \cdot 6 & 4 \cdot 7 & 4 \cdot 8 \\ -3 \cdot 9 & -3 \cdot 8 & -3 \cdot 7 & -3 \cdot 6 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 20 & 24 & 28 & 32 \\ -27 & -24 & -21 & -18 \end{bmatrix}.$$

Now, let's remove the last column to turn this into a square matrix:

$$\rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 20 & 24 & 28 \\ -27 & -24 & -21 \end{bmatrix}$$

□

Exercise 19. At what conditions does the inverse of a diagonal matrix exist?

Proof. The inverse of a diagonal matrix exists only if the matrix is square and every diagonal entry is nonzero, and in that case: $\text{diag}(\mathbf{v})^{-1} = \text{diag}\left(\left[\frac{1}{\mathbf{v}_1}, \dots, \frac{1}{\mathbf{v}_n}\right]\right)^\top$

Example. Say $\text{diag}(\mathbf{v}) = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 7 \end{bmatrix}$, so $\text{diag}(\mathbf{v})^{-1} = \begin{bmatrix} \frac{1}{3} & 0 & 0 \\ 0 & \frac{1}{5} & 0 \\ 0 & 0 & \frac{1}{7} \end{bmatrix}$.

To prove that $\text{diag}(\mathbf{v})^{-1}$ is in fact the inverse: $\text{diag}(\mathbf{v})^{-1}\text{diag}(\mathbf{v}) =$

$$\begin{bmatrix} 3 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 7 \end{bmatrix} \begin{bmatrix} \frac{1}{3} & 0 & 0 \\ 0 & \frac{1}{5} & 0 \\ 0 & 0 & \frac{1}{7} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \mathbf{I}_3$$

Note, that the diagonal matrices are symmetric.

□

Exercise 20. What is a symmetric matrix?

Proof. A symmetric matrix is any matrix that is equal to its own transpose: $\mathbf{A} = \mathbf{A}^\top$.

Example. $\mathbf{A} = \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix}$, so $\mathbf{A}^\top = \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix}$. Thus, $\mathbf{A} = \mathbf{A}^\top$.

□

Exercise 21. What is a unit vector?

Proof. A unit vector is a vector with unit norm: $\|\mathbf{x}\|_2 = 1$

□

Exercise 22. When are two vectors x and y orthogonal?

Proof. A vector \mathbf{x} and a vector \mathbf{y} are orthogonal to each other if $\mathbf{x}^\top \mathbf{y} = 0$. If both vectors have nonzero norm, this means that they are at a 90° angle to each other.

Note. You can represent the dot product of two vectors in terms of norms, which helps prove when two vectors are orthogonal:

$$\mathbf{x}^\top \mathbf{y} = \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \cos \theta.$$

Hence, this equation holds true when $\theta = \frac{\pi}{2}$ or 90° , which means that vectors \mathbf{x}, \mathbf{y} must be perpendicular to each other to be orthogonal.

Note. The dot product of two vectors is also called the **inner product**, which is denoted as:

$$\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^\top \mathbf{y} = \sum_{i=1}^m x_i y_i. \tag{1}$$

Note the output is a scalar not a matrix.

□

Exercise 23. At \mathbb{R}^n , what is the maximum possible number of orthogonal vectors with non-zero norm?

Proof. In \mathbb{R}^n , at most n vectors may be mutually orthogonal with nonzero norm.

□

Exercise 24. When are two vectors \mathbf{x} and \mathbf{y} orthonormal?

Proof. If two vectors are not only orthogonal but also have unit norm, then they are orthonormal.

□

Exercise 25. What is an orthogonal matrix? Why is computationally preferred?

Proof. An orthogonal matrix is a square matrix whose rows are mutually orthonormal and whose columns are mutually orthonormal: $\mathbf{A}^\top \mathbf{A} = \mathbf{A} \mathbf{A}^\top = \mathbf{I} \rightarrow \mathbf{A}^{-1} = \mathbf{A}^\top$.

Orthogonal matrices are of interest because their inverse is very cheap to compute.

Example. Orthogonal matrix $\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$

□

Exercise 26. What is eigendecomposition, eigenvectors and eigenvalues?

Proof. Eigendecomposition decomposes a matrix into a set of eigenvectors and eigenvalues (kind of like how we can break down a number into its prime factors).

An eigenvector of a square matrix \mathbf{A} is a nonzero vector \mathbf{v} such that multiplication by \mathbf{A} alters only on the scale of \mathbf{v} : $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ (the scalar *lambda* is known as the eigenvalue corresponding to the eigenvector). □

Exercise 27. How to find eigenvalues of a matrix?

Proof. To find the eigenvalues of a matrix, we find the values of λ which satisfy the characteristic equation of the matrix \mathbf{A} : $\det(\mathbf{A} - I\lambda) = 0$. Then, plug the solved lambda values in to the matrix \mathbf{A} to get the augmented matrix. Then row-reduce to find the eigenvectors.

Note: matrix \mathbf{A} must be a square matrix.

Example. Let's find the eigenvalues and then the eigenvectors of matrix $\mathbf{A} = \begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix}$.

$$\begin{aligned} \det(\mathbf{A} - I_2\lambda) = 0 &\rightarrow \det\left(\begin{bmatrix} -\lambda & 1 \\ -2 & -3 - \lambda \end{bmatrix}\right) = 0 \\ &\rightarrow \begin{vmatrix} -\lambda & 1 \\ -2 & -3 - \lambda \end{vmatrix} = 0 \\ (-\lambda)(-3 - \lambda) - (-2)(1) &= 0 \\ (\lambda)(3 + \lambda) + 2 &= 0 \\ 3\lambda + \lambda^2 + 2 &= 0 \\ (\lambda + 1)(\lambda + 2) &= 0 \\ \lambda &= -1, -2 \end{aligned}$$

These are our eigenvalues. Now, let's find their corresponding eigenvectors.

Using the equation $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$, isolate \mathbf{v} (our eigenvector): $(\mathbf{A} - I_2\lambda)\mathbf{v} = 0$. Plug in the first λ value of -1 into the equation $\rightarrow (\mathbf{A} - I_2(-1))\mathbf{v} = (\mathbf{A} + I_2)\mathbf{v} = 0$

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

Row-reducing this matrix:

Row Operation: $2R_1 + R_2 \rightarrow R_2$

$$\Rightarrow \left[\begin{array}{cc|c} 1 & 1 & 0 \\ 0 & 0 & 0 \end{array} \right].$$

Now,

$$\begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \mathbf{v} = 0$$

$$\begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0 \Rightarrow x + y = 0 \Rightarrow x = -y,$$

so $\mathbf{v}_1 = \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -y \\ y \end{bmatrix} = y \begin{bmatrix} -1 \\ 1 \end{bmatrix} \Rightarrow k \begin{bmatrix} -1 \\ 1 \end{bmatrix}$. To normalize vector, set $k = \frac{1}{\sqrt{(-1)^2 + (1)^2}} = \frac{1}{\sqrt{2}}$.

Plug in the second λ value of -2 into the equation $\rightarrow (\mathbf{A} - I_2(-2))\mathbf{v} = (\mathbf{A} + 2I_2)\mathbf{v} = 0$

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

Row-reducing this matrix:

Row Operation: $R_1 + R_2 \rightarrow R_2$

$$\Rightarrow \left[\begin{array}{cc|c} 2 & 1 & 0 \\ 0 & 0 & 0 \end{array} \right].$$

Now,

$$\begin{bmatrix} 2 & 1 \\ 0 & 0 \end{bmatrix} \mathbf{v} = 0$$

$$\begin{bmatrix} 2 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0 \Rightarrow 2x + y = 0 \Rightarrow x = -\frac{y}{2},$$

so $\mathbf{v}_2 = \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -\frac{y}{2} \\ y \end{bmatrix} = y \begin{bmatrix} -\frac{1}{2} \\ 1 \end{bmatrix} \Rightarrow k \begin{bmatrix} -\frac{1}{2} \\ 1 \end{bmatrix}$. To normalize vector, set $k = \frac{1}{\sqrt{\left(-\frac{1}{2}\right)^2 + (1)^2}} = \frac{2}{\sqrt{5}}$

□

Exercise 28. Write the eigendecomposition formula for a matrix. If the matrix is real symmetric, how will this change?

Proof. The eigendecomposition of \mathbf{A} is given by:

$$\mathbf{A} = \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1}, \quad (2)$$

where matrix \mathbf{V} is all the eigenvectors with one eigenvector per column: $\mathbf{V} = [\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}]$ and likewise concatenating all the eigenvalues to form a vector $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_n]^T$

Example. Using our previous example from above, find the eigendecomposition of $\mathbf{A} = \begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix}$.

Now, as calculated from above, both the eigenvalues ($\lambda = -1, -2$) and the eigenvectors ($\mathbf{v}_1 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \mathbf{v}_2 = \begin{bmatrix} -\frac{1}{2} \\ 1 \end{bmatrix}$) are known. So to create matrix \mathbf{V} , stack the eigenvector corresponding to the largest eigenvalue first followed by the second.

$$\mathbf{V} = \begin{bmatrix} -1 & -\frac{1}{2} \\ 1 & 1 \end{bmatrix}.$$

So,

$$\det(\mathbf{V}) = \begin{vmatrix} -1 & -\frac{1}{2} \\ 1 & 1 \end{vmatrix} = -1(1) - (-\frac{1}{2})(1) = -1 + \frac{1}{2} = -\frac{1}{2}.$$

Because the determinant is nonzero, \mathbf{V}^{-1} does exist. $\mathbf{V}^{-1} = \frac{1}{\det(\mathbf{V})} \begin{bmatrix} 1 & \frac{1}{2} \\ -1 & -1 \end{bmatrix} = \begin{bmatrix} -2 & -1 \\ 2 & 2 \end{bmatrix}$

and $\text{diag}(\boldsymbol{\lambda}) = \begin{bmatrix} -1 & 0 \\ 0 & -2 \end{bmatrix}$.

Plugging this all into the equation:

$$\begin{aligned} \mathbf{A} &= \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1} = \begin{bmatrix} -1 & -\frac{1}{2} \\ 1 & 1 \end{bmatrix} \begin{bmatrix} -1 & 0 \\ 0 & -2 \end{bmatrix} \begin{bmatrix} -2 & -1 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} -1(-1) + 0 & 0 + (-2)(-\frac{1}{2}) \\ -1(1) + 0 & 0 - 2 \end{bmatrix} \begin{bmatrix} -2 & -1 \\ 2 & 2 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 1 \\ -1 & -2 \end{bmatrix} \begin{bmatrix} -2 & -1 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} -2(1) + 2(1) & -1(1) + 2(1) \\ -2(-1) + (-2)(2) & (-1)(-1) + 2(-2) \end{bmatrix} \\ &= \begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix}. \end{aligned}$$

Note: there are infinite eigenvectors corresponding with a specific eigenvalues. While this may change the determinant, the important point here is that the ratio of \mathbf{v}_1 and \mathbf{v}_2 remains the same.

For real symmetric matrices,

$$\mathbf{A} = \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^T, \quad (3)$$

where \mathbf{Q} is an orthogonal matrix composed of eigenvectors of \mathbf{A} , and $\boldsymbol{\Lambda}$ is a diagonal matrix. The eigenvalue $\Lambda_{i,j}$ is associated with the eigenvector of column i of $\mathbf{Q}_{:,i}$. Because \mathbf{Q} is an orthogonal matrix think of \mathbf{A} as scaling space by λ_i in the direction $\mathbf{v}^{(i)}$.

Example. Matrix $\mathbf{A} = \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix}$ is real symmetric. So, let's find the eigendecomposition.

First, solve the characteristic equation: $\det(\mathbf{A} - I_2\lambda) = 0$.

$$\rightarrow \det\left(\begin{bmatrix} 3-\lambda & 2 \\ 2 & 3-\lambda \end{bmatrix}\right) = 0 \rightarrow \begin{vmatrix} 3-\lambda & 2 \\ 2 & 3-\lambda \end{vmatrix} = 0$$

Difference of Squares: $\rightarrow (3 - \lambda)^2 - (2)^2 = 0$

$$(3 - \lambda + 2)(3 - \lambda - 2) = 0$$

$$(5 - \lambda)(1 - \lambda) = 0$$

$$\lambda = 5, 1$$

These are our eigenvalues. Let's find their corresponding eigenvectors.

$\lambda = 5 : \mathbf{A} - I_2(5)$

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

Row-reducing this matrix:

Row Operation: $R_1 + R_2 \rightarrow R_2$

$$\Rightarrow \left[\begin{array}{cc|c} -2 & 2 & 0 \\ 0 & 0 & 0 \end{array} \right]$$

Row Operation: $-\frac{1}{2}R_1 \rightarrow R_1$

$$\Rightarrow \left[\begin{array}{cc|c} 1 & -1 & 0 \\ 0 & 0 & 0 \end{array} \right]$$

Now, remember the equation $(\mathbf{A} - \mathbf{I}_2\lambda)\mathbf{v} = 0$? Time to use it again:

$$\begin{bmatrix} 1 & -1 \\ 0 & 0 \end{bmatrix} \mathbf{v} = 0 \rightarrow \begin{bmatrix} 1 & -1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0 \rightarrow x = y.$$

Thus, $\mathbf{v}_1 = \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} y \\ y \end{bmatrix} = y \begin{bmatrix} 1 \\ 1 \end{bmatrix} = k \begin{bmatrix} 1 \\ 1 \end{bmatrix}$. To normalize vector, let $k = \frac{1}{\sqrt{2}}$
 $\lambda = 1 : \mathbf{A} - \mathbf{I}_2(1)$

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

Row-reducing this matrix:

Row Operation: $R_1 - R_2 \rightarrow R_2$

$$\Rightarrow \left[\begin{array}{cc|c} 2 & 2 & 0 \\ 0 & 0 & 0 \end{array} \right]$$

Row Operation: $\frac{1}{2}R_1 \rightarrow R_1$

$$\Rightarrow \left[\begin{array}{cc|c} 1 & 1 & 0 \\ 0 & 0 & 0 \end{array} \right]$$

$$\begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \mathbf{v} = 0 \rightarrow \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0 \rightarrow x = -y.$$

Therefore, $\mathbf{v}_2 = \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -y \\ y \end{bmatrix} = y \begin{bmatrix} -1 \\ 1 \end{bmatrix} = k \begin{bmatrix} -1 \\ 1 \end{bmatrix}$. To normalize vector, let $k = \frac{1}{\sqrt{2}}$.

Now, stacking these eigenvectors together (from highest eigenvalue to lowest),

$$\mathbf{Q} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}$$

$$\mathbf{Q}^\top = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$$

So,

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$$

$$\begin{aligned} \Rightarrow \mathbf{A} &= \left(\frac{1}{\sqrt{2}}\right)^2 \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 5 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 5 & -1 \\ 5 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 6 & 4 \\ 4 & 6 \end{bmatrix} \\ &= \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix} \end{aligned}$$

□

Exercise 29. Is the Eigendecomposition guaranteed to be unique? If not, then how do we represent it?

Proof. While any real symmetric matrix \mathbf{A} is guaranteed to have an eigendecomposition, the eigendecomposition may not be unique. If any two or more eigenvectors share the same eigenvalue, then any set of orthogonal vectors lying in their span are also eigenvectors with that eigenvalue, and we could equivalently choose a \mathbf{Q} using those eigenvectors instead. By convention, we usually sort the entries of $\mathbf{\Lambda}$ in descending order. Under this convention, the eigendecomposition is unique only if all of the eigenvalues are unique. □

Exercise 30. What are positive definite, negative definite, positive semi definite and negative semi definite matrices?

Proof. A matrix whose eigenvalues are all positive is called positive definite. A matrix whose eigenvalues are all positive or zero-valued is called positive semidefinite. Likewise, if all eigenvalues are negative, the matrix is negative definite, and if all eigenvalues are negative or zero-valued, it is negative semidefinite.

Positive semidefinite matrices are interesting because they guarantee that $\forall \mathbf{x}, \mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$.

Positive definite matrices additionally guarantee that $\mathbf{x}^\top \mathbf{A} \mathbf{x} = 0, \rightarrow \mathbf{x} = 0$.

Note. So why do we care?

Imagine, there is a vector \mathbf{z} , which will have a certain direction. When we multiply matrix \mathbf{M} with \mathbf{z} , \mathbf{z} no longer points in the same direction. Therefore, the direction of \mathbf{z} is transformed by \mathbf{M} . Would it not be nice in an abstract sense to be able to multiply some matrices multiple times and they will not change the sign of the vectors? If you multiply positive numbers to other positive numbers, it does not change its sign. So let's apply the same logic to Linear Algebra using the eigenvalues.

Example. Say, we have a matrix \mathbf{A} and eigenvalues λ . Now, each eigenvector of \mathbf{A} can be represented by vector \mathbf{x} :

$$\mathbf{A} \mathbf{x} = \lambda \mathbf{x} \rightarrow \mathbf{x}^\top \mathbf{A} \mathbf{x} = \mathbf{x}^\top \lambda \mathbf{x}.$$

Now using norms instead of dot products,

$$\mathbf{x}^\top \lambda \mathbf{x} = \|\mathbf{x}\|_2 \|\lambda \mathbf{x}\|_2 \cos \theta = \lambda \|\mathbf{x}\|_2^2 \cos(0) = \lambda \|\mathbf{x}\|_2^2.$$

And,

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \|\mathbf{x}\|_2 \|\mathbf{A} \mathbf{x}\|_2 \cos \theta$$

Plugging this all back in, we have:

$$\Rightarrow \|\mathbf{x}\|_2 \|\mathbf{A} \mathbf{x}\|_2 \cos \theta = \lambda \|\mathbf{x}\|_2^2$$

$$\Rightarrow \lambda = (\|\mathbf{x}\|_2^2)^{-1} \|\mathbf{x}\|_2 \|\mathbf{Ax}\|_2 \cos \theta.$$

On an intuitive note, if \mathbf{A} is a positive definite matrix, the new direction will always point in the “same general” direction (here “same general” means $\theta < \frac{\pi}{2}$). \mathbf{A} is “semi” definite, if $\theta \leq \frac{\pi}{2}$. However, it will reverse the original direction if $\theta \geq \frac{\pi}{2}$.

Now how does this relate to eigenvalues? Well, if we assume $\|\mathbf{x}\|_2 \|\mathbf{Ax}\|_2 \cos \theta > 0$ and $\|\mathbf{x}\|_2^2 > 0$, then eigenvalues (λ) must be greater than 0! So, a positive definite matrix must have positive eigenvalues.

If $\mathbf{A} = \begin{bmatrix} 3 & 2 \\ 2 & 3 \end{bmatrix}$, then our eigenvalues (λ) = 5,1. So, matrix \mathbf{A} is a positive definite matrix.

If $\mathbf{A} = \begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix}$, then since our eigenvalues (λ) = -1,-2. So, therefore, matrix \mathbf{A} is a negative definite matrix.

Note. Positive definite matrices are extremely useful tools. For example, if the Hessian of a function is a positive definite matrix then the function is convex. Otherwise if the Hessian is a negative definite matrix, then the function is non-convex.

□

Exercise 31. What is Singular Value Decomposition? Why do we use it? Why not just use ED?

Proof. The **singular value decomposition** (SVD) provides another way to factorize a matrix, into singular vectors and singular values. Like eigendecomposition (ED), the SVD allows one to discover some of the same kind of information. However, the SVD is generally more applicable. Every real matrix has a singular value decomposition, but the same is not true of the eigenvalue decomposition. For example, if a matrix is not square, the eigendecomposition is not defined, and we must use a singular value decomposition instead. □

Exercise 32. Given a matrix \mathbf{A} , how will you calculate its Singular Value Decomposition?

Proof. Suppose that \mathbf{A} is an $m \times n$ matrix. Then \mathbf{U} is defined to be an $m \times m$ matrix, \mathbf{D} to be an $m \times n$ matrix, and \mathbf{V} to be an $n \times n$ matrix, then the Singular Value Decomposition (SVD):

$$\mathbf{A} = \mathbf{UDV}^\top \tag{4}$$

Note, matrices \mathbf{U} and \mathbf{V} are both defined to be orthogonal matrices and matrix \mathbf{D} is defined to be a diagonal matrix (note it does not have to be a square).

Example. Say matrix $\mathbf{A} = \begin{bmatrix} 2 & 2 \\ 1 & 1 \end{bmatrix}$, so $\mathbf{A}^\top = \begin{bmatrix} 2 & 1 \\ 2 & 1 \end{bmatrix}$.

$$\mathbf{A}^\top \mathbf{A} = \begin{bmatrix} 2 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 2 & 2 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 5 \\ 5 & 5 \end{bmatrix}, \quad \mathbf{AA}^\top = \begin{bmatrix} 2 & 2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 2 & 1 \end{bmatrix} = \begin{bmatrix} 8 & 4 \\ 4 & 2 \end{bmatrix}.$$

So,

$$\det(\mathbf{A}^\top \mathbf{A} - \mathbf{I}_2 \lambda) = \begin{vmatrix} 5 - \lambda & 5 \\ 5 & 5 - \lambda \end{vmatrix} = 0 \rightarrow (5 - \lambda)^2 - 25 = 0 \rightarrow \lambda^2 - 10\lambda = 0 \\ \rightarrow \lambda = 0, 10$$

$$\det(\mathbf{A}\mathbf{A}^\top - \mathbf{I}_2\lambda) = \begin{vmatrix} 8-\lambda & 4 \\ 4 & 2-\lambda \end{vmatrix} = 0 \rightarrow (8-\lambda)(2-\lambda) - 16 = 0 \rightarrow \lambda^2 - 10\lambda = 0$$

$$\rightarrow \lambda = 0, 10$$

These are our eigenvalues. Let's find their corresponding eigenvectors.

$$\lambda = 10: \mathbf{A}\mathbf{A}^\top - \mathbf{I}_2(10)$$

$$\rightarrow \left[\begin{array}{cc|c} 8-10 & 4 & 0 \\ 4 & 2-10 & 0 \end{array} \right] = \left[\begin{array}{cc|c} -2 & 4 & 0 \\ 4 & -8 & 0 \end{array} \right]$$

$$\text{Row Operation: } 2R_1 + R_2 \rightarrow R_2$$

$$\rightarrow \left[\begin{array}{cc|c} -2 & 4 & 0 \\ 0 & 0 & 0 \end{array} \right]$$

$$\text{Row Operation: } -\frac{1}{2}R_1 \rightarrow R_1$$

$$\rightarrow \left[\begin{array}{cc|c} 1 & -2 & 0 \\ 0 & 0 & 0 \end{array} \right] \Rightarrow \begin{bmatrix} 1 & -2 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \Rightarrow x = 2y.$$

So, $\mathbf{v}_1 = \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 2y \\ y \end{bmatrix} = y \begin{bmatrix} 2 \\ 1 \end{bmatrix} = k \begin{bmatrix} 2 \\ 1 \end{bmatrix}$. To normalize vector, set $k = \frac{1}{\sqrt{5}}$.

$$\lambda = 0: \mathbf{A}\mathbf{A}^\top - \mathbf{I}_2(0)$$

$$\rightarrow \left[\begin{array}{cc|c} 8-0 & 4 & 0 \\ 4 & 2-0 & 0 \end{array} \right] = \left[\begin{array}{cc|c} 8 & 4 & 0 \\ 4 & 2 & 0 \end{array} \right]$$

$$\text{Row Operation: } \frac{1}{2}R_1 - R_2 \rightarrow R_2$$

$$\rightarrow \left[\begin{array}{cc|c} 8 & 4 & 0 \\ 0 & 0 & 0 \end{array} \right]$$

$$\text{Row Operation: } \frac{1}{4}R_1 \rightarrow R_1$$

$$\rightarrow \left[\begin{array}{cc|c} 2 & 1 & 0 \\ 0 & 0 & 0 \end{array} \right] \Rightarrow \begin{bmatrix} 2 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \Rightarrow x = -\frac{y}{2}.$$

So, $\mathbf{v}_2 = \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -\frac{y}{2} \\ y \end{bmatrix} = y \begin{bmatrix} -\frac{1}{2} \\ 1 \end{bmatrix} = k \begin{bmatrix} -\frac{1}{2} \\ 1 \end{bmatrix}$. To normalize, set $k = \frac{2}{\sqrt{5}}$.

$$\lambda = 10: \mathbf{A}^\top \mathbf{A} - \mathbf{I}_2(10)$$

$$\rightarrow \left[\begin{array}{cc|c} 5-10 & 5 & 0 \\ 5 & 5-10 & 0 \end{array} \right] = \left[\begin{array}{cc|c} -5 & 5 & 0 \\ 5 & -5 & 0 \end{array} \right]$$

$$\text{Row Operation: } R_1 + R_2 \rightarrow R_2$$

$$\rightarrow \left[\begin{array}{cc|c} -5 & 5 & 0 \\ 0 & 0 & 0 \end{array} \right]$$

$$\text{Row Operation: } -\frac{1}{5}R_1 \rightarrow R_1$$

$$\rightarrow \left[\begin{array}{cc|c} 1 & -1 & 0 \\ 0 & 0 & 0 \end{array} \right] \Rightarrow \begin{bmatrix} 1 & -1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \Rightarrow x = y.$$

So, $\mathbf{v}_1 = \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} y \\ y \end{bmatrix} = y \begin{bmatrix} 1 \\ 1 \end{bmatrix} = k \begin{bmatrix} 1 \\ 1 \end{bmatrix}$. To normalize, set $k = \frac{1}{\sqrt{2}}$.

$$\lambda = 0: \mathbf{A}^\top \mathbf{A} - \mathbf{I}_2(0)$$

$$\rightarrow \left[\begin{array}{cc|c} 5-0 & 5 & 0 \\ 5 & 5-0 & 0 \end{array} \right] = \left[\begin{array}{cc|c} 5 & 5 & 0 \\ 5 & 5 & 0 \end{array} \right]$$

Row Operation: $R_1 - R_2 \rightarrow R_2$

$$\rightarrow \left[\begin{array}{cc|c} 5 & 5 & 0 \\ 0 & 0 & 0 \end{array} \right]$$

Row Operation: $\frac{1}{5}R_1 \rightarrow R_1$

$$\rightarrow \left[\begin{array}{cc|c} 1 & 1 & 0 \\ 0 & 0 & 0 \end{array} \right] \Rightarrow \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \Rightarrow x = -y.$$

So, $\mathbf{v}_2 = \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -y \\ y \end{bmatrix} = y \begin{bmatrix} -1 \\ 1 \end{bmatrix} = k \begin{bmatrix} -1 \\ 1 \end{bmatrix}$. To normalize, set $k = \frac{1}{\sqrt{2}}$.

$$\text{Now, since } \mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top \rightarrow \mathbf{A}^\top = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top)^\top = \mathbf{V}\mathbf{\Sigma}^\top\mathbf{U}^\top$$

$$\Rightarrow \mathbf{A}^\top \mathbf{A} = \mathbf{V}\mathbf{\Sigma}^\top\mathbf{U}^\top\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top = \mathbf{V}\mathbf{\Sigma}^\top\mathbf{\Sigma}\mathbf{V}^\top = \mathbf{V}\mathbf{\Sigma}^2\mathbf{V}^\top.$$

$$\Rightarrow \mathbf{A}\mathbf{A}^\top = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top\mathbf{V}\mathbf{\Sigma}^\top\mathbf{U}^\top = \mathbf{U}\mathbf{\Sigma}^\top\mathbf{\Sigma}\mathbf{U}^\top = \mathbf{U}\mathbf{\Sigma}^2\mathbf{U}^\top.$$

Now, since we calculated 2 pairs of eigenvalues for $\mathbf{A}\mathbf{A}^\top$ and $\mathbf{A}^\top\mathbf{A}$, we know which respective matrices we were building. $\mathbf{\Sigma}\mathbf{\Sigma}^\top$ or $\mathbf{\Sigma}^\top\mathbf{\Sigma}$ is just the diagonal matrix of eigenvalues for $\mathbf{A}^\top\mathbf{A}$ and $\mathbf{A}\mathbf{A}^\top$

$$\Rightarrow \left[\begin{array}{cc} 10 & 0 \\ 0 & 0 \end{array} \right] \Rightarrow \mathbf{\Sigma} = \begin{bmatrix} \sqrt{10} & 0 \\ 0 & \sqrt{0} \end{bmatrix}$$

$$\text{Therefore, } \mathbf{U} = \frac{1}{\sqrt{5}} \begin{bmatrix} 2 & -1 \\ 1 & 2 \end{bmatrix} \text{ and } \mathbf{V} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \rightarrow \mathbf{V}^\top = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}.$$

$$\Rightarrow \frac{1}{\sqrt{10}} \begin{bmatrix} 2 & -1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} \sqrt{10} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} = \mathbf{A}$$

Note. Now, what does this show intuitively? More importantly, what does this mean? Think of the \mathbf{A} as a type of linear transformation. Therefore, SVD breaks down the transformation into three simple steps:

1) The initial rotation (or reflection) \mathbf{V}^\top .

2) The $\mathbf{\Sigma}$ matrix has to do with rescaling vertically or horizontally (really in any dimension depending on the matrix).

3) The final rotation (or reflection) \mathbf{U} .

By looking at the SVD, we can tell which matrix has a significant impact on transforming the matrix and how much as well.

If we apply that intuition to this problem, we explain break up the transformation:

$$\mathbf{V}^\top = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} = \frac{1}{\sqrt{2}} \left(\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \right) = \frac{1}{\sqrt{2}}(\mathbf{I}_2 + \mathbf{R}(\theta = 270^\circ)).$$

So \mathbf{V}^\top boils down to summing the initial vector with the rotation of the initial vector by 270° .

$$\Sigma = \begin{bmatrix} \sqrt{10} & 0 \\ 0 & 0 \end{bmatrix} \rightarrow \Sigma \text{ basically scales the } x \text{ coordinate by } \sqrt{10}.$$

$$\mathbf{U} = \frac{1}{\sqrt{5}} \begin{bmatrix} 2 & 1 \\ -1 & 2 \end{bmatrix} = \frac{1}{\sqrt{5}} \left(\begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} + \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \right) = \frac{1}{\sqrt{5}}(2\mathbf{I}_2 + \mathbf{R}(\theta = 90^\circ)).$$

\mathbf{U} boils down to summing two times the initial vector with the rotation of the initial vector by 90° .

□

Exercise 33. What are singular values, left singulars, and right singulars?

Proof. The elements along the diagonal of \mathbf{D} are known as the **singular values** of the matrix \mathbf{A} . The columns of \mathbf{U} are known as the **left-singular vectors**. The columns of \mathbf{V} are known as as the **right-singular vectors**. □

Exercise 34. What is the connection of Singular Value Decomposition of \mathbf{A} with functions of \mathbf{A} ?

Proof. We can actually interpret the singular value decomposition of \mathbf{A} in terms of the eigendecomposition of functions of \mathbf{A} . The left-singular vectors of \mathbf{A} are the eigenvectors of $\mathbf{A}\mathbf{A}^\top$. The right-singular vectors of \mathbf{A} are the eigenvectors of $\mathbf{A}^\top\mathbf{A}$. The non-zero singular values of \mathbf{A} are the square roots of the eigenvalues of $\mathbf{A}^\top\mathbf{A}$. The same is true for $\mathbf{A}\mathbf{A}^\top$. □

Exercise 35. Why are singular values always non-negative?

Proof. Let's assume matrix \mathbf{A} has real entries—otherwise consider $\mathbf{A}^H\mathbf{A}$ —then $\mathbf{A}^\top\mathbf{A}$ is positive semidefinite because it is an inner product. One of the properties of an inner product is that for any vector \mathbf{u} , $\langle \mathbf{u}, \mathbf{u} \rangle \geq 0$ or $\langle \mathbf{u}, \mathbf{u} \rangle = 0$ if and only if $\mathbf{u} = \mathbf{0}$. Now, since $\mathbf{A}^\top\mathbf{A}$ is matrix (and we want a vector), if \mathbf{v} is a nonzero vector in \mathbb{R}^n , then we have

$$\begin{aligned} \langle \mathbf{A}^\top\mathbf{A}\mathbf{v}, \mathbf{v} \rangle &:= (\mathbf{A}^\top\mathbf{A}\mathbf{v})^\top\mathbf{v} = \mathbf{v}^\top\mathbf{A}^\top\mathbf{A}\mathbf{v} \\ &= \langle \mathbf{A}\mathbf{v}, \mathbf{A}\mathbf{v} \rangle \geq 0, \forall \mathbf{v}. \end{aligned}$$

Therefore, the eigenvalues of $\mathbf{A}^\top\mathbf{A}$ are non-negative, and the singular values of \mathbf{A} are the square roots of the eigenvalues. □

Exercise 36. What is the Moore Penrose pseudo inverse and how to calculate it?

Proof. The Moore Penrose pseudo inverse definition is:

$$\mathbf{A}^+ = \lim_{\alpha \rightarrow 0} (\mathbf{A}^\top\mathbf{A} + \alpha\mathbf{I})^{-1}\mathbf{A}^\top. \quad (5)$$

Practical algorithms for computing the pseudo inverse are not based on this definition, but rather the formula:

$$\mathbf{A}^+ = \mathbf{V}\mathbf{D}^+\mathbf{U}^\top \quad (6)$$

where \mathbf{U} , \mathbf{D} , and \mathbf{V} are the singular value decomposition of \mathbf{A} . The pseudo inverse \mathbf{D}^+ of a diagonal matrix \mathbf{D} is obtained by taking the reciprocal of its non-zero elements then taking the transpose of the resulting matrix.

Example. Say $\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ \sqrt{3} & 0 \end{bmatrix}$, so $\rightarrow \mathbf{A}^\top = \begin{bmatrix} 1 & 0 & \sqrt{3} \\ 0 & 1 & 0 \end{bmatrix}$.

$$\mathbf{A}^\top\mathbf{A} = \begin{bmatrix} 1 & 0 & \sqrt{3} \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ \sqrt{3} & 0 \end{bmatrix} = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{A}\mathbf{A}^\top = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ \sqrt{3} & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & \sqrt{3} \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & \sqrt{3} \\ 0 & 1 & 0 \\ \sqrt{3} & 0 & 3 \end{bmatrix}.$$

So,

$$\det(\mathbf{A}^\top\mathbf{A} - \mathbf{I}_2\lambda) = \begin{vmatrix} 4-\lambda & 0 \\ 0 & 1-\lambda \end{vmatrix} = 0 \rightarrow (4-\lambda)(1-\lambda) - 0 = 0 \rightarrow \lambda = 1, 4$$

$$\det(\mathbf{A}\mathbf{A}^\top - \mathbf{I}_3\lambda) = \begin{vmatrix} 1-\lambda & 0 & \sqrt{3} \\ 0 & 1-\lambda & 0 \\ \sqrt{3} & 0 & 3-\lambda \end{vmatrix} = 0 \rightarrow (1-\lambda)[(1-\lambda)(3-\lambda) - 0] - 0 + \sqrt{3}[0 - \sqrt{3}(1-\lambda)] = 0$$

$$(1-\lambda)[(1-\lambda)(3-\lambda)] - 3(1-\lambda) = 0$$

$$(1-\lambda)[(1-\lambda)(3-\lambda) - 3] = 0$$

$$(1-\lambda)[3 - \lambda - 3\lambda + \lambda^2 - 3] = 0$$

$$(1-\lambda)[-4\lambda + \lambda^2] = 0$$

$$(1-\lambda)(-\lambda)(4\lambda) = 0$$

$$\rightarrow \lambda = 0, 1, 4$$

These are our eigenvalues. Let's find their corresponding eigenvectors.

$\lambda = 4$: $\mathbf{A}\mathbf{A}^\top - \mathbf{I}_3(4)$

$$\rightarrow \left[\begin{array}{ccc|c} (1-4) & 0 & \sqrt{3} & 0 \\ 0 & (1-4) & 0 & 0 \\ \sqrt{3} & 0 & (3-4) & 0 \end{array} \right] = \left[\begin{array}{ccc|c} -3 & 0 & \sqrt{3} & 0 \\ 0 & -3 & 0 & 0 \\ \sqrt{3} & 0 & -1 & 0 \end{array} \right]$$

Row Operation: $-\sqrt{3}R_3 - R_1 \rightarrow R_1, -\frac{1}{3}R_2 \rightarrow R_2$

$$\Rightarrow \left[\begin{array}{ccc|c} 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ \sqrt{3} & 0 & -1 & 0 \end{array} \right]$$

Row Operation: $R_3 \leftrightarrow R_1$

$$\Rightarrow \left[\begin{array}{ccc|c} \sqrt{3} & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right] \Rightarrow \left[\begin{array}{ccc|c} \sqrt{3} & 0 & -1 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{array} \right] \begin{bmatrix} x \\ y \\ z \end{bmatrix} = 0$$

$$\Rightarrow x\sqrt{3} - z = 0 \Rightarrow z = x\sqrt{3}$$

$$\Rightarrow y = 0.$$

Thus, $\mathbf{v}_1 = \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} x \\ 0 \\ x\sqrt{3} \end{bmatrix} = x \begin{bmatrix} 1 \\ 0 \\ \sqrt{3} \end{bmatrix} = k \begin{bmatrix} 1 \\ 0 \\ \sqrt{3} \end{bmatrix}$. To normalize, let $k = \frac{1}{2}$.

$\lambda = 1$: $\mathbf{A}\mathbf{A}^\top - \mathbf{I}_3(1)$

$$\rightarrow \left[\begin{array}{ccc|c} (1-1) & 0 & \sqrt{3} & 0 \\ 0 & (1-1) & 0 & 0 \\ \sqrt{3} & 0 & (3-1) & 0 \end{array} \right] = \left[\begin{array}{ccc|c} 0 & 0 & \sqrt{3} & 0 \\ 0 & 0 & 0 & 0 \\ \sqrt{3} & 0 & 2 & 0 \end{array} \right]$$

Row Operation: $\frac{1}{\sqrt{3}}R_1 \rightarrow R_1$

$$\Rightarrow \left[\begin{array}{ccc|c} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ \sqrt{3} & 0 & 2 & 0 \end{array} \right]$$

Row Operation: $R_3 \leftrightarrow R_1$

$$\Rightarrow \left[\begin{array}{ccc|c} \sqrt{3} & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{array} \right]$$

Row Operation: $2R_3 - R_1$

$$\Rightarrow \left[\begin{array}{ccc|c} \sqrt{3} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{array} \right] \Rightarrow \left[\begin{array}{ccc|c} \sqrt{3} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{array} \right] \begin{bmatrix} x \\ y \\ z \end{bmatrix} = 0$$

$$\Rightarrow x\sqrt{3} = 0 \Rightarrow x = 0$$

$$\Rightarrow z = 0.$$

Thus, $\mathbf{v}_2 = \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 0 \\ y \\ 0 \end{bmatrix} = y \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = k \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$.

$\lambda = 0$: $\mathbf{A}\mathbf{A}^\top$

$$\rightarrow \left[\begin{array}{ccc|c} (1-0) & 0 & \sqrt{3} & 0 \\ 0 & (1-0) & 0 & 0 \\ \sqrt{3} & 0 & (3-0) & 0 \end{array} \right] = \left[\begin{array}{ccc|c} 1 & 0 & \sqrt{3} & 0 \\ 0 & 1 & 0 & 0 \\ \sqrt{3} & 0 & 3 & 0 \end{array} \right]$$

Row Operation: $\sqrt{3}R_1 - R_3 \rightarrow R_3$

$$\Rightarrow \left[\begin{array}{ccc|c} 1 & 0 & \sqrt{3} & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right] \Rightarrow \begin{bmatrix} 1 & 0 & \sqrt{3} \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = 0$$

$$\Rightarrow x + z\sqrt{3} = 0 \Rightarrow x = -z\sqrt{3}$$

$$\Rightarrow y = 0.$$

Thus, $\mathbf{v}_3 = \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} -z\sqrt{3} \\ 0 \\ z \end{bmatrix} = z \begin{bmatrix} -\sqrt{3} \\ 0 \\ 1 \end{bmatrix} = k \begin{bmatrix} -\sqrt{3} \\ 0 \\ 1 \end{bmatrix}$. To normalize, $k = \frac{1}{2}$.

$\lambda = 4$: $\mathbf{A}^\top \mathbf{A} - \mathbf{I}_2(4)$

$$\rightarrow \left[\begin{array}{cc|c} (4-4) & 0 & 0 \\ 0 & (1-4) & 0 \end{array} \right] = \left[\begin{array}{cc|c} 0 & 0 & 0 \\ 0 & -3 & 0 \end{array} \right] \Rightarrow \begin{bmatrix} 0 & 0 \\ 0 & -3 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0 \rightarrow y = 0.$$

Thus, $\mathbf{v}_1 = \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x \\ 0 \end{bmatrix} = x \begin{bmatrix} 1 \\ 0 \end{bmatrix} = k \begin{bmatrix} 1 \\ 0 \end{bmatrix}$.

$\lambda = 1$: $\mathbf{A}^\top \mathbf{A} - \mathbf{I}_2(1)$

$$\rightarrow \left[\begin{array}{cc|c} (4-1) & 0 & 0 \\ 0 & (1-1) & 0 \end{array} \right] = \left[\begin{array}{cc|c} 3 & 0 & 0 \\ 0 & 0 & 0 \end{array} \right] \Rightarrow \begin{bmatrix} 3 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0 \rightarrow x = 0.$$

Thus, $\mathbf{v}_2 = \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ y \end{bmatrix} = y \begin{bmatrix} 0 \\ 1 \end{bmatrix} = k \begin{bmatrix} 0 \\ 1 \end{bmatrix}$.

Now, we calculated 2 eigenvalues for $\mathbf{A}^\top \mathbf{A}$ and 3 eigenvalues for $\mathbf{A}\mathbf{A}^\top$. However, our Σ is not a square matrix, so it cannot be symmetric. Note, that since this is pseudoinverse problem, we will now denote \mathbf{D} to be Σ . Because we have an extra eigenvalue (0), we will add to it to the subsequent row.

Therefore our \mathbf{D} matrix is: $\begin{bmatrix} \sqrt{4} & 0 \\ 0 & \sqrt{1} \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \rightarrow \mathbf{D}^+ = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{1} \\ 0 & 0 \end{bmatrix}^\top = \begin{bmatrix} \frac{1}{2} & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$.

Therefore, $\mathbf{U} = \begin{bmatrix} \frac{1}{2} & 0 & -\frac{\sqrt{3}}{2} \\ 0 & 1 & 0 \\ \frac{\sqrt{3}}{2} & 0 & \frac{1}{2} \end{bmatrix} \rightarrow \mathbf{U}^\top = \begin{bmatrix} \frac{1}{2} & 0 & \frac{\sqrt{3}}{2} \\ 0 & 1 & 0 \\ -\frac{\sqrt{3}}{2} & 0 & \frac{1}{2} \end{bmatrix}$ and $\mathbf{V} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \mathbf{V}^\top$.

Thus, $\mathbf{A}^\top = \mathbf{V}\mathbf{D}^+\mathbf{U}^\top$

$$\Rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 0 & \frac{\sqrt{3}}{2} \\ 0 & 1 & 0 \\ -\frac{\sqrt{3}}{2} & 0 & \frac{1}{2} \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 1 & 0 & \sqrt{3} \\ 0 & 4 & 0 \end{bmatrix}.$$

□

Exercise 37. If we do Moore Penrose pseudo inverse on $\mathbf{A}\mathbf{x} = \mathbf{b}$, what solution is provided is \mathbf{A} is fat? Moreover, what solution is provided if \mathbf{A} is tall?

Proof. When \mathbf{A} has more columns than rows (fat matrix), then solving a linear equation using the pseudoinverse provides one of the many possible solutions. Specifically, it provides the solution $\mathbf{x} = \mathbf{A}^+\mathbf{b}$ with minimal Euclidean norm $\|\mathbf{x}\|_2$ among all possible solutions.

When \mathbf{A} has more rows than columns (tall matrix), it is possible for there to be no solution. In this case, using the pseudoinverse gives us the \mathbf{x} for which \mathbf{Ax} is as close as possible to \mathbf{b} in terms of Euclidean norm $\|\mathbf{Ax} - \mathbf{b}\|_2$. \square

Exercise 38. Which matrices can be decomposed by ED?

Proof. Only square matrices have eigenvalue decomposition. \square

Exercise 39. Which matrices can be decomposed by SVD?

Proof. Every real matrix has a singular value decomposition square. However, matrix \mathbf{V} can either be conjugate transpose or normal transpose depending on whether matrix \mathbf{A} is complex or real. \square

Exercise 40. What is the trace of a matrix?

$$\text{Tr}(\mathbf{A}) = \sum_i \mathbf{A}_{i,i} \quad (7)$$

Exercise 41. How to write Frobenius norm of a matrix \mathbf{A} in terms of trace?

$$\|\mathbf{A}\|_F = \sqrt{\text{Tr}(\mathbf{AA}^\top)} \quad (8)$$

Exercise 42. Why is trace of a multiplication of matrices invariant to cyclic permutations?

Proof. Because the trace of a matrix is only concerned with the diagonal entries, regardless of the cyclic permutations (even if the two matrices are different in shape), the sum of the diagonals will stay the same. \square

Exercise 43. What is the trace of a scalar?

$$a = \text{Tr}(a) \quad (9)$$

2 Extra Stuff

Exercise 44. What is a Hermitian Matrix?

Proof. A **Hermitian matrix** is a complex square matrix that is equal to its own conjugate transpose—more specifically the element of the i -th row and the j -th column is equal to the complex conjugate of the element in the j -th row and the i -column.

$$\mathbf{A} = \mathbf{A}^H = \overline{\mathbf{A}^\top}, \quad (10)$$

where \mathbf{A}^H denotes the **conjugate transpose**, which is done by first taking the transpose of \mathbf{A} and then taking the complex conjugate of each entry in \mathbf{A}^\top (denoted as $\overline{\mathbf{A}^\top}$).

Example. Say $\mathbf{A} = \begin{bmatrix} 2 & 2+i & 4 \\ 2-i & 3 & i \\ 4 & -i & 1 \end{bmatrix} \rightarrow \mathbf{A}^\top = \begin{bmatrix} 2 & 2-i & 4 \\ 2+i & 3 & -i \\ 4 & i & 1 \end{bmatrix} \rightarrow \overline{\mathbf{A}^\top} = \begin{bmatrix} 2 & 2+i & 4 \\ 2-i & 3 & i \\ 4 & -i & 1 \end{bmatrix}.$

$$\overline{\mathbf{A}^\top} = \mathbf{A}^H = \mathbf{A}.$$

Note. Hermitian matrices can be thought of as a complex extension of real symmetric matrices because they always have real eigenvalues. Sometimes in SVD, \mathbf{U}, \mathbf{V} matrices can be complex— so $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H \rightarrow \mathbf{A}^H = \mathbf{V}\mathbf{\Sigma}\mathbf{U}^H.$

More properties of Hermitian matrices can be found [here](#). □

Exercise 45. What is an outer product?

Proof. The outer product $\mathbf{u} \otimes \mathbf{v}$ is equivalent to a matrix multiplication $\mathbf{u}\mathbf{v}^\top$, provided that \mathbf{u} is represented as a $m \times 1$ column vector and \mathbf{v} as a $n \times 1$ column vector (which makes \mathbf{v}^\top a row vector). For instance, if $m = 4$ and $n = 3$, then

$$\mathbf{u} \otimes \mathbf{v} = \mathbf{u}\mathbf{v}^\top = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{bmatrix} \begin{bmatrix} v_1 & v_2 & v_3 \end{bmatrix} = \begin{bmatrix} u_1v_1 & u_1v_2 & u_1v_3 \\ u_2v_1 & u_2v_2 & u_2v_3 \\ u_3v_1 & u_3v_2 & u_3v_3 \\ u_4v_1 & u_4v_2 & u_4v_3 \end{bmatrix}. \quad (11)$$

Note. This returns a matrix, not a scalar. □

Exercise 46. What is the Rotation Matrix?

Proof.

$$\mathbf{R}(\theta) = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \quad (12)$$

□

Exercise 47. What is a projection? What is a projection matrix? Why do we care?

Proof. If we have a vector \mathbf{b} and a line determined by a vector \mathbf{a} , how do we find the point on the line that is closest to \mathbf{b} ?

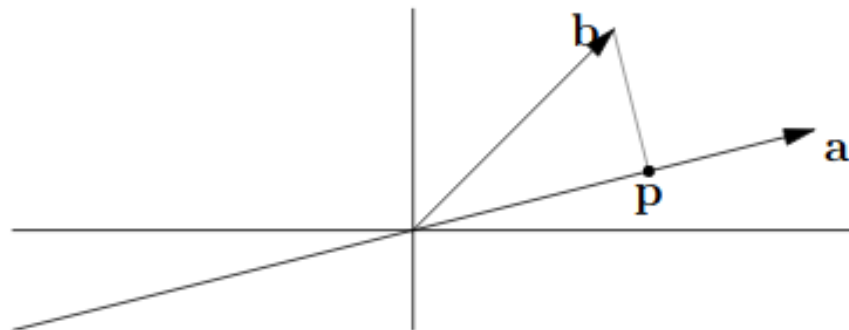


Figure 1: The point closest to \mathbf{b} on the line determined by \mathbf{a} .

As we see from Figure 1, the closest point \mathbf{p} is at the intersection formed by a line through \mathbf{b} that is orthogonal to \mathbf{a} . If we think of \mathbf{p} as an approximation of \mathbf{b} , then the length of $\mathbf{e} = \mathbf{b} - \mathbf{p}$ is the error in that approximation. We could try to find \mathbf{p} using trigonometry or calculus, but it's easier to use linear algebra. Since \mathbf{p} lies on the line through \mathbf{a} , we know $\mathbf{p} = \mathbf{a}x$ for some number x . We also know that \mathbf{a} is perpendicular to $\mathbf{e} = \mathbf{b} - \mathbf{a}x$, so

$$\begin{aligned}\mathbf{a} \cdot \mathbf{e} = 0 &\rightarrow \mathbf{a}^\top (\mathbf{b} - \mathbf{a}x) = 0 \\ \mathbf{a}^\top \mathbf{b} - \mathbf{a}^\top \mathbf{a}x &= 0 \\ x = \frac{\mathbf{a}^\top \mathbf{b}}{\mathbf{a}^\top \mathbf{a}} &= \frac{\langle \mathbf{a}, \mathbf{b} \rangle}{\|\mathbf{a}\|^2},\end{aligned}$$

and $\mathbf{p} = \mathbf{a}x = \mathbf{a} \frac{\mathbf{a}^\top \mathbf{b}}{\mathbf{a}^\top \mathbf{a}}$. Note, doubling \mathbf{b} doubles \mathbf{p} . Doubling \mathbf{a} does not affect \mathbf{p} . \mathbf{p} is called the **projection**.

We'd like to write this projection in terms of a **projection matrix** $P : \mathbf{p} = P\mathbf{b}$

$$\mathbf{a} \frac{\mathbf{a}^\top \mathbf{b}}{\mathbf{a}^\top \mathbf{a}} = P\mathbf{b} \rightarrow \left[\frac{\mathbf{a}\mathbf{a}^\top}{\|\mathbf{a}\|^2} \right] \mathbf{b} \rightarrow P = \frac{\mathbf{a}\mathbf{a}^\top}{\|\mathbf{a}\|^2}.$$

Note that $\mathbf{a}\mathbf{a}^\top$ is a three by three matrix, not a number; matrix multiplication is not commutative. The column space of P is spanned by \mathbf{a} because for any \mathbf{b} , $P\mathbf{b}$ lies on the line determined by \mathbf{a} . The rank of P is 1. P is symmetric. $P^2\mathbf{b} = P\mathbf{b}$ because the projection of a vector already on the line through \mathbf{a} is just that vector. In general, projection matrices have the properties:

$$P^\top = P \text{ and } P^2 = P$$

As we know, the equation $\mathbf{A}\mathbf{x} = \mathbf{b}$ may have no solution. The vector $\mathbf{A}\mathbf{x}$ is always in the column space of \mathbf{A} , and \mathbf{b} is unlikely to be in the column space. So, we project \mathbf{b} onto a vector \mathbf{p} in the column space of \mathbf{A} and solve $\mathbf{A}\hat{\mathbf{x}} = \mathbf{p}$.

More info can be found [here](#). □