

Applications (This lecture and next 2-3)

- Numerical linear algebra: rank, condition number, bases for fundamental subspaces, pseudoinverse
- Compression: Principle Component Analysis
- Recommender Systems
- Ranking sports teams.

Topics

- Singular values and singular value decomposition.

Notes based on ALA 8.7 and LAA 7.4.

Warmup

The diagonalization theorems we've seen for complete and symmetric matrices have played a role in many interesting applications. Unfortunately, not all matrices can be factored as $A = PDP^{-1}$ for a diagonal matrix D ; for example, such a factorization makes no sense if A is not square. Fortunately, a factorization $A = PDQ^{-1}$ is possible for any $m \times n$ matrix A ! A special factorization of this type, called the **singular value decomposition**, is one of the most useful and widely applicable matrix factorizations in linear algebra.

The singular value decomposition is based on the following key property of matrix diagonalization which we'll show can be captured in general rectangular matrices:

Key observation: The absolute values of the eigenvalues of a symmetric matrix A measure the amounts that A stretches or shrinks certain vectors (the eigenvectors). If $Ax = \lambda x$ and $\|x\| = 1$, then

$$\|Ax\| = \|\lambda x\| = |\lambda| \|x\| = |\lambda|.$$

If λ_1 is the eigenvalue with the greatest magnitude, i.e., if $|\lambda_1| \geq |\lambda_i|$ for $i=2, \dots, n$, then a corresponding unit eigenvector v_1 identifies the direction in which stretching is greatest. That is, the length of Ax is maximized when $x = v_1$, and $\|Ax\| = |\lambda_1|$.

This description is reminiscent of the optimization principle we saw for characterizing eigenvalues of symmetric matrices, albeit with a focus on maximizing length $\|Ax\|$ rather than the quadratic form $x^T Ax$. What we'll see next is that this description of v_1 and $|\lambda_1|$ has an analogue for rectangular matrices that will lead to the singular value decomposition.

Example: The matrix $A = \begin{bmatrix} 4 & 11 & 14 \\ 8 & 7 & -2 \end{bmatrix}$ defines a linear map $x \mapsto Ax$

from \mathbb{R}^3 to \mathbb{R}^2 . If we consider the effects of this map on the unit sphere $\{x \in \mathbb{R}^3 \mid \|x\| = 1\}$, we observe that multiplication by A transforms this sphere in \mathbb{R}^3 into an ellipse in \mathbb{R}^2 :

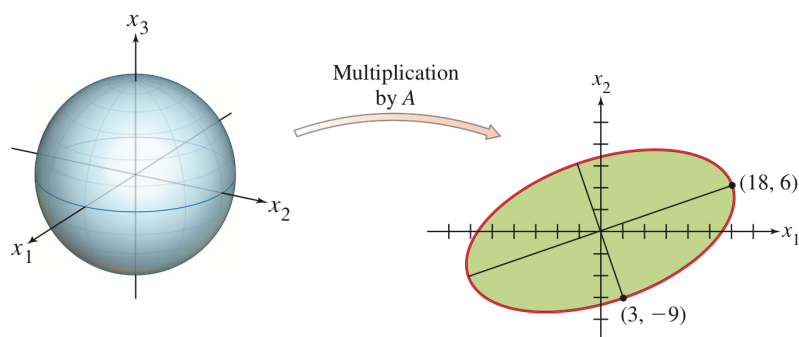


FIGURE 1 A transformation from \mathbb{R}^3 to \mathbb{R}^2 .

(LAA Ch 9.4)

Our task is to find a unit vector \pm at which the length $\|A\pm\|$ is maximized, and compute this maximum length. That is, we want to solve the optimization problem

$$\text{maximize } \|A\pm\|$$

over choices of \pm satisfying $\|\pm\|=1$. Our first observation is that the quantity $\|A\pm\|^2$ is maximized by the same \pm that maximizes $\|A\pm\|$, but that $\|A\pm\|^2$ is easier to work with. Specifically, note that

$$\|A\pm\|^2 = \langle A\pm, A\pm \rangle = (A\pm)^T (A\pm) = \pm^T A^T A \pm = \pm^T (A^T A) \pm.$$

So our task is to now find a unit vector $\|\pm\|=1$ that maximizes the quadratic form $\pm^T (A^T A) \pm$ defined by the symmetric (positive semidefinite) matrix $A^T A$: we know how to do this. By our theorem characterizing eigenvalues of symmetric matrices from an optimization perspective, we know the maximum value is the largest eigenvalue λ_1 of the matrix $A^T A$, and is attained at the unit eigenvector \underline{v}_1 of $A^T A$ corresponding to λ_1 .

For the matrix in this example:

$$A^T A = \begin{bmatrix} 4 & 8 \\ 11 & 7 \\ 14 & -2 \end{bmatrix} \begin{bmatrix} 4 & 11 & 14 \\ 8 & 7 & -2 \end{bmatrix} = \begin{bmatrix} 80 & 100 & 40 \\ 100 & 120 & 140 \\ 40 & 140 & 200 \end{bmatrix}$$

and the eigenvalue/vector pairs are:

$$\lambda_1 = 360, \underline{v}_1 = \begin{bmatrix} 1/3 \\ 2/3 \\ 2/3 \end{bmatrix}, \quad \lambda_2 = 90, \underline{v}_2 = \begin{bmatrix} -2/3 \\ -1/3 \\ 2/3 \end{bmatrix}, \quad \lambda_3 = 0, \underline{v}_3 = \begin{bmatrix} 2/3 \\ -2/3 \\ 1/3 \end{bmatrix}.$$

The maximum value of $\pm^T (A^T A) \pm = \|A\pm\|^2$ is thus $\lambda_1 = 360$, and attained when $\pm = \underline{v}_1$. The vector $A\underline{v}_1$ is a point on the ellipse in Fig 7 above farthest from the origin, namely

$$A\underline{v}_1 = \begin{bmatrix} 4 & 11 & 14 \\ 8 & 7 & -2 \end{bmatrix} \begin{bmatrix} 1/3 \\ 2/3 \\ 2/3 \end{bmatrix} = \begin{bmatrix} 18 \\ 6 \end{bmatrix}.$$

For $\|\pm\|=1$, the maximum value of $\|A\pm\|$ is $\|A\underline{v}_1\| = \sqrt{360} = 6\sqrt{10}$.

This example suggests that the effect of a matrix A on the unit sphere in \mathbb{R}^3 is related to the quadratic form $\pm^T A^T A \pm$. What we'll see next is that the entire geometric behavior of the map $\pm \mapsto A\pm$ is captured by this quadratic form.

The Singular Values of an $m \times n$ Matrix

Consider an $m \times n$ real matrix $A \in \mathbb{R}^{m \times n}$. Then $A^T A$ is an $n \times n$ symmetric matrix, and can be orthogonally diagonalized. Let $V = [v_1 \dots v_n]$ be an orthogonal matrix composed of orthonormal eigenvectors of $A^T A$, and let $\lambda_1, \dots, \lambda_n$ be the associated eigenvalues of $A^T A$. Then for $i=1, \dots, n$,

$$\begin{aligned}\|A v_i\|^2 &= (A v_i)^T (A v_i) = v_i^T (A^T A v_i) \\ &= v_i^T (\lambda_i v_i) \\ &= \lambda_i v_i^T v_i = \lambda_i \|v_i\|^2 \\ &= \lambda_i.\end{aligned}$$

This tells us that all of the eigenvalues $\lambda_i = \|A v_i\|^2 \geq 0$, since norms can only take on nonnegative values, i.e., $A^T A$ is a positive semidefinite matrix. Let's assume that we've ordered our eigenvalues in decreasing order:

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0.$$

The **singular values** of A are the positive square roots of the nonzero eigenvalues $\lambda_i > 0$ of $A^T A$, denoted σ_i . That is, let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$, and $\lambda_{r+1} = \lambda_{r+2} = \dots = \lambda_n = 0$ be a partition of the eigenvalues such that $\lambda_i > 0$ for $i=1, \dots, r$, and $\lambda_i = 0$ for $i=r+1, \dots, n$. Then A has r singular values, defined as

$$\sigma_i = \sqrt{\lambda_i}, \quad i=1, \dots, r.$$

WARNING: Some texts include the zero eigenvalues $\lambda_{r+1}, \dots, \lambda_n$ of $A^T A$ as singular values of A . This is simply a different convention, and is mathematically equivalent. However, we find our definition to be more natural for our purposes.

Example: Using the same $A = \begin{bmatrix} 4 & 11 & 14 \\ 8 & 7 & -2 \end{bmatrix}$ as the previous example,

we have $\sigma_1 = \sqrt{360} = 6\sqrt{10}$, $\sigma_2 = \sqrt{90} = 3\sqrt{10}$. In this case, A only has two singular values as $\lambda_3 = 0$. For this example, $r=2$, and $\lambda_1 = 360 > \lambda_2 = 90 > \lambda_3 = 0$. From the previous example, the first singular value of A is the maximum of $\|A x\|$ over all $\|x\|=1$, attained at v_1 . Our optimization based characterization of eigenvalues of symmetric matrices tells us that the second singular value of A is the maximum of $\|A x\|$ over all unit vectors **orthogonal** to v_1 : this is attained by the second eigenvector v_2 of $A^T A$. For v_2 from the previous example,

$$A v_2 = \begin{bmatrix} 4 & 11 & 14 \\ 8 & 7 & -2 \end{bmatrix} \begin{bmatrix} -2/3 \\ -1/3 \\ 2/3 \end{bmatrix} = \begin{bmatrix} 3 \\ 9 \end{bmatrix}.$$

This point is on the minor axis of the ellipse in Fig. 1 above, just as $A v_1$ is on the major axis (see Fig. 2 below). The two singular values of A are the lengths

of the major and minor semiaxes of the ellipse.

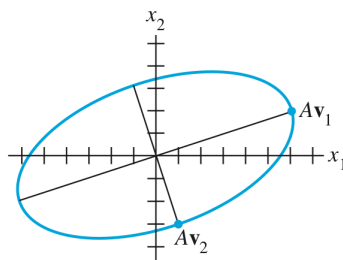


FIGURE 2

That Av_1 and Av_2 are orthogonal is no accident, as the next theorem shows.

Theorem: Suppose that u_1, \dots, u_n is an orthonormal basis for \mathbb{R}^n composed of the eigenvectors of $A^T A$, ordered so that the corresponding eigenvalues of $A^T A$ satisfy

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > \lambda_{r+1} = \dots = \lambda_n = 0,$$

where r denotes the number of nonzero eigenvalues of $A^T A$, i.e., the number of singular values $\sigma_i = \sqrt{\lambda_i} > 0$, $i=1, \dots, r$, of A . Then, Av_1, \dots, Av_r is an orthogonal basis for $\text{Col}(A)$, and $\text{rank}(A) = r$.

Proof: Because v_i and $\lambda_j v_j$ are orthogonal for $i \neq j$,

$$(Av_i)^T (Av_j) = v_i^T A^T A v_j = v_i^T \lambda_j v_j = 0.$$

Thus, Av_1, \dots, Av_r are mutually orthogonal, and hence linearly independent. They are also clearly contained in $\text{Col}(A)$. Now, for any $y \in \text{Col}(A)$, there must be an $x \in \mathbb{R}^n$ such that $y = Ax$. Expanding x in the basis v_1, \dots, v_n , as $x = c_1 v_1 + \dots + c_n v_n$ for some $c_1, \dots, c_n \in \mathbb{R}$, we have:

$$\begin{aligned} y = Ax &= A(c_1 v_1 + \dots + c_n v_n) = c_1 Av_1 + \dots + c_r Av_r + c_{r+1} Av_{r+1} + \dots + c_n Av_n \\ &= c_1 Av_1 + \dots + c_r Av_r. \end{aligned}$$

We used that $\|Av_i\|^2 = \lambda_i = 0$ for $i=r+1, \dots, n$ ($\Rightarrow Av_i = 0$ for $i=r+1, \dots, n$) in the last equality.

Therefore, we have that $y \in \text{span}\{Av_1, \dots, Av_r\}$. Thus Av_1, \dots, Av_r is both linearly independent and a spanning set for $\text{Col}(A)$, meaning it is an orthogonal basis for $\text{Col}(A)$. Hence, by the Fundamental Theorem of Linear Algebra:

$$\text{rank}(A) = \dim(\text{Col}(A)) = r.$$

Numerical Note: In certain cases, the rank of A may be very sensitive to small changes in the entries of A . The naive approach of counting the # of pivot columns in A does not work well if A is row reduced by a computer, as roundoff errors often create a row echelon form with full rank. In practice, the most reliable way of computing the rank of a large matrix A is to count the number of singular values larger than a small threshold ϵ (typically on the order of 10^{-12} , but can vary depend on applications). In this case, singular values smaller than ϵ are treated as zeros for all practical purposes, and the effective rank of A is computed by counting the remaining nonzero singular values.

The Singular Value Decomposition

The decomposition of A involves an $r \times r$ diagonal matrix Σ of the form

$$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r).$$

We note that because $r = \dim \text{Col}(A) = \dim \text{Row}(A)$ by the FTLA, we must have that $r \leq \min\{m, n\}$ if $A \in \mathbb{R}^{m \times n}$.

Theorem: Let $A \in \mathbb{R}^{m \times n}$ be an $m \times n$ matrix of rank $r > 0$. Then A can be factored as

$$A = U \Sigma V^T,$$

where $U \in \mathbb{R}^{m \times r}$ has orthonormal columns, so $U^T U = I_r$, $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r)$ is a diagonal matrix with the singular values of A σ_i along the diagonal, and $V \in \mathbb{R}^{n \times r}$ has orthonormal columns, so $V^T V = I_r$.

Such a factorization of A is called its singular value decomposition, and the columns of U are called the left singular vectors of A , while the columns of V are called the right singular vectors of A .

Proof: Let λ_i and \underline{v}_i be the eigenvalues/vectors of $A^T A$ as described previously, so that $A\underline{v}_1, \dots, A\underline{v}_r$ is an orthogonal basis for $\text{col}(A)$. Normalize each $A\underline{v}_i$ to form an orthonormal basis for $\text{col}(A)$:

$$\underline{u}_i = \frac{1}{\|A\underline{v}_i\|} A\underline{v}_i = \frac{1}{\sigma_i} A\underline{v}_i$$

and hence $A\underline{v}_i = \sigma_i \underline{u}_i$ for $i = 1, \dots, r$. Define the matrices

$$U = [\underline{u}_1 \dots \underline{u}_r] \in \mathbb{R}^{m \times r} \quad \text{and} \quad V = [\underline{v}_1 \dots \underline{v}_r] \in \mathbb{R}^{n \times r}$$

By construction, the columns of U are orthonormal: $U^T U = I_r$, and similarly for the columns of V : $V^T V = I_r$.

Let's define the following "full" matrices:

$$\hat{U} = [U \ U^\perp] \in \mathbb{R}^{m \times m} \quad \text{and} \quad \hat{V} = [V \ V^\perp] \in \mathbb{R}^{n \times n}.$$

Here, $V^\perp = [v_{r+1} \ \dots \ v_n]$ has orthonormal columns spanning the orthogonal complement of $\text{span}\{v_1, \dots, v_r\}$, so that the columns of \hat{V} form an orthonormal basis of \mathbb{R}^n .

Similarly, let U^\perp have orthonormal columns spanning the orthogonal complement of $\text{span}\{u_1, \dots, u_r\}$, so the columns of \hat{U} form an orthonormal basis for \mathbb{R}^m .

Finally, define $\hat{\Sigma} = \begin{bmatrix} \underbrace{\Sigma}_{r \times r} & \underbrace{0}_{r \times (n-r)} \\ 0 & 0 \end{bmatrix} \begin{matrix} \text{ } \\ \text{ } \end{matrix} \begin{matrix} \text{ } \\ \text{ } \end{matrix}$. We first show that

$$A = \hat{U} \hat{\Sigma} \hat{V}^T, \quad \text{or equivalently (since } \hat{V} \text{ is orthogonal), } A\hat{V} = \hat{U} \hat{\Sigma}.$$

$$A\hat{V} = [Av_1 \ \dots \ Av_r \ Av_{r+1} \ \dots \ Av_n] = [G_1 u_1 \ \dots \ G_r u_r \ 0 \ \dots \ 0],$$

Then, notice:

$$\hat{U} \hat{\Sigma} = [u_1 \ \dots \ u_r \ u_{r+1} \ \dots \ u_m] \begin{bmatrix} \overbrace{G_1 \ 0 \ \dots \ 0}^r & \overbrace{0 \ \dots \ 0}^{n-r} \\ \vdots & \vdots \\ 0 \ 0 \ \dots \ 0 & 0 \ \dots \ 0 \\ \hline 0 \ \dots \ 0 & \\ \vdots & \\ 0 \ \dots \ 0 & \end{bmatrix} \begin{matrix} \text{ } \\ \text{ } \end{matrix} \begin{matrix} \text{ } \\ \text{ } \end{matrix}$$

$$= [G_1 u_1 \ \dots \ G_r u_r \ 0 \ \dots \ 0].$$

So that $A\hat{V} = \hat{U} \hat{\Sigma}$, or equivalently, $A = \hat{U} \hat{\Sigma} \hat{V}^T$. But, now, notice:

$$A = \hat{U} \hat{\Sigma} \hat{V}^T = [U \ U^\perp] \begin{bmatrix} \Sigma & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V^T \\ (V^\perp)^T \end{bmatrix} = U \Sigma V^T, \quad \text{proving our result.}$$

NOTE: Some textbooks define the singular value decomposition of A as $A = \hat{U} \hat{\Sigma} \hat{V}^T$ — this is necessary when allowing for singular values equal to zero. When only considering nonzero singular values, as we do, $A = U \Sigma V^T$ is the appropriate definition of the SVD. This is sometimes called the compact SVD of A , but we will just call it the SVD.

Example: Let's use the results of the previous examples to construct the SVD of

$$A = \begin{bmatrix} 4 & 11 & 14 \\ 8 & 7 & -2 \end{bmatrix}.$$

Step 1: Find an orthogonal diagonalization of $A^T A$. In general, for A with many columns, this is done numerically; here we use the data from before.

$$A^T A = \hat{V} \hat{\Lambda} \hat{V}^T = [v_1 \ v_2 \ v_3] \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_3 \end{bmatrix} \begin{bmatrix} v_1^T \\ v_2^T \\ v_3^T \end{bmatrix}$$

with (λ_i, v_i) as specified above. with $\lambda_1 = 360$, $\lambda_2 = 90$, $\lambda_3 = 0$.

Step 2: Setup V and Σ : Arrange the nonzero eigenvalues of $A^T A$ in decreasing order and compute the singular values. For this example:

$$G_1 = 6\sqrt{10} \quad \text{and} \quad G_2 = 3\sqrt{10},$$

and

$$\Sigma = \text{diag}(G_1, G_2) = \begin{bmatrix} 6\sqrt{10} & 0 \\ 0 & 3\sqrt{10} \end{bmatrix}. \text{ Hence } \text{rank } A = 2, \text{ and } V \in \mathbb{R}^{3 \times 2}.$$

The corresponding eigenvectors define the columns of V :

$$V = [v_1 \ v_2] = \begin{bmatrix} 1/3 & -2/3 \\ 2/3 & -1/3 \\ 2/3 & 2/3 \end{bmatrix}.$$

Step 3: Construct U : Since $\text{rank } A = 2$, $U \in \mathbb{R}^{2 \times 2}$. The columns of U are given by the normalized vectors obtained from Av_1 and Av_2 . Recall that we showed above that $\|Av_1\| = G_1$ and $\|Av_2\| = G_2$, so $U = [u_1 \ u_2]$ with

$$u_1 = \frac{Av_1}{G_1} = \frac{1}{6\sqrt{10}} \begin{bmatrix} 18 \\ 6 \end{bmatrix} = \begin{bmatrix} 3/\sqrt{10} \\ 1/\sqrt{10} \end{bmatrix} \quad \text{and}$$

$$u_2 = \frac{Av_2}{G_2} = \frac{1}{3\sqrt{10}} \begin{bmatrix} 3 \\ -9 \end{bmatrix} = \begin{bmatrix} 1/\sqrt{10} \\ -3/\sqrt{10} \end{bmatrix}.$$

Finally, the SVD of A is:

$$A = \underbrace{\begin{bmatrix} 3/\sqrt{10} & 1/\sqrt{10} \\ 1/\sqrt{10} & -3/\sqrt{10} \end{bmatrix}}_U \underbrace{\begin{bmatrix} 6\sqrt{10} & 0 \\ 0 & 3\sqrt{10} \end{bmatrix}}_\Sigma \underbrace{\begin{bmatrix} 1/3 & 2/3 & 2/3 \\ -2/3 & -1/3 & 2/3 \end{bmatrix}}_{V^T}.$$

You can check that indeed $A = U \Sigma V^T$ here, and that $U^T U = V^T V = I_2$.

ONLINE NOTES: Please add example 4 from LAA 7.4, suitably modified to use the compact SVD form.

Linear Algebra Applications of the SVD

The next few classes will focus on engineering and AI applications of the SVD. For now, we highlight some more technical linear algebraic applications: these are all immensely important from a practical perspective, and form subroutines for most real-world applications of linear algebra.

The Condition Number Most numerical calculations that require solving a linear equation $Ax=b$ are as reliable as possible when the SVD of A is used. Since the matrices U and V have orthonormal columns, they do not affect lengths or vectors between vectors. For example, for $U \in \mathbb{R}^{m \times m}$, we have:

$$\langle Ux, Uy \rangle = x^T U^T U y = x^T y = \langle x, y \rangle$$

for any $x, y \in \mathbb{R}^n$. Therefore, any numerical issues that arise will be due to the diagonal entries of Σ , i.e., due to the singular values of A .

In particular, if some of singular values are much larger than others, this means certain directions are stretched out much more than others, which can lead to roundoff errors. A natural way to quantify this notion is using the singular values of A . If $A \in \mathbb{R}^{m \times n}$ is an $n \times n$ invertible matrix, so that $r = \text{rank}(A) = n$, we define the Condition Number of A to be the ratio $\kappa(A) = \frac{\sigma_1}{\sigma_n}$ of the largest to smallest singular values of A . If A is not invertible, it is on convention to set $\kappa(A) = \infty$, although the ratio σ_1/σ_r is a useful measure of the numerical stability of computing with a rectangular matrix $A \in \mathbb{R}^{m \times n}$.

ONLINE NOTES: Using numpy, show that ill-conditioned A can lead to bad solutions to $Ax=b$ even if A is invertible. Case 1: assume we use $\tilde{b} = b + n$ for n some small measurement noise. Case 2: just make A super ill conditional and show that $\tilde{x} = A^{-1}\tilde{b}$ computed using numpy doesn't actually satisfy $A\tilde{x} = \tilde{b}$.

Computing Bases of Fundamental Subspaces

Given an SVD for an $m \times n$ matrix $A \in \mathbb{R}^{m \times n}$ with $\text{rank}(A) = r$, let u_1, \dots, u_r be the left singular vectors, v_1, \dots, v_r the right singular vectors, and $\sigma_1, \dots, \sigma_r$ the singular values.

Recall that we showed that u_1, \dots, u_r forms a basis for $\text{Col}(A)$. Let u_{r+1}, \dots, u_m be an orthonormal basis for $\text{Col}(A)^\perp$ so that u_1, \dots, u_m form a basis for \mathbb{R}^m , completed for example using the Gram-Schmidt Process. Then, by the FTLA, we have that $\text{Col}(A)^\perp = \text{Null}(A^T) = \text{span}\{u_{r+1}, \dots, u_m\}$, i.e., these vectors form an orthonormal basis for $\text{Null}(A^T) = \text{LNull}(A)$.

Next, recall that $v_1, \dots, v_r, v_{r+1}, \dots, v_n$, the eigenvectors of $A^T A$, form an orthonormal

basis of \mathbb{R}^n . Since $Av_i = 0$ for $i=r+1, \dots, n$, the vectors v_{r+1}, \dots, v_n span a subspace of $\text{Null}(A)$ of dimension $n-r$. But, by the FTLA, $\dim \text{Null}(A) = n - \text{rank}(A) = n-r$.

Therefore, v_{r+1}, \dots, v_n are an orthonormal basis for $\text{Null}(A)$.

Finally, $\text{Null}(A)^\perp = \text{Col}(A^T) = \text{Row}(A)$. But $\text{Null}(A)^\perp = \text{span}\{v_1, \dots, v_r\}$ since the v_i are an orthonormal basis for \mathbb{R}^n , and thus v_1, \dots, v_r are an orthonormal basis for $\text{Row}(A)$.

Summarizing, we have:

- $\text{Col}(A) = \text{span}\{u_1, \dots, u_r\}$
- $\text{Col}(A)^\perp = \text{Null}(A^T) = \text{LNull}(A) = \text{span}\{u_{r+1}, \dots, u_m\}$
- $\text{Col}(A^T) = \text{Row}(A) = \text{span}\{v_1, \dots, v_r\}$
- $\text{Col}(A^T)^\perp = \text{Null}(A) = \text{span}\{v_{r+1}, \dots, v_n\}$

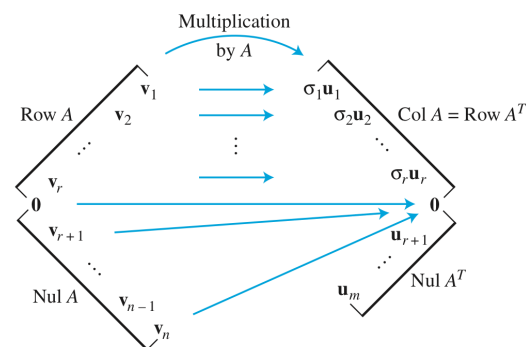


FIGURE 4 The four fundamental subspaces and the action of A .

Specializing these observations to square matrices, we have the following theorem characterizing invertible matrices:

Theorem: The following statements are equivalent for a square $n \times n$ matrix $A \in \mathbb{R}^{n \times n}$:

- $\text{Col}(A)^\perp = \text{Null}(A^T) = \{0\}$
- $\text{Null}(A)^\perp = \text{Col}(A^T) = \text{Row}(A) = \mathbb{R}^n$
- $\text{Col}(A^T)^\perp = \text{Null}(A) = \{0\}$
- $\text{Null}(A^T)^\perp = \text{Col}(A) = \mathbb{R}^n$
- A has $\text{rank} = n$
- A has n (nonzero) singular values

The Pseudoinverse of A

Recall the least squares problem of finding a vector \underline{x} that minimizes the objective $\|A\underline{x} - \underline{b}\|^2$. We saw that the least squares solution is given by the solution to the normal equations

$$A^T A \underline{x} = A^T \underline{b}. \quad (\text{NE})$$

Let's rewrite (NE) using the SVD $A = U \Sigma V^T$, $A^T = V \Sigma^T U^T = V \Sigma U^T$ ($\Sigma = \Sigma^T$)

$$A^T A \underline{x} = \underbrace{V}_{(a)} \underbrace{\Sigma^T U^T U \Sigma}_{(c)} \underbrace{V^T}_{(b)} \underline{x} = \underbrace{V}_{(a)} \underbrace{\Sigma^2}_{(c)} \underbrace{V^T}_{(b)} \underline{x} = \underbrace{V}_{(a)} \underbrace{\Sigma}_{(c)} \underbrace{U^T}_{(b)} \underline{b}$$

Let's start by left multiplying (a) and (b) by V^T to take advantage of $V^T V = I$.

$$V^T (V \Sigma^2 V^T \underline{x}) = V^T (V \Sigma U^T \underline{b}) \Rightarrow \Sigma^2 V^T \underline{x} = \Sigma U^T \underline{b}.$$

Now, let's isolate $V^T \underline{x}$ by multiplying both sides by Σ^{-2} :

$$V^T \underline{x} = \Sigma^{-1} U^T \underline{b}. \quad (*)$$

Finally, note that \underline{x} satisfies star if $\hat{\underline{x}} = V \Sigma^{-1} U^T \underline{b}$ (again since $V^T V = I$)

for any $\underline{n} \in \text{Null}(V^T) = \text{Col}(V)^\perp$. The special solution $\underline{x}^* = V \Sigma^{-1} U^T \underline{b}$ can be shown to be the minimum norm least squares solution when several \underline{x} exist such that $A \underline{x} = \underline{b}$. The matrix

$$A^+ = V \Sigma^{-1} U^T$$

is called the **pseudoinverse of A**, and is also known as the **Moore-Penrose Inverse of A**.

If we look at $A \underline{x}^* = A A^+ \underline{b}$, we observe that:

$$A \underline{x}^* = U \underbrace{\Sigma V^T V}_{I} \Sigma^{-1} U^T \underline{b} = U U^T \underline{b},$$

i.e., $A \underline{x}^*$ is the orthogonal projection $\hat{\underline{b}}$ of \underline{b} onto $\text{Col}(A)$.

ONLINE NOTES: Please add the Practice Problems on p.471 of LAA at the end of Ch. 7.4 + Solutions