# The Singular Value Decomposition

We are interested in more than just sym+def matrices. But the eigenvalue decompositions discussed in the last section of notes will play a major role in solving general systems of equations

$$\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}, \qquad \boldsymbol{y} \in \mathbb{R}^M, \quad \boldsymbol{A} \text{ is } M \times N, \quad \boldsymbol{x} \in \mathbb{R}^N.$$

We have seen that a symmetric positive definite matrix can be decomposed as $\boldsymbol{A} = \boldsymbol{V}\boldsymbol{\Lambda}\boldsymbol{V}^{\mathrm{T}}$, where $\boldsymbol{V}$ is an orthogonal matrix ($\boldsymbol{V}^{\mathrm{T}}\boldsymbol{V} = \boldsymbol{V}\boldsymbol{V}^{\mathrm{T}} = \boldsymbol{I}$) whose columns are the eigenvectors of $\boldsymbol{A}$, and $\boldsymbol{\Lambda}$ is a diagonal matrix containing the eigenvalues of $\boldsymbol{A}$. Because both orthogonal and diagonal matrices are trivial to invert, this eigenvalue decomposition makes it very easy to solve systems of equations $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}$ and analyze the stability of these solutions.

The **singular value decomposition** (SVD) takes apart an arbitrary $M \times N$ matrix $\boldsymbol{A}$ in a similar manner. The SVD of a real-valued $M \times N$ matrix $\boldsymbol{A}$ with rank[1] $R$ is

$$\boldsymbol{A} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}}$$

where

1. $\boldsymbol{U}$ is an $M \times R$ matrix

$$\boldsymbol{U} = \begin{bmatrix} \boldsymbol{u}_1 & | & \boldsymbol{u}_2 & | & \cdots & | & \boldsymbol{u}_R \end{bmatrix},$$

whose columns $\boldsymbol{u}_m \in \mathbb{R}^M$ are orthonormal. Note that while $\boldsymbol{U}^{\mathrm{T}}\boldsymbol{U} = \boldsymbol{I}$, in general $\boldsymbol{U}\boldsymbol{U}^{\mathrm{T}} \neq \boldsymbol{I}$ when $R < M$. The columns of $\boldsymbol{U}$ are an orthobasis for the range space of $\boldsymbol{A}$.

---

[1] Recall that the rank of a matrix is the number of linearly independent columns of a matrix (which is always equal to the number of linearly independent rows).

2. $V$ is an $N \times R$ matrix

$$V = \begin{bmatrix} v_1 & | & v_2 & | & \cdots & | & v_R \end{bmatrix},$$

whose columns $v_n \in \mathbb{R}^N$ are orthonormal. Again, while $V^{\mathrm{T}}V = I$, in general $VV^{\mathrm{T}} \neq I$ when $R < N$. The columns of $V$ are an orthobasis for the range space of $A^{\mathrm{T}}$ (recall that $\mathrm{Range}(A^{\mathrm{T}})$ consists of everything orthogonal to the nullspace of $A$).

3. $\Sigma$ is an $R \times R$ diagonal matrix with positive entries:

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & \cdots \\ 0 & \sigma_2 & 0 & \cdots \\ \vdots & & \ddots & \\ 0 & \cdots & \cdots & \sigma_R \end{bmatrix}.$$

We call the $\sigma_r$ the **singular values** of $A$. By convention, we will order them such that $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_R$.

4. The $v_1, \ldots, v_R$ are eigenvectors of the positive semi-definite matrix $A^{\mathrm{T}}A$. Note that

$$A^{\mathrm{T}}A = V\Sigma U^{\mathrm{T}}U\Sigma V^{\mathrm{T}} = V\Sigma^2 V^{\mathrm{T}},$$

and so the singular values $\sigma_1, \ldots, \sigma_R$ are the square roots of the non-zero eigenvalues of $A^{\mathrm{T}}A$.

5. Similarly,

$$AA^{\mathrm{T}} = U\Sigma^2 U^{\mathrm{T}},$$

and so the $u_1, \ldots, u_R$ are eigenvectors of the positive semi-definite matrix $AA^{\mathrm{T}}$. Since the non-zero eigenvalues of $A^{\mathrm{T}}A$ and $AA^{\mathrm{T}}$ are the same, the $\sigma_r$ are also square roots of the eigenvalues of $AA^{\mathrm{T}}$.

The rank $R$ is the dimension of the space spanned by the columns of $\boldsymbol{A}$, this is the same as the dimension of the space spanned by the rows. Thus $R \leq \min(M, N)$. We say $\boldsymbol{A}$ is **full rank** if $R = \min(M, N)$.

As before, we will often times find it useful to write the SVD as the sum of $R$ rank-1 matrices:

$$\boldsymbol{A} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}} = \sum_{r=1}^{R} \sigma_r\, \boldsymbol{u}_r \boldsymbol{v}_r^{\mathrm{T}}.$$

When $\boldsymbol{A}$ is **overdetermined** $(M > N)$, the decomposition looks like this

$$\begin{bmatrix} \boldsymbol{A} \end{bmatrix} = \begin{bmatrix} \boldsymbol{U} \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_R \end{bmatrix} \begin{bmatrix} \boldsymbol{V}^{\mathrm{T}} \end{bmatrix}.$$

When $\boldsymbol{A}$ is **underdetermined** $(M < N)$, the SVD looks like this

$$\begin{bmatrix} \boldsymbol{A} \end{bmatrix} = \begin{bmatrix} \boldsymbol{U} \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_R \end{bmatrix} \begin{bmatrix} \boldsymbol{V}^{\mathrm{T}} \end{bmatrix}.$$

When $\boldsymbol{A}$ is **square** and full rank $(M = N = R)$, the SVD looks like

$$\begin{bmatrix} \boldsymbol{A} \end{bmatrix} = \begin{bmatrix} \boldsymbol{U} \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_N \end{bmatrix} \begin{bmatrix} \boldsymbol{V}^{\mathrm{T}} \end{bmatrix}.$$

# The Least-Squares Problem

We can use the SVD to "solve" the general system of linear equations

$$\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}$$

where $\boldsymbol{y} \in \mathbb{R}^M$, $\boldsymbol{x} \in \mathbb{R}^N$, and $\boldsymbol{A}$ is an $M \times N$ matrix.

Given $\boldsymbol{y}$, we want to find $\boldsymbol{x}$ in such a way that

1. when there is a unique solution, we return it;
2. when there is no solution, we return something reasonable;
3. when there are an infinite number of solutions, we choose one to return in a "smart" way.

The **least-squares** framework revolves around finding an $\boldsymbol{x}$ that minimizes the length of the residual

$$\boldsymbol{r} = \boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}.$$

That is, we want to solve the optimization problem

$$\underset{\boldsymbol{x} \in \mathbb{R}^N}{\text{minimize}} \ \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_2^2, \tag{1}$$

where $\|\cdot\|_2$ is the standard Euclidean norm. We will see that the SVD of $\boldsymbol{A}$:

$$\boldsymbol{A} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}}, \tag{2}$$

plays a pivotal role in solving this problem.

To start, note that we can write any $\boldsymbol{x} \in \mathbb{R}^N$ as

$$\boldsymbol{x} = \boldsymbol{V}\boldsymbol{\alpha} + \boldsymbol{V}_0\boldsymbol{\alpha}_0. \tag{3}$$

Here, $\boldsymbol{V}$ is the $N \times R$ matrix appearing in the SVD decomposition (2), and $\boldsymbol{V}_0$ is a $N \times (N - R)$ matrix whose columns are orthogonal to one another and to the columns in $\boldsymbol{V}$. We have the relations

$$\boldsymbol{V}^{\mathrm{T}}\boldsymbol{V} = \mathbf{I}, \quad \boldsymbol{V}_0^{\mathrm{T}}\boldsymbol{V}_0 = \mathbf{I}, \quad \boldsymbol{V}^{\mathrm{T}}\boldsymbol{V}_0 = \mathbf{0}.$$

You can think of $\boldsymbol{V}_0$ as an orthobasis for the null space of $\boldsymbol{A}$. Of course, $\boldsymbol{V}_0$ is not unique, as there are many orthobases for $\mathrm{Null}(\boldsymbol{A})$, but any such set of vectors will serve our purposes here. The decomposition (3) is possible since $\mathrm{Range}(\boldsymbol{A}^{\mathrm{T}})$ and $\mathrm{Null}(\boldsymbol{A})$ partition $\mathbb{R}^N$ for any $M \times N$ matrix $\boldsymbol{A}$. Taking

$$\boldsymbol{\alpha} = \boldsymbol{V}^{\mathrm{T}}\boldsymbol{x}, \quad \boldsymbol{\alpha}_0 = \boldsymbol{V}_0^{\mathrm{T}}\boldsymbol{x},$$

we see that (3) holds since

$$\boldsymbol{x} = \boldsymbol{V}\boldsymbol{V}^{\mathrm{T}}\boldsymbol{x} + \boldsymbol{V}_0\boldsymbol{V}_0^{\mathrm{T}}\boldsymbol{x} = (\boldsymbol{V}\boldsymbol{V}^{\mathrm{T}} + \boldsymbol{V}_0\boldsymbol{V}_0^{\mathrm{T}})\boldsymbol{x} = \boldsymbol{x},$$

where we have made use of the fact that $\boldsymbol{V}\boldsymbol{V}^{\mathrm{T}} + \boldsymbol{V}_0\boldsymbol{V}_0^{\mathrm{T}} = \mathbf{I}$, because $\boldsymbol{V}\boldsymbol{V}^{\mathrm{T}}$ and $\boldsymbol{V}_0\boldsymbol{V}_0^{\mathrm{T}}$ are ortho-projectors onto complementary subspaces[2] of $\mathbb{R}^N$. So we can solve for $\boldsymbol{x} \in \mathbb{R}^N$ by solving for the pair $\boldsymbol{\alpha} \in \mathbb{R}^R$, $\boldsymbol{\alpha}_0 \in \mathbb{R}^{N-R}$.

Similarly, we can decompose $\boldsymbol{y}$ as

$$\boldsymbol{y} = \boldsymbol{U}\boldsymbol{\beta} + \boldsymbol{U}_0\boldsymbol{\beta}_0, \tag{4}$$

where $\boldsymbol{U}$ is the $M \times R$ matrix from the SVD decomposition, and $\boldsymbol{U}_0$ is a $M \times (M - R)$ complementary orthogonal basis. Again,

$$\boldsymbol{U}^{\mathrm{T}}\boldsymbol{U} = \mathbf{I}, \quad \boldsymbol{U}_0^{\mathrm{T}}\boldsymbol{U}_0 = \mathbf{I}, \quad \boldsymbol{U}^{\mathrm{T}}\boldsymbol{U}_0 = \mathbf{0},$$

---

[2]Subspaces $\mathcal{S}_1$ and $\mathcal{S}_2$ are **complementary** in $\mathbb{R}^N$ if $\mathcal{S}_1 \perp \mathcal{S}_2$ (everything in $\mathcal{S}_1$ is orthogonal to everything in $\mathcal{S}_2$) and $\mathcal{S}_1 \oplus \mathcal{S}_2 = \mathbb{R}^N$. You can think of $\mathcal{S}_1, \mathcal{S}_2$ as a partition of $\mathbb{R}^N$ into two orthogonal subspaces.

and we can think of $\boldsymbol{U}_0$ as an orthogonal basis for everything in $\mathbb{R}^M$ that is not in the range of $\boldsymbol{A}$. As before, we can calculate the decomposition above using

$$\boldsymbol{\beta} = \boldsymbol{U}^{\mathrm{T}}\boldsymbol{y}, \quad \boldsymbol{\beta}_0 = \boldsymbol{U}_0^{\mathrm{T}}\boldsymbol{y}.$$

Using the decompositions (2), (3), and (4) for $\boldsymbol{A}$, $\boldsymbol{x}$, and $\boldsymbol{y}$, we can write the residual $\boldsymbol{r} = \boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}$ as

$$\begin{aligned}
\boldsymbol{r} &= \boldsymbol{U}\boldsymbol{\beta} + \boldsymbol{U}_0\boldsymbol{\beta}_0 - \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}}(\boldsymbol{V}\boldsymbol{\alpha} + \boldsymbol{V}_0\boldsymbol{\alpha}_0) \\
&= \boldsymbol{U}\boldsymbol{\beta} + \boldsymbol{U}_0\boldsymbol{\beta}_0 - \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{\alpha} \quad (\text{since } \boldsymbol{V}^{\mathrm{T}}\boldsymbol{V} = \boldsymbol{I} \text{ and } \boldsymbol{V}^{\mathrm{T}}\boldsymbol{V}_0 = \boldsymbol{0}) \\
&= \boldsymbol{U}_0\boldsymbol{\beta}_0 + \boldsymbol{U}(\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}).
\end{aligned}$$

We want to choose $\boldsymbol{\alpha}$ that minimizes the energy of $\boldsymbol{r}$:

$$\begin{aligned}
\|\boldsymbol{r}\|_2^2 &= \langle \boldsymbol{U}_0\boldsymbol{\beta}_0 + \boldsymbol{U}(\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}), \ \boldsymbol{U}_0\boldsymbol{\beta}_0 + \boldsymbol{U}(\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}) \rangle \\
&= \langle \boldsymbol{U}_0\boldsymbol{\beta}_0, \boldsymbol{U}_0\boldsymbol{\beta}_0 \rangle \ + \ 2\langle \boldsymbol{U}_0\boldsymbol{\beta}_0, \boldsymbol{U}(\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}) \rangle \\
&\qquad\qquad + \langle \boldsymbol{U}(\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}), \boldsymbol{U}(\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}) \rangle \\
&= \|\boldsymbol{\beta}_0\|_2^2 + \|\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}\|_2^2
\end{aligned}$$

where the last equality comes from the facts that $\boldsymbol{U}_0^{\mathrm{T}}\boldsymbol{U}_0 = \boldsymbol{I}, \boldsymbol{U}^{\mathrm{T}}\boldsymbol{U} = \boldsymbol{I}$, and $\boldsymbol{U}^{\mathrm{T}}\boldsymbol{U}_0 = \boldsymbol{0}$. We have no control over $\|\boldsymbol{\beta}_0\|_2^2$, since it is determined entirely by our observations $\boldsymbol{y}$. Therefore, our problem has been reduced to finding $\boldsymbol{\alpha}$ that minimizes the second term $\|\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\alpha}\|_2^2$ above, which is non-negative. We can make it zero (i.e. as small as possible) by taking

$$\hat{\boldsymbol{\alpha}} = \boldsymbol{\Sigma}^{-1}\boldsymbol{\beta}.$$

Finally, the $\boldsymbol{x}$ which minimizes the residual (solves (1)) is

$$\hat{\boldsymbol{x}} = \boldsymbol{V}\hat{\boldsymbol{\alpha}} = \boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{\beta} = \boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{U}^{\mathrm{T}}\boldsymbol{y}. \tag{5}$$

Thus we can calculate the solution to (1) simply by applying the linear operator $\boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{U}^{\mathrm{T}}$ to the input data $\boldsymbol{y}$. There are two interesting facts about the solution $\hat{\boldsymbol{x}}$ in (5):

1. When $\boldsymbol{y} \in \mathrm{span}(\{\boldsymbol{u}_1, \ldots, \boldsymbol{u}_M\})$, we have $\boldsymbol{\beta}_0 = \boldsymbol{U}_0^{\mathrm{T}}\boldsymbol{y} = \boldsymbol{0}$, and so the residual $\boldsymbol{r} = \boldsymbol{0}$. In this case, there is at least one exact solution, and the one we choose satisfies $\boldsymbol{A}\hat{\boldsymbol{x}} = \boldsymbol{y}$.

2. Note that if $R < N$, then the solution is not unique. In this case, $\boldsymbol{V}_0$ has at least one column, and any part of a vector $\boldsymbol{x}$ in the range of $\boldsymbol{V}_0$ is not seen by $\boldsymbol{A}$, since

$$\boldsymbol{A}\boldsymbol{V}_0\boldsymbol{\alpha}_0 = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}}\boldsymbol{V}_0\boldsymbol{\alpha}_0 = \boldsymbol{0} \quad (\text{since } \boldsymbol{V}^{\mathrm{T}}\boldsymbol{V}_0 = \boldsymbol{0}).$$

As such,
$$\boldsymbol{x}' = \hat{\boldsymbol{x}} + \boldsymbol{V}_0\boldsymbol{\alpha}_0$$

for *any* $\boldsymbol{\alpha}_0 \in \mathbb{R}^{N-R}$ will have exactly the same residual, since $\boldsymbol{A}\boldsymbol{x}' = \boldsymbol{A}\hat{\boldsymbol{x}}$. In this case, our solution $\hat{\boldsymbol{x}}$ is the solution with smallest norm, since

$$\begin{aligned}
\|\boldsymbol{x}'\|_2^2 &= \langle \hat{\boldsymbol{x}} + \boldsymbol{V}_0\boldsymbol{\alpha}_0, \ \hat{\boldsymbol{x}} + \boldsymbol{V}_0\boldsymbol{\alpha}_0 \rangle \\
&= \langle \hat{\boldsymbol{x}}, \hat{\boldsymbol{x}} \rangle + 2\langle \hat{\boldsymbol{x}}, \boldsymbol{V}_0\boldsymbol{\alpha}_0 \rangle + \langle \boldsymbol{V}_0\boldsymbol{\alpha}, \boldsymbol{V}_0\boldsymbol{\alpha} \rangle \\
&= \|\hat{\boldsymbol{x}}\|_2^2 + 2\langle \boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{U}^{\mathrm{T}}\boldsymbol{y}, \boldsymbol{V}_0\boldsymbol{\alpha}_0 \rangle + \|\boldsymbol{\alpha}_0\|_2^2 \quad (\text{since } \boldsymbol{V}_0^{\mathrm{T}}\boldsymbol{V}_0 = \boldsymbol{I}) \\
&= \|\hat{\boldsymbol{x}}\|_2^2 + \|\boldsymbol{\alpha}_0\|_2^2 \quad (\text{since } \boldsymbol{V}^{\mathrm{T}}\boldsymbol{V}_0 = \boldsymbol{0})
\end{aligned}$$

which is minimized by taking $\boldsymbol{\alpha}_0 = \boldsymbol{0}$.

To summarize, $\hat{\boldsymbol{x}} = \boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{U}^{\mathrm{T}}\boldsymbol{y}$ has the desired properties stated at the beginning of this module, since

1. when $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}$ has a unique exact solution, it must be $\hat{\boldsymbol{x}}$,

2. when an exact solution is not available, $\hat{\boldsymbol{x}}$ is the solution to (1),

43

3. when there are an infinite number of minimizers to (1), $\hat{\boldsymbol{x}}$ is the one with smallest norm.

Because the matrix $\boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{U}^{\mathrm{T}}$ gives us such an elegant solution to this problem, we give it a special name: the **pseudo-inverse**.

## The Pseudo-Inverse

The **pseudo-inverse** of a matrix $\boldsymbol{A}$ with singular value decomposition (SVD) $\boldsymbol{A} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}}$ is

$$\boldsymbol{A}^{\dagger} = \boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{U}^{\mathrm{T}}. \tag{6}$$

Other names for $\boldsymbol{A}^{\dagger}$ include **natural inverse**, **Lanczos inverse**, and **Moore-Penrose inverse**.

Given an observation $\boldsymbol{y}$, taking $\hat{\boldsymbol{x}} = \boldsymbol{A}^{\dagger}\boldsymbol{y}$ gives us the **least squares** solution to $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}$. The pseudo-inverse $\boldsymbol{A}^{\dagger}$ always exists, since every matrix (with rank $R$) has an SVD decomposition $\boldsymbol{A} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}}$ with $\boldsymbol{\Sigma}$ as an $R \times R$ diagonal matrix with $\Sigma[r, r] > 0$.

When $\boldsymbol{A}$ is full rank ($R = \min(M, N)$), then we can calculate the pseudo-inverse without using the SVD. There are three cases:

- When $\boldsymbol{A}$ is square and invertible ($R = M = N$), then

$$\boldsymbol{A}^{\dagger} = \boldsymbol{A}^{-1}.$$

  This is easy to check, as here

$$\boldsymbol{A} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}} \quad \text{where both } \boldsymbol{U}, \boldsymbol{V} \text{ are } N \times N,$$

44

and since in this case $\boldsymbol{V}\boldsymbol{V}^{\mathrm{T}} = \boldsymbol{V}^{\mathrm{T}}\boldsymbol{V} = \mathbf{I}$ and $\boldsymbol{U}\boldsymbol{U}^{\mathrm{T}} = \boldsymbol{U}^{\mathrm{T}}\boldsymbol{U} = \mathbf{I}$,

$$\begin{aligned}
\boldsymbol{A}^{\dagger}\boldsymbol{A} &= \boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{U}^{\mathrm{T}}\boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}} \\
&= \boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}} \\
&= \boldsymbol{V}\boldsymbol{V}^{\mathrm{T}} \\
&= \mathbf{I}.
\end{aligned}$$

Similarly, $\boldsymbol{A}\boldsymbol{A}^{\dagger} = \mathbf{I}$, and so $\boldsymbol{A}^{\dagger}$ is both a left and right inverse of $\boldsymbol{A}$, and thus $\boldsymbol{A}^{\dagger} = \boldsymbol{A}^{-1}$.

- When $\boldsymbol{A}$ more rows than columns and has full column rank ($R = N \leq M$), then $\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}$ is invertible, and

$$\boldsymbol{A}^{\dagger} = (\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A})^{-1}\boldsymbol{A}^{\mathrm{T}}. \tag{7}$$

This type of $\boldsymbol{A}$ is "tall and skinny"

$$\begin{bmatrix} \boldsymbol{A} \\ \phantom{x} \end{bmatrix},$$

and its columns are linearly independent. To verify equation (7), recall that

$$\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A} = \boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{U}^{\mathrm{T}}\boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}} = \boldsymbol{V}\boldsymbol{\Sigma}^{2}\boldsymbol{V}^{\mathrm{T}},$$

and so

$$(\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A})^{-1}\boldsymbol{A}^{\mathrm{T}} = \boldsymbol{V}\boldsymbol{\Sigma}^{-2}\boldsymbol{V}^{\mathrm{T}}\boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{U}^{\mathrm{T}} = \boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{U}^{\mathrm{T}},$$

which is exactly the content of (6).

- When $\boldsymbol{A}$ has more columns than rows and has full row rank ($R = M \leq N$), then $\boldsymbol{A}\boldsymbol{A}^{\mathrm{T}}$ is invertible, and

$$\boldsymbol{A}^{\dagger} = \boldsymbol{A}^{\mathrm{T}}(\boldsymbol{A}\boldsymbol{A}^{\mathrm{T}})^{-1}. \qquad (8)$$

  This occurs when $\boldsymbol{A}$ is "short and fat"

$$\begin{bmatrix} & & \\ & \boldsymbol{A} & \\ & & \end{bmatrix},$$

  and its rows are linearly independent. To verify equation (8), recall that

$$\boldsymbol{A}\boldsymbol{A}^{\mathrm{T}} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}}\boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{U}^{\mathrm{T}} = \boldsymbol{U}\boldsymbol{\Sigma}^2\boldsymbol{U}^{\mathrm{T}},$$

  and so

$$\boldsymbol{A}^{\mathrm{T}}(\boldsymbol{A}\boldsymbol{A}^{\mathrm{T}})^{-1} = \boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{U}^{\mathrm{T}}\boldsymbol{U}\boldsymbol{\Sigma}^{-2}\boldsymbol{U}^{\mathrm{T}} = \boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{U}^{\mathrm{T}},$$

  which again is exactly (6).

## $\boldsymbol{A}^{\dagger}$ is as close to an inverse of $\boldsymbol{A}$ as possible

As discussed in above, when $\boldsymbol{A}$ is square and invertible, $\boldsymbol{A}^{\dagger}$ is exactly the inverse of $\boldsymbol{A}$. When $\boldsymbol{A}$ is not square, we can ask if there is a better right or left inverse. We will argue that there is not.

**Left inverse** Given $\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}$, we would like $\boldsymbol{A}^{\dagger}\boldsymbol{y} = \boldsymbol{A}^{\dagger}\boldsymbol{A}\boldsymbol{x} = \boldsymbol{x}$ for any $\boldsymbol{x}$. That is, we would like $\boldsymbol{A}^{\dagger}$ to be a *left inverse* of $\boldsymbol{A}$: $\boldsymbol{A}^{\dagger}\boldsymbol{A} = \boldsymbol{I}$. Of course, this is not always possible, especially when $\boldsymbol{A}$ has more columns than rows, $M < N$. But we can ask if any other matrix $\boldsymbol{H}$ comes closer to being a left inverse

than $\boldsymbol{A}^\dagger$. To find the "best" left-inverse, we look for the matrix which minimizes

$$\min_{\boldsymbol{H} \in \mathbb{R}^{N \times M}} \|\boldsymbol{H}\boldsymbol{A} - \mathbf{I}\|_F^2. \tag{9}$$

Here, $\| \cdot \|_F$ is the *Frobenius norm*, defined for an $N \times M$ matrix $\boldsymbol{Q}$ as the sum of the squares of the entries:[3]

$$\|\boldsymbol{Q}\|_F^2 = \sum_{n=1}^{M} \sum_{n=1}^{N} |Q[m, n]|^2$$

With (9), we are finding $\boldsymbol{H}$ such that $\boldsymbol{H}\boldsymbol{A}$ is as close to the identity as possible in the least-squares sense.

The pseudo-inverse $\boldsymbol{A}^\dagger$ minimizes (9). To see this, recognize (see the exercise below) that the solution $\hat{\boldsymbol{H}}$ to (9) must obey

$$\boldsymbol{A}\boldsymbol{A}^{\mathrm{T}}\hat{\boldsymbol{H}}^{\mathrm{T}} = \boldsymbol{A}. \tag{10}$$

We can see that this is indeed true for $\hat{\boldsymbol{H}} = \boldsymbol{A}^\dagger$:

$$\boldsymbol{A}\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}^{\dagger\mathrm{T}} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}}\boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{U}^{\mathrm{T}}\boldsymbol{U}\boldsymbol{\Sigma}^{-1}\boldsymbol{V}^{\mathrm{T}} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}} = \boldsymbol{A}.$$

So there is no $N \times M$ matrix that is closer to being a left inverse than $\boldsymbol{A}^\dagger$.

---

[3]It is also true that $\|\boldsymbol{Q}\|_F^2$ is the sum of the squares of the singular values of $\boldsymbol{Q}$: $\|\boldsymbol{Q}\|_F^2 = \lambda_1^2 + \cdots + \lambda_p^2$. This is something that you will prove on the next homework.

**Right inverse** If we re-apply $\boldsymbol{A}$ to our solution $\hat{\boldsymbol{x}} = \boldsymbol{A}^\dagger \boldsymbol{y}$, we would like it to be as close as possible to our observations $\boldsymbol{y}$. That is, we would like $\boldsymbol{A}\boldsymbol{A}^\dagger$ to be as close to the identity as possible. Again, achieving this goal exactly is not always possible, especially if $\boldsymbol{A}$ has more rows that columns. But we can attempt to find the "best" right inverse, in the least-squares sense, by solving

$$\underset{\boldsymbol{H} \in \mathbb{R}^{N \times M}}{\text{minimize}} \ \|\boldsymbol{A}\boldsymbol{H} - \mathbf{I}\|_F^2. \tag{11}$$

The solution $\hat{\boldsymbol{H}}$ to (11) (see the exercise below) must obey

$$\boldsymbol{A}^\mathrm{T}\boldsymbol{A}\hat{\boldsymbol{H}} = \boldsymbol{A}^\mathrm{T}. \tag{12}$$

Again, we show that $\boldsymbol{A}^\dagger$ satisfies (12), and hence is a minimizer to (11):

$$\boldsymbol{A}^\mathrm{T}\boldsymbol{A}\boldsymbol{A}^\dagger = \boldsymbol{V}\boldsymbol{\Sigma}^2\boldsymbol{V}^\mathrm{T}\boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{U}^\mathrm{T} = \boldsymbol{V}\boldsymbol{\Sigma}\boldsymbol{U}^\mathrm{T} = \boldsymbol{A}^\mathrm{T}.$$

---

Moral:
$\boldsymbol{A}^\dagger = \boldsymbol{V}\boldsymbol{\Sigma}^{-1}\boldsymbol{U}^\mathrm{T}$ **is as close (in the least-squares sense) to an inverse of $\boldsymbol{A}$ as you could possibly have**.

---

**Exercise:**

Show that the minimizer $\hat{\boldsymbol{H}}$ to (9) must obey (10). Do this by using the fact that the derivative of the functional $\|\boldsymbol{H}\boldsymbol{A} - \mathbf{I}\|_F^2$ with respect to an entry $H[k, \ell]$ in $\boldsymbol{H}$ must obey

$$\frac{\partial \|\boldsymbol{H}\boldsymbol{A} - \mathbf{I}\|_F^2}{\partial H[k, \ell]} = 0, \quad \text{for all } 1 \le k \le N, \ 1 \le \ell \le M,$$

to be a solution to (9). Do the same for (11) and (12).

# Technical Details: Existence of the SVD

In this section we will prove that any $M \times N$ matrix $\boldsymbol{A}$ with $\mathrm{rank}(\boldsymbol{A}) = R$ can be written as

$$\boldsymbol{A} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathrm{T}}$$

where $\boldsymbol{U}, \boldsymbol{\Sigma}, \boldsymbol{V}$ have the five properties listed at the beginning of the last section.

Since $\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}$ is symmetric positive semi-definite, we can write:

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

where the $\boldsymbol{v}_n$ are orthonormal and the $\lambda_n$ are real and non-negative. Since $\mathrm{rank}(\boldsymbol{A}) = R$, we also have $\mathrm{rank}(\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}) = R$, and so $\lambda_1, \ldots, \lambda_R$ are all strictly positive above, and $\lambda_{R+1} = \cdots = \lambda_N = 0$.

Set

$$\boldsymbol{u}_m = \frac{1}{\sqrt{\lambda_m}}\boldsymbol{A}\boldsymbol{v}_m, \quad \text{for } m = 1, \ldots, R, \qquad \boldsymbol{U} = \begin{bmatrix} \boldsymbol{u}_1 & \cdots & \boldsymbol{u}_R \end{bmatrix}.$$

Notice that these $\boldsymbol{u}_m$ are orthonormal, as

$$\langle \boldsymbol{u}_m, \boldsymbol{u}_\ell \rangle = \frac{1}{\sqrt{\lambda_m \lambda_\ell}} \boldsymbol{v}_\ell^{\mathrm{T}} \boldsymbol{A}^{\mathrm{T}} \boldsymbol{A} \boldsymbol{v}_m = \sqrt{\frac{\lambda_m}{\lambda_\ell}} \boldsymbol{v}_\ell^{\mathrm{T}} \boldsymbol{v}_m = \begin{cases} 1, & m = \ell, \\ 0, & m \neq \ell. \end{cases}$$

These $\boldsymbol{u}_m$ also happen to be eigenvectors of $\boldsymbol{A}\boldsymbol{A}^{\mathrm{T}}$, as

$$\boldsymbol{A}\boldsymbol{A}^{\mathrm{T}}\boldsymbol{u}_m = \frac{1}{\sqrt{\lambda_m}}\boldsymbol{A}\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}\boldsymbol{v}_m = \sqrt{\lambda_m}\boldsymbol{A}\boldsymbol{v}_m = \lambda_m \boldsymbol{u}_m.$$

Now let $\boldsymbol{u}_{R+1}, \ldots, \boldsymbol{u}_M$ be an orthobasis for the null space of $\boldsymbol{U}^{\mathrm{T}}$ — concatenating these two sets into $\boldsymbol{u}_1, \ldots, \boldsymbol{u}_M$ forms an orthobasis for all of $\mathbb{R}^M$.

Let $\boldsymbol{V} = \begin{bmatrix} \boldsymbol{v}_1 & \boldsymbol{v}_2 & \cdots & \boldsymbol{v}_R \end{bmatrix}$. In addition, let

$$\boldsymbol{V}_0 = \begin{bmatrix} \boldsymbol{v}_{R+1} & \boldsymbol{v}_{R+2} & \cdots & \boldsymbol{v}_N \end{bmatrix}, \quad \boldsymbol{V}_{\text{full}} = \begin{bmatrix} \boldsymbol{V} & \boldsymbol{V}_0 \end{bmatrix}$$

and

$$\boldsymbol{U}_0 = \begin{bmatrix} \boldsymbol{u}_{R+1} & \boldsymbol{u}_{R+2} & \cdots & \boldsymbol{u}_M \end{bmatrix}, \quad \boldsymbol{U}_{\text{full}} = \begin{bmatrix} \boldsymbol{U} & \boldsymbol{U}_0 \end{bmatrix}.$$

It should be clear that $\boldsymbol{V}_{\text{full}}$ is an $N \times N$ orthonormal matrix and $\boldsymbol{U}_{\text{full}}$ is a $M \times M$ orthonormal matrix. Consider the $M \times N$ matrix $\boldsymbol{U}_{\text{full}}^{\mathrm{T}} \boldsymbol{A} \boldsymbol{V}_{\text{full}}$ — the entry in the $m^{\text{th}}$ rows and $n^{\text{th}}$ column of this matrix is

$$(\boldsymbol{U}_{\text{full}}^{\mathrm{T}} \boldsymbol{A} \boldsymbol{V}_{\text{full}})[m,n] = \boldsymbol{u}_m^{\mathrm{T}} \boldsymbol{A} \boldsymbol{v}_n = \begin{cases} \sqrt{\lambda_n}\, \boldsymbol{u}_m^{\mathrm{T}} \boldsymbol{u}_n & n = 1, \ldots, R \\ 0, & n = R+1, \ldots, N. \end{cases}$$

$$= \begin{cases} \sqrt{\lambda_n}, & m = n = 1, \ldots, R \\ 0, & \text{otherwise.} \end{cases}$$

Thus

$$\boldsymbol{U}_{\text{full}}^{\mathrm{T}} \boldsymbol{A} \boldsymbol{V}_{\text{full}} = \boldsymbol{\Sigma}_{\text{full}}$$

where

$$\Sigma_{\text{full}}[m,n] = \begin{cases} \sqrt{\lambda_n}, & m = n = 1, \ldots, R \\ 0, & \text{otherwise.} \end{cases}$$

Since $\boldsymbol{U}_{\text{full}} \boldsymbol{U}_{\text{full}}^{\mathrm{T}} = \mathbf{I}$ and $\boldsymbol{V}_{\text{full}} \boldsymbol{V}_{\text{full}}^{\mathrm{T}} = \mathbf{I}$, we have

$$\boldsymbol{A} = \boldsymbol{U}_{\text{full}} \boldsymbol{\Sigma}_{\text{full}} \boldsymbol{V}_{\text{full}}^{\mathrm{T}}.$$

Since $\boldsymbol{\Sigma}_{\text{full}}$ is non-zero only in the first $R$ locations along its main diagonal, the above reduces to

$$\boldsymbol{A} = \boldsymbol{U} \boldsymbol{\Sigma} \boldsymbol{V}^{\mathrm{T}}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \sqrt{\lambda_1} & & & \\ & \sqrt{\lambda_2} & & \\ & & \ddots & \\ & & & \sqrt{\lambda_R} \end{bmatrix}.$$