

Linear Algebra

Garrett Thomas

May 29, 2018

1 About

This document is part of a series of notes about math and machine learning. You are free to distribute it as you wish. The latest version can be found at <http://gwthomas.github.io/notes>. Please report any errors to gwthomas@stanford.edu.

1.1 Vector spaces

Vector spaces are the basic setting in which linear algebra happens. A vector space over a field \mathbb{F} consists of a set V (the elements of which are called **vectors**) on which two operations are defined: vectors in V can be added together, and they can be multiplied by **scalars** from \mathbb{F} . For simplicity we will assume $\mathbb{F} = \mathbb{R}$ or $\mathbb{F} = \mathbb{C}$, although in general it could be any field. A vector space must satisfy

- (i) There exists an additive identity (denoted 0) in V such that $v + 0 = v$ for all $v \in V$
- (ii) For each $v \in V$, there exists an additive inverse (denoted $-v$) such that $v + (-v) = 0$
- (iii) There exists a multiplicative identity (denoted 1) in \mathbb{F} such that $1v = v$ for all $v \in V$
- (iv) Commutativity: $u + v = v + u$ for all $u, v \in V$
- (v) Associativity: $(u + v) + w = u + (v + w)$ and $\alpha(\beta v) = (\alpha\beta)v$ for all $u, v, w \in V$ and $\alpha, \beta \in \mathbb{F}$
- (vi) Distributivity: $\alpha(u + v) = \alpha u + \alpha v$ and $(\alpha + \beta)v = \alpha v + \beta v$ for all $u, v \in V$ and $\alpha, \beta \in \mathbb{F}$

A set of vectors $v_1, \dots, v_n \in V$ is said to be **linearly independent** if

$$\alpha_1 v_1 + \dots + \alpha_n v_n = 0 \quad \text{implies} \quad \alpha_1 = \dots = \alpha_n = 0.$$

The **span** of $v_1, \dots, v_n \in V$ is the set of all vectors that can be expressed of a **linear combination** of them:

$$\text{span}\{v_1, \dots, v_n\} \triangleq \{\alpha_1 v_1 + \dots + \alpha_n v_n : \alpha_1, \dots, \alpha_n \in \mathbb{F}\}$$

If a set of vectors is linearly independent and its span is the whole of V , those vectors are said to be a **basis** for V . In fact, every linearly independent set of vectors forms a basis for its span.

If a vector space is spanned by a finite number of vectors, it is said to be **finite-dimensional**. Otherwise it is **infinite-dimensional**. The number of vectors in a basis for a finite-dimensional vector space V is called the **dimension** of V and denoted $\dim V$.

1.1.1 Euclidean space

The quintessential vector space is **Euclidean space**, which we denote \mathbb{R}^n . The vectors in this space consist of n -tuples of real numbers:

$$x = (x_1, x_2, \dots, x_n)$$

For our purposes, it will be useful to think of them as $n \times 1$ matrices, or **column vectors**:

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Addition and scalar multiplication are defined component-wise on vectors in \mathbb{R}^n :

$$x + y = \begin{bmatrix} x_1 + y_1 \\ \vdots \\ x_n + y_n \end{bmatrix}, \quad \alpha x = \begin{bmatrix} \alpha x_1 \\ \vdots \\ \alpha x_n \end{bmatrix}$$

Euclidean space is used to mathematically represent physical space, with notions such as distance, length, and angles. Although it becomes hard to visualize for $n > 3$, these concepts generalize mathematically in obvious ways. Even when you're working in more general settings than \mathbb{R}^n , it is often useful to visualize vector addition and scalar multiplication in terms of 2D vectors in the plane or 3D vectors in space.

1.1.2 Subspaces

Vector spaces can contain other vector spaces. If V is a vector space, then $S \subseteq V$ is said to be a **subspace** of V if

- (i) $0 \in S$
- (ii) S is closed under addition: $u, v \in S$ implies $u + v \in S$
- (iii) S is closed under scalar multiplication: $v \in S, \alpha \in \mathbb{F}$ implies $\alpha v \in S$

Note that V is always a subspace of V , as is the trivial vector space which contains only 0.

As a concrete example, a line passing through the origin is a subspace of Euclidean space.

If U and W are subspaces of V , then their sum is defined as

$$U + W \triangleq \{u + w : u \in U, w \in W\}$$

It is straightforward to verify that this set is also a subspace of V . If $U \cap W = \{0\}$, the sum is said to be a **direct sum** and written $U \oplus W$. Every vector in $U \oplus W$ can be written uniquely as $u + w$ for some $u \in U$ and $w \in W$. (This is both a necessary and sufficient condition for a direct sum.)

The dimensions of sums of subspaces obey a friendly relationship:

$$\dim(U + W) = \dim U + \dim W - \dim(U \cap W)$$

It follows that

$$\dim(U \oplus W) = \dim U + \dim W$$

since $\dim(U \cap W) = \dim(\{0\}) = 0$ if the sum is direct.

1.2 Linear maps

A **linear map** is a function $T : V \rightarrow W$, where V and W are vector spaces, that satisfies

- (i) $T(u + v) = Tu + Tv$ for all $u, v \in V$
- (ii) $T(\alpha v) = \alpha Tv$ for all $v \in V, \alpha \in \mathbb{F}$

The standard notational convention for linear maps (which we follow here) is to drop unnecessary parentheses, writing Tv rather than $T(v)$ if there is no risk of ambiguity, and denote composition of linear maps by ST rather than the usual $S \circ T$.

Observe that the definition of a linear map is suited to reflect the structure of vector spaces, since it preserves vector spaces' two main operations, addition and scalar multiplication. In algebraic terms, a linear map is called a **homomorphism** of vector spaces. An invertible homomorphism (where the inverse is also a homomorphism) is called an **isomorphism**. If there exists an isomorphism from V to W , then V and W are said to be **isomorphic**, and we write $V \cong W$. Isomorphic vector spaces are essentially “the same” in terms of their algebraic structure. It is an interesting fact that finite-dimensional vector spaces of the same dimension over the same field are always isomorphic; if V, W are real vector spaces with $\dim V = \dim W = n$, then we have the natural isomorphism

$$\begin{aligned} \varphi : V &\rightarrow W \\ \alpha_1 v_1 + \cdots + \alpha_n v_n &\mapsto \alpha_1 w_1 + \cdots + \alpha_n w_n \end{aligned}$$

where v_1, \dots, v_n and w_1, \dots, w_n are any bases for V and W . This map is well-defined because every vector in V can be expressed uniquely as a linear combination of v_1, \dots, v_n . It is straightforward to verify that φ is an isomorphism, so in fact $V \cong W$. In particular, every real n -dimensional vector space is isomorphic to \mathbb{R}^n .

1.2.1 The matrix of a linear map

Vector spaces are fairly abstract. To represent and manipulate vectors and linear maps on a computer, we use rectangular arrays of numbers known as **matrices**.

Suppose V and W are finite-dimensional vector spaces with bases v_1, \dots, v_n and w_1, \dots, w_m , respectively, and $T : V \rightarrow W$ is a linear map. Then the matrix of T , with entries A_{ij} where $i = 1, \dots, m$, $j = 1, \dots, n$, is defined by

$$Tv_j = A_{1j}w_1 + \cdots + A_{mj}w_m$$

That is, the j th column of A consists of the coordinates of Tv_j in the chosen basis for W .

Conversely, every matrix $A \in \mathbb{R}^{m \times n}$ induces a linear map $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ given by

$$Tx = Ax$$

and the matrix of this map with respect to the standard bases of \mathbb{R}^n and \mathbb{R}^m is of course simply A .

If $A \in \mathbb{R}^{m \times n}$, its **transpose** $A^\top \in \mathbb{R}^{n \times m}$ is given by $(A^\top)_{ij} = A_{ji}$ for each (i, j) . In other words, the columns of A become the rows of A^\top , and the rows of A become the columns of A^\top .

The transpose has several nice algebraic properties that can be easily verified from the definition:

- (i) $(A^\top)^\top = A$
- (ii) $(A + B)^\top = A^\top + B^\top$
- (iii) $(\alpha A)^\top = \alpha A^\top$
- (iv) $(AB)^\top = B^\top A^\top$

1.2.2 Nullspace, range

Some of the most important subspaces are those induced by linear maps. If $T : V \rightarrow W$ is a linear map, the **nullspace**¹ of T is defined as

$$\text{null}(T) \triangleq \{v \in V : Tv = 0\}$$

and the **range** of T as

$$\text{range}(T) \triangleq \{w \in W : Tv = w \text{ for some } v \in V\}$$

It is a good exercise to verify that the nullspace and range of a linear map are always subspaces of its domain and codomain, respectively.

The **columnspace** of a matrix $A \in \mathbb{R}^{m \times n}$ is the span of its columns (considered as vectors in \mathbb{R}^m), and similarly the **rowspace** of A is the span of its rows (considered as vectors in \mathbb{R}^n). It is not hard to see that the columnspace of A is exactly the range of the linear map from \mathbb{R}^n to \mathbb{R}^m which is induced by A , so we denote it by $\text{range}(A)$ in a slight abuse of notation. Similarly, the rowspace is denoted $\text{range}(A^\top)$.

It is a remarkable fact that the dimension of the columnspace of A is the same as the dimension of the rowspace of A . This quantity is called the **rank** of A , and defined as

$$\text{rank}(A) \triangleq \dim \text{range}(A)$$

1.3 Normed spaces

Norms generalize the notion of length from Euclidean space.

A **norm** on a vector space V is a function $\|\cdot\| : V \rightarrow \mathbb{R}$ that satisfies

- (i) $\|v\| \geq 0$, with equality if and only if $x = 0$
- (ii) $\|\alpha v\| = |\alpha| \|v\|$
- (iii) $\|u + v\| \leq \|u\| + \|v\|$ (the **triangle inequality**)

for all $u, v \in V$ and all $\alpha \in \mathbb{F}$. A vector space endowed with a norm is called a **normed vector space**, or simply a **normed space**.

Note that any norm on V induces a distance metric on V :

$$d(u, v) \triangleq \|u - v\|$$

One can verify that the axioms for metrics are satisfied under this definition and follow directly from the axioms for norms. Therefore any normed space is also a metric space. If a normed space is complete with respect to the distance metric induced by its norm, it is said to be a **Banach space**.

¹ It is sometimes called the **kernel** by algebraists, but we eschew this terminology because the word “kernel” has another meaning in machine learning.

We will typically only be concerned with a few specific norms on \mathbb{R}^n :

$$\begin{aligned}\|x\|_1 &= \sum_{i=1}^n |x_i| \\ \|x\|_2 &= \sqrt{\sum_{i=1}^n x_i^2} \\ \|x\|_p &= \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} \quad (p \geq 1) \\ \|x\|_\infty &= \max_{1 \leq i \leq n} |x_i|\end{aligned}$$

Note that the 1- and 2-norms are special cases of the p -norm, and the ∞ -norm is the limit of the p -norm as p tends to infinity. We require $p \geq 1$ for the general definition of the p -norm because the triangle inequality fails to hold if $p < 1$. (Try to find a counterexample!)

Here's a fun fact: for any given finite-dimensional vector space V , all norms on V are equivalent in the sense that for two norms $\|\cdot\|_A, \|\cdot\|_B$, there exist constants $\alpha, \beta > 0$ such that

$$\alpha\|v\|_A \leq \|v\|_B \leq \beta\|v\|_A$$

for all $v \in V$. Therefore convergence in one norm implies convergence in any other norm. This rule may not apply in infinite-dimensional vector spaces such as function spaces, though.

1.4 Inner product spaces

An **inner product** on a vector space V is a function $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{F}$ satisfying

- (i) $\langle v, v \rangle \geq 0$, with equality if and only if $v = 0$
- (ii) Linearity in the first slot: $\langle u + v, w \rangle = \langle u, w \rangle + \langle v, w \rangle$ and $\langle \alpha u, v \rangle = \alpha \langle u, v \rangle$
- (iii) **Conjugate symmetry**: $\langle u, v \rangle = \overline{\langle v, u \rangle}$

for all $u, v, w \in V$ and all $\alpha \in \mathbb{F}$. A vector space endowed with an inner product is called an **inner product space**.

Note that any inner product on V induces a norm on V :

$$\|v\| \triangleq \sqrt{\langle v, v \rangle}$$

One can verify that the axioms for norms are satisfied under this definition and follow (almost) directly from the axioms for inner products. Therefore any inner product space is also a normed space (and hence also a metric space). If an inner product space is complete with respect to the distance metric induced by its inner product, we say that it is a **Hilbert space**.

Two vectors u and v are said to be **orthogonal** if $\langle u, v \rangle = 0$; we write $u \perp v$ for shorthand. Orthogonality generalizes the notion of perpendicularity from Euclidean space. If two orthogonal vectors u and v additionally have unit length (i.e. $\|u\| = \|v\| = 1$), then they are described as **orthonormal**.

The standard inner product on \mathbb{R}^n is given by

$$\langle x, y \rangle \triangleq \sum_{i=1}^n x_i y_i = x^\top y$$

The matrix notation on the righthand side arises because this inner product is a special case of matrix multiplication where we regard the resulting 1×1 matrix as a scalar. The inner product on \mathbb{R}^n is also often written $x \cdot y$ (hence the alternate name **dot product**). The reader can verify that the two-norm $\|\cdot\|_2$ on \mathbb{R}^n is induced by this inner product.

1.4.1 Pythagorean Theorem

The well-known Pythagorean theorem generalizes naturally to arbitrary inner product spaces.

Theorem 1. *If $u \perp v$, then*

$$\|u + v\|^2 = \|u\|^2 + \|v\|^2$$

Proof. Suppose $u \perp v$, i.e. $\langle u, v \rangle = 0$. Then

$$\|u + v\|^2 = \langle u + v, u + v \rangle = \langle u, u \rangle + \langle v, u \rangle + \langle u, v \rangle + \langle v, v \rangle = \|u\|^2 + \|v\|^2$$

as claimed. □

1.4.2 Cauchy-Schwarz inequality

This inequality is sometimes useful in proving bounds:

$$|\langle u, v \rangle| \leq \|u\| \|v\|$$

for all $u, v \in V$. Equality holds exactly when u and v are scalar multiples of each other (or equivalently, when they are linearly dependent).

1.4.3 Orthogonal complements and projections

If $S \subseteq V$ where V is an inner product space, then the **orthogonal complement** of S , denoted S^\perp , is the set of all vectors in V that are orthogonal to every element of S :

$$S^\perp = \{v \in V : v \perp s \text{ for all } s \in S\}$$

It is easy to verify that S^\perp is a subspace of V for any $S \subseteq V$. Note that there is no requirement that S itself be a subspace of V . However, if S is a (finite-dimensional) subspace of V , we have the following important decomposition.

Proposition 1. *Let V be an inner product space and S be a finite-dimensional subspace of V . Then every $v \in V$ can be written uniquely in the form*

$$v = v_S + v_\perp$$

where $v_S \in S$ and $v_\perp \in S^\perp$.

Proof. Let u_1, \dots, u_m be an orthonormal basis for S , and suppose $v \in V$. Define

$$v_S = \langle v, u_1 \rangle u_1 + \dots + \langle v, u_m \rangle u_m$$

and

$$v_\perp = v - v_S$$

It is clear that $v_S \in S$ since it is in the span of the chosen basis. We also have, for all $i = 1, \dots, m$,

$$\begin{aligned}\langle v_\perp, u_i \rangle &= \langle v - (\langle v, u_1 \rangle u_1 + \dots + \langle v, u_m \rangle u_m), u_i \rangle \\ &= \langle v, u_i \rangle - \langle v, u_1 \rangle \langle u_1, u_i \rangle - \dots - \langle v, u_m \rangle \langle u_m, u_i \rangle \\ &= \langle v, u_i \rangle - \langle v, u_i \rangle \\ &= 0\end{aligned}$$

which implies $v_\perp \in S^\perp$.

It remains to show that this decomposition is unique, i.e. doesn't depend on the choice of basis. To this end, let u'_1, \dots, u'_m be another orthonormal basis for S , and define v'_S and v'_\perp analogously. We claim that $v'_S = v_S$ and $v'_\perp = v_\perp$.

By definition,

$$v_S + v_\perp = v = v'_S + v'_\perp$$

so

$$\underbrace{v_S - v'_S}_{\in S} = \underbrace{v'_\perp - v_\perp}_{\in S^\perp}$$

From the orthogonality of these subspaces, we have

$$0 = \langle v_S - v'_S, v'_\perp - v_\perp \rangle = \langle v_S - v'_S, v_S - v'_S \rangle = \|v_S - v'_S\|^2$$

It follows that $v_S - v'_S = 0$, i.e. $v_S = v'_S$. Then $v'_\perp = v - v'_S = v - v_S = v_\perp$ as well. \square

The existence and uniqueness of the decomposition above mean that

$$V = S \oplus S^\perp$$

whenever S is a subspace.

Since the mapping from v to v_S in the decomposition above always exists and is unique, we have a well-defined function

$$\begin{aligned}P_S : V &\rightarrow S \\ v &\mapsto v_S\end{aligned}$$

which is called the **orthogonal projection** onto S . We give the most important properties of this function below.

Proposition 2. *Let S be a finite-dimensional subspace of V . Then*

(i) *For any $v \in V$ and orthonormal basis u_1, \dots, u_m of S ,*

$$P_S v = \langle v, u_1 \rangle u_1 + \dots + \langle v, u_m \rangle u_m$$

(ii) *For any $v \in V$, $v - P_S v \perp S$.*

(iii) *P_S is a linear map.*

(iv) *P_S is the identity when restricted to S (i.e. $P_S s = s$ for all $s \in S$).*

(v) *$\text{range}(P_S) = S$ and $\text{null}(P_S) = S^\perp$.*

(vi) *$P_S^2 = P_S$.*

(vii) For any $v \in V$, $\|P_S v\| \leq \|v\|$.

(viii) For any $v \in V$ and $s \in S$,

$$\|v - P_S v\| \leq \|v - s\|$$

with equality if and only if $s = P_S v$. That is,

$$P_S v = \arg \min_{s \in S} \|v - s\|$$

Proof. The first two statements are immediate from the definition of P_S and the work done in the proof of the previous proposition.

In this proof, we abbreviate $P = P_S$ for brevity.

(iii) Suppose $u, v \in V$ and $\alpha \in \mathbb{R}$. Write $u = u_S + u_\perp$ and $v = v_S + v_\perp$, where $u_S, v_S \in S$ and $u_\perp, v_\perp \in S^\perp$. Then

$$u + v = \underbrace{u_S + v_S}_{\in S} + \underbrace{u_\perp + v_\perp}_{\in S^\perp}$$

so $P(u + v) = u_S + v_S = Pu + Pv$, and

$$\alpha v = \underbrace{\alpha v_S}_{\in S} + \underbrace{\alpha v_\perp}_{\in S^\perp}$$

so $P(\alpha v) = \alpha v_S = \alpha Pv$. Thus P is linear.

(iv) If $s \in S$, then we can write $s = s + 0$ where $s \in S$ and $0 \in S^\perp$, so $Ps = s$.

(v) $\text{range}(P) \subseteq S$: By definition.

$\text{range}(P) \supseteq S$: Using the previous result, any $s \in S$ satisfies $s = Pv$ for some $v \in V$ (specifically, $v = s$).

$\text{null}(P) \subseteq S^\perp$: Suppose $v \in \text{null}(P)$. Write $v = v_S + v_\perp$ where $v_S \in S$ and $v_\perp \in S^\perp$. Then $0 = Pv = v_S$, so $v = v_\perp \in S^\perp$.

$\text{null}(P) \supseteq S^\perp$: If $v \in S^\perp$, then $v = 0 + v$ where $0 \in S$ and $v \in S^\perp$, so $Pv = 0$.

(vi) For any $v \in V$,

$$P^2 v = P(Pv) = Pv$$

since $Pv \in S$ and P is the identity on S . Hence $P^2 = P$.

(vii) Suppose $v \in V$. Then by the Pythagorean theorem,

$$\|v\|^2 = \|Pv + (v - Pv)\|^2 = \|Pv\|^2 + \|v - Pv\|^2 \geq \|Pv\|^2$$

The result follows by taking square roots.

(viii) Suppose $v \in V$ and $s \in S$. Then by the Pythagorean theorem,-

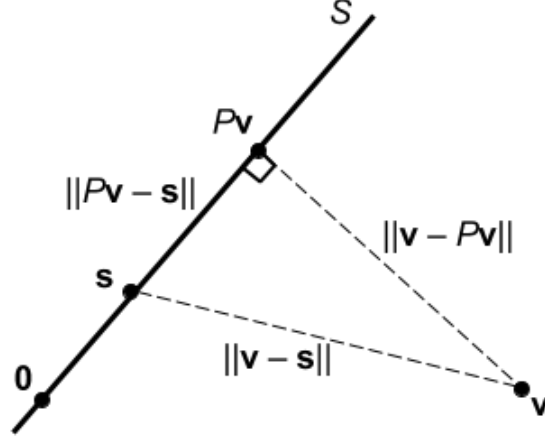
$$\|v - s\|^2 = \|(v - Pv) + (Pv - s)\|^2 = \|v - Pv\|^2 + \|Pv - s\|^2 \geq \|v - Pv\|^2$$

We obtain $\|v - s\| \geq \|v - Pv\|$ by taking square roots. Equality holds iff $\|Pv - s\|^2 = 0$, which is true iff $Pv = s$.

□

Any linear map P that satisfies $P^2 = P$ is called a **projection**, so we have shown that P_S is a projection (hence the name).

The last part of the previous result shows that orthogonal projection solves the optimization problem of finding the closest point in S to a given $v \in V$. This makes intuitive sense from a pictorial representation of the orthogonal projection:



Let us now consider the specific case where S is a subspace of \mathbb{R}^n with orthonormal basis u_1, \dots, u_m . Then

$$P_S x = \sum_{i=1}^m \langle x, u_i \rangle u_i = \sum_{i=1}^m x^\top u_i u_i = \sum_{i=1}^m u_i u_i^\top x = \left(\sum_{i=1}^m u_i u_i^\top \right) x$$

so the operator P_S can be expressed as a matrix

$$P_S = \sum_{i=1}^m u_i u_i^\top = U U^\top$$

where U has u_1, \dots, u_m as its columns. Here we have used the sum-of-outer-products identity.

1.5 Eigenthings

For a square matrix $A \in \mathbb{R}^{n \times n}$, there may be vectors which, when A is applied to them, are simply scaled by some constant. We say that a nonzero vector $x \in \mathbb{R}^n$ is an **eigenvector** of A corresponding to **eigenvalue** λ if

$$Ax = \lambda x$$

The zero vector is excluded from this definition because $A0 = 0 = \lambda 0$ for every λ .

We now give some useful results about how eigenvalues change after various manipulations.

Proposition 3. *Let x be an eigenvector of A with corresponding eigenvalue λ . Then*

- (i) *For any $\gamma \in \mathbb{R}$, x is an eigenvector of $A + \gamma I$ with eigenvalue $\lambda + \gamma$.*
- (ii) *If A is invertible, then x is an eigenvector of A^{-1} with eigenvalue λ^{-1} .*
- (iii) *$A^k x = \lambda^k x$ for any $k \in \mathbb{Z}$ (where $A^0 = I$ by definition).*

Proof. (i) follows readily:

$$(A + \gamma I)x = Ax + \gamma Ix = \lambda x + \gamma x = (\lambda + \gamma)x$$

(ii) Suppose A is invertible. Then

$$x = A^{-1}Ax = A^{-1}(\lambda x) = \lambda A^{-1}x$$

Dividing by λ , which is valid because the invertibility of A implies $\lambda \neq 0$, gives $\lambda^{-1}x = A^{-1}x$.

(iii) The case $k \geq 0$ follows immediately by induction on k . Then the general case $k \in \mathbb{Z}$ follows by combining the $k \geq 0$ case with (ii). \square

1.6 Trace

The **trace** of a square matrix is the sum of its diagonal entries:

$$\text{tr}(A) = \sum_{i=1}^n A_{ii}$$

The trace has several nice algebraic properties:

- (i) $\text{tr}(A + B) = \text{tr}(A) + \text{tr}(B)$
- (ii) $\text{tr}(\alpha A) = \alpha \text{tr}(A)$
- (iii) $\text{tr}(A^\top) = \text{tr}(A)$
- (iv) $\text{tr}(AB) = \text{tr}(BA)$

The first three properties follow readily from the definition. The last is known as **invariance under cyclic permutations**.

Interestingly, the trace of a matrix is equal to the sum of its eigenvalues (repeated according to multiplicity):

$$\text{tr}(A) = \sum_i \lambda_i(A)$$

1.7 Determinant

The **determinant** of a square matrix can be defined in several different confusing ways, none of which are particularly important for our purposes; go look at an introductory linear algebra text (or Wikipedia) if you need a definition. But it's good to know the properties:

- (i) $\det(I) = 1$
- (ii) $\det(A^\top) = \det(A)$
- (iii) $\det(AB) = \det(A) \det(B)$
- (iv) $\det(A^{-1}) = \det(A)^{-1}$
- (v) $\det(\alpha A) = \alpha^n \det(A)$

Interestingly, the determinant of a matrix is equal to the product of its eigenvalues (repeated according to multiplicity):

$$\det(A) = \prod_i \lambda_i(A)$$

1.8 Orthogonal matrices

A matrix $Q \in \mathbb{R}^{n \times n}$ is said to be **orthogonal** if its columns are pairwise orthonormal. This definition implies that

$$Q^\top Q = QQ^\top = I$$

or equivalently, $Q^\top = Q^{-1}$. A nice thing about orthogonal matrices is that they preserve inner products:

$$(Qx)^\top(Qy) = x^\top Q^\top Qy = x^\top Iy = x^\top y$$

A direct result of this fact is that they also preserve 2-norms:

$$\|Qx\|_2 = \sqrt{(Qx)^\top(Qx)} = \sqrt{x^\top x} = \|x\|_2$$

Therefore multiplication by an orthogonal matrix can be considered as a transformation that preserves length, but may rotate or reflect the vector about the origin.

1.9 Symmetric matrices

A matrix $A \in \mathbb{R}^{n \times n}$ is said to be **symmetric** if it is equal to its own transpose ($A = A^\top$), meaning that $A_{ij} = A_{ji}$ for all (i, j) . This definition seems harmless enough but turns out to have some strong implications. We summarize the most important of these as

Theorem 2. (*Spectral Theorem*) *If $A \in \mathbb{R}^{n \times n}$ is symmetric, then there exists an orthonormal basis for \mathbb{R}^n consisting of eigenvectors of A .*

The practical application of this theorem is a particular factorization of symmetric matrices, referred to as the **eigendecomposition** or **spectral decomposition**. Denote the orthonormal basis of eigenvectors q_1, \dots, q_n and their eigenvalues $\lambda_1, \dots, \lambda_n$. Let Q be an orthogonal matrix with q_1, \dots, q_n as its columns, and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. Since by definition $Aq_i = \lambda_i q_i$ for every i , the following relationship holds:

$$AQ = Q\Lambda$$

Right-multiplying by Q^\top , we arrive at the decomposition

$$A = Q\Lambda Q^\top$$

1.9.1 Rayleigh quotients

Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix. The expression $x^\top Ax$ is called a **quadratic form**.

There turns out to be an interesting connection between the quadratic form of a symmetric matrix and its eigenvalues. This connection is provided by the **Rayleigh quotient**

$$R_A(x) = \frac{x^\top Ax}{x^\top x}$$

The Rayleigh quotient has a couple of important properties which the reader can (and should!) easily verify from the definition:

- (i) **Scale invariance:** for any vector $x \neq 0$ and any scalar $\alpha \neq 0$, $R_A(x) = R_A(\alpha x)$.
- (ii) If x is an eigenvector of A with eigenvalue λ , then $R_A(x) = \lambda$.

We can further show that the Rayleigh quotient is bounded by the largest and smallest eigenvalues of A . But first we will show a useful special case of the final result.

Proposition 4. *For any x such that $\|x\|_2 = 1$,*

$$\lambda_{\min}(A) \leq x^\top A x \leq \lambda_{\max}(A)$$

with equality if and only if x is a corresponding eigenvector.

Proof. We show only the max case because the argument for the min case is entirely analogous.

Since A is symmetric, we can decompose it as $A = Q\Lambda Q^\top$. Then use the change of variable $y = Q^\top x$, noting that the relationship between x and y is one-to-one and that $\|y\|_2 = 1$ since Q is orthogonal. Hence

$$\max_{\|x\|_2=1} x^\top A x = \max_{\|y\|_2=1} y^\top \Lambda y = \max_{y_1^2 + \dots + y_n^2 = 1} \sum_{i=1}^n \lambda_i y_i^2$$

Written this way, it is clear that y maximizes this expression exactly if and only if it satisfies $\sum_{i \in I} y_i^2 = 1$ where $I = \{i : \lambda_i = \max_{j=1, \dots, n} \lambda_j = \lambda_{\max}(A)\}$ and $y_j = 0$ for $j \notin I$. That is, I contains the index or indices of the largest eigenvalue. In this case, the maximal value of the expression is

$$\sum_{i=1}^n \lambda_i y_i^2 = \sum_{i \in I} \lambda_i y_i^2 = \lambda_{\max}(A) \sum_{i \in I} y_i^2 = \lambda_{\max}(A)$$

Then writing q_1, \dots, q_n for the columns of Q , we have

$$x = Q Q^\top x = Q y = \sum_{i=1}^n y_i q_i = \sum_{i \in I} y_i q_i$$

where we have used the matrix-vector product identity.

Recall that q_1, \dots, q_n are eigenvectors of A and form an orthonormal basis for \mathbb{R}^n . Therefore by construction, the set $\{q_i : i \in I\}$ forms an orthonormal basis for the eigenspace of $\lambda_{\max}(A)$. Hence x , which is a linear combination of these, lies in that eigenspace and thus is an eigenvector of A corresponding to $\lambda_{\max}(A)$.

We have shown that $\max_{\|x\|_2=1} x^\top A x = \lambda_{\max}(A)$, from which we have the general inequality $x^\top A x \leq \lambda_{\max}(A)$ for all unit-length x . \square

By the scale invariance of the Rayleigh quotient, we immediately have as a corollary (since $x^\top A x = R_A(x)$ for unit x)

Theorem 3. (*Min-max theorem*) *For all $x \neq 0$,*

$$\lambda_{\min}(A) \leq R_A(x) \leq \lambda_{\max}(A)$$

with equality if and only if x is a corresponding eigenvector.

This is sometimes referred to as a **variational characterization of eigenvalues** because it expresses the smallest/largest eigenvalue of A in terms of a minimization/maximization problem:

$$\lambda_{\min / \max}(A) = \min_{x \neq 0} / \max_{x \neq 0} R_A(x)$$

1.10 Positive (semi-)definite matrices

A symmetric matrix A is **positive semi-definite** if for all $x \in \mathbb{R}^n$, $x^\top Ax \geq 0$. Sometimes people write $A \succeq 0$ to indicate that A is positive semi-definite.

A symmetric matrix A is **positive definite** if for all nonzero $x \in \mathbb{R}^n$, $x^\top Ax > 0$. Sometimes people write $A \succ 0$ to indicate that A is positive definite. Note that positive definiteness is a strictly stronger property than positive semi-definiteness, in the sense that every positive definite matrix is positive semi-definite but not vice-versa.

These properties are related to eigenvalues in the following way.

Proposition 5. *A symmetric matrix is positive semi-definite if and only if all of its eigenvalues are nonnegative, and positive definite if and only if all of its eigenvalues are positive.*

Proof. Suppose A is positive semi-definite, and let x be an eigenvector of A with eigenvalue λ . Then

$$0 \leq x^\top Ax = x^\top (\lambda x) = \lambda x^\top x = \lambda \|x\|_2^2$$

Since $x \neq 0$ (by the assumption that it is an eigenvector), we have $\|x\|_2^2 > 0$, so we can divide both sides by $\|x\|_2^2$ to arrive at $\lambda \geq 0$. If A is positive definite, the inequality above holds strictly, so $\lambda > 0$. This proves one direction.

To simplify the proof of the other direction, we will use the machinery of Rayleigh quotients. Suppose that A is symmetric and all its eigenvalues are nonnegative. Then for all $x \neq 0$,

$$0 \leq \lambda_{\min}(A) \leq R_A(x)$$

Since $x^\top Ax$ matches $R_A(x)$ in sign, we conclude that A is positive semi-definite. If the eigenvalues of A are all strictly positive, then $0 < \lambda_{\min}(A)$, whence it follows that A is positive definite. \square

As an example of how these matrices arise, consider

Proposition 6. *Suppose $A \in \mathbb{R}^{m \times n}$. Then $A^\top A$ is positive semi-definite. If $\text{null}(A) = \{0\}$, then $A^\top A$ is positive definite.*

Proof. For any $x \in \mathbb{R}^n$,

$$x^\top (A^\top A)x = (Ax)^\top (Ax) = \|Ax\|_2^2 \geq 0$$

so $A^\top A$ is positive semi-definite. If $\text{null}(A) = \{0\}$, then $Ax \neq 0$ whenever $x \neq 0$, so $\|Ax\|_2^2 > 0$, and thus $A^\top A$ is positive definite. \square

Positive definite matrices are invertible (since their eigenvalues are nonzero), whereas positive semi-definite matrices might not be. However, if you already have a positive semi-definite matrix, it is possible to perturb its diagonal slightly to produce a positive definite matrix.

Proposition 7. *If A is positive semi-definite and $\epsilon > 0$, then $A + \epsilon I$ is positive definite.*

Proof. Assuming A is positive semi-definite and $\epsilon > 0$, we have for any $x \neq 0$ that

$$x^\top (A + \epsilon I)x = x^\top Ax + \epsilon x^\top Ix = \underbrace{x^\top Ax}_{\geq 0} + \underbrace{\epsilon \|x\|_2^2}_{> 0} > 0$$

as claimed. \square

An obvious but frequently useful consequence of the two propositions we have just shown is that $A^\top A + \epsilon I$ is positive definite (and in particular, invertible) for *any* matrix A and any $\epsilon > 0$.

1.10.1 The geometry of positive definite quadratic forms

A useful way to understand quadratic forms is by the geometry of their level sets. A **level set** or **isocontour** of a function is the set of all inputs such that the function applied to those inputs yields a given output. Mathematically, the c -isocontour of f is $\{x \in \text{dom } f : f(x) = c\}$.

Let us consider the special case $f(x) = x^\top A x$ where A is a positive definite matrix. Since A is positive definite, it has a unique matrix square root $A^{\frac{1}{2}} = Q\Lambda^{\frac{1}{2}}Q^\top$, where $Q\Lambda Q^\top$ is the eigendecomposition of A and $\Lambda^{\frac{1}{2}} = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$. It is easy to see that this matrix $A^{\frac{1}{2}}$ is positive definite (consider its eigenvalues) and satisfies $A^{\frac{1}{2}}A^{\frac{1}{2}} = A$. Fixing a value $c \geq 0$, the c -isocontour of f is the set of $x \in \mathbb{R}^n$ such that

$$c = x^\top A x = x^\top A^{\frac{1}{2}} A^{\frac{1}{2}} x = \|A^{\frac{1}{2}} x\|_2^2$$

where we have used the symmetry of $A^{\frac{1}{2}}$. Making the change of variable $z = A^{\frac{1}{2}} x$, we have the condition $\|z\|_2 = \sqrt{c}$. That is, the values z lie on a sphere of radius \sqrt{c} . These can be parameterized as $z = \sqrt{c}\hat{z}$ where \hat{z} has $\|\hat{z}\|_2 = 1$. Then since $A^{-\frac{1}{2}} = Q\Lambda^{-\frac{1}{2}}Q^\top$, we have

$$x = A^{-\frac{1}{2}} z = Q\Lambda^{-\frac{1}{2}}Q^\top \sqrt{c}\hat{z} = \sqrt{c}Q\Lambda^{-\frac{1}{2}}\tilde{z}$$

where $\tilde{z} = Q^\top \hat{z}$ also satisfies $\|\tilde{z}\|_2 = 1$ since Q is orthogonal. Using this parameterization, we see that the solution set $\{x \in \mathbb{R}^n : f(x) = c\}$ is the image of the unit sphere $\{\tilde{z} \in \mathbb{R}^n : \|\tilde{z}\|_2 = 1\}$ under the invertible linear map $x = \sqrt{c}Q\Lambda^{-\frac{1}{2}}\tilde{z}$.

What we have gained with all these manipulations is a clear algebraic understanding of the c -isocontour of f in terms of a sequence of linear transformations applied to a well-understood set. We begin with the unit sphere, then scale every axis i by $\lambda_i^{-\frac{1}{2}}$, resulting in an axis-aligned ellipsoid. Observe that the axis lengths of the ellipsoid are proportional to the inverse square roots of the eigenvalues of A . Hence larger eigenvalues correspond to shorter axis lengths, and vice-versa.

Then this axis-aligned ellipsoid undergoes a rigid transformation (i.e. one that preserves length and angles, such as a rotation/reflection) given by Q . The result of this transformation is that the axes of the ellipse are no longer along the coordinate axes in general, but rather along the directions given by the corresponding eigenvectors. To see this, consider the unit vector $e_i \in \mathbb{R}^n$ that has $(e_i)_j = \delta_{ij}$. In the pre-transformed space, this vector points along the axis with length proportional to $\lambda_i^{-\frac{1}{2}}$. But after applying the rigid transformation Q , the resulting vector points in the direction of the corresponding eigenvector q_i , since

$$Qe_i = \sum_{j=1}^n (e_i)_j q_j = q_i$$

where we have used the matrix-vector product identity from earlier.

In summary: the isocontours of $f(x) = x^\top A x$ are ellipsoids such that the axes point in the directions of the eigenvectors of A , and the radii of these axes are proportional to the inverse square roots of the corresponding eigenvalues.

1.11 Singular value decomposition

Singular value decomposition (SVD) is a widely applicable tool in linear algebra. Its strength stems partially from the fact that *every matrix* $A \in \mathbb{R}^{m \times n}$ has an SVD (even non-square matrices)! The decomposition goes as follows:

$$A = U\Sigma V^\top$$

where $U \in \mathbb{R}^{m \times m}$ and $V \in \mathbb{R}^{n \times n}$ are orthogonal matrices and $\Sigma \in \mathbb{R}^{m \times n}$ is a diagonal matrix with the **singular values** of A (denoted σ_i) on its diagonal.

Only the first $r = \text{rank}(A)$ singular values are nonzero, and by convention, they are in non-increasing order, i.e.

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > \sigma_{r+1} = \cdots = \sigma_{\min(m,n)} = 0$$

Another way to write the SVD (cf. the sum-of-outer-products identity) is

$$A = \sum_{i=1}^r \sigma_i u_i v_i^\top$$

where u_i and v_i are the i th columns of U and V , respectively.

Observe that the SVD factors provide eigendecompositions for $A^\top A$ and AA^\top :

$$\begin{aligned} A^\top A &= (U \Sigma V^\top)^\top U \Sigma V^\top = V \Sigma^\top U^\top U \Sigma V^\top = V \Sigma^\top \Sigma V^\top \\ AA^\top &= U \Sigma V^\top (U \Sigma V^\top)^\top = U \Sigma V^\top V \Sigma^\top U^\top = U \Sigma \Sigma^\top U^\top \end{aligned}$$

It follows immediately that the columns of V (the **right-singular vectors** of A) are eigenvectors of $A^\top A$, and the columns of U (the **left-singular vectors** of A) are eigenvectors of AA^\top .

The matrices $\Sigma^\top \Sigma$ and $\Sigma \Sigma^\top$ are not necessarily the same size, but both are diagonal with the squared singular values σ_i^2 on the diagonal (plus possibly some zeros). Thus the singular values of A are the square roots of the eigenvalues of $A^\top A$ (or equivalently, of AA^\top)².

1.12 Fundamental Theorem of Linear Algebra

Despite its fancy name, the “Fundamental Theorem of Linear Algebra” is not a universally-agreed-upon theorem; there is some ambiguity as to exactly what statements it includes. The version we present here is sufficient for our purposes.

Theorem 4. *If $A \in \mathbb{R}^{m \times n}$, then*

- (i) $\text{null}(A) = \text{range}(A^\top)^\perp$
- (ii) $\text{null}(A) \oplus \text{range}(A^\top) = \mathbb{R}^n$
- (iii) $\underbrace{\dim \text{range}(A)}_{\text{rank}(A)} + \dim \text{null}(A) = n$.³
- (iv) *If $A = U \Sigma V^\top$ is the singular value decomposition of A , then the columns of U and V form orthonormal bases for the four “fundamental subspaces” of A :*

Subspace	Columns
$\text{range}(A)$	The first r columns of U
$\text{range}(A^\top)$	The first r columns of V
$\text{null}(A^\top)$	The last $m - r$ columns of U
$\text{null}(A)$	The last $n - r$ columns of V

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

² Recall that $A^\top A$ and AA^\top are positive semi-definite, so their eigenvalues are nonnegative, and thus taking square roots is always well-defined.

³ This result is sometimes referred to by itself as the **rank-nullity theorem**.

Proof. (i) Let a_1, \dots, a_m denote the rows of A . Then

$$\begin{aligned}
x \in \text{null}(A) &\iff Ax = 0 \\
&\iff a_i^\top x = 0 \text{ for all } i = 1, \dots, m \\
&\iff (\alpha_1 a_1 + \dots + \alpha_m a_m)^\top x = 0 \text{ for all } \alpha_1, \dots, \alpha_m \\
&\iff v^\top x = 0 \text{ for all } v \in \text{range}(A^\top) \\
&\iff x \in \text{range}(A^\top)^\perp
\end{aligned}$$

which proves the result.

- (ii) Recall our previous result on orthogonal complements: if S is a finite-dimensional subspace of V , then $V = S \oplus S^\perp$. Thus the claim follows from the previous part (take $V = \mathbb{R}^n$ and $S = \text{range}(A^\top)$).
- (iii) Recall that if U and W are subspaces of a finite-dimensional vector space V , then $\dim(U \oplus W) = \dim U + \dim W$. Thus the claim follows from the previous part, using the fact that $\dim \text{range}(A) = \dim \text{range}(A^\top)$.

□

A direct result of (ii) is that every $x \in \mathbb{R}^n$ can be written (uniquely) in the form

$$x = A^\top v + w$$

for some $v \in \mathbb{R}^m, w \in \mathbb{R}^n$, where $Aw = 0$.

Note that there is some asymmetry in the theorem, but analogous statements can be obtained by applying the theorem to A^\top .

1.13 Operator and matrix norms

If V and W are vector spaces, then the set of linear maps from V to W forms another vector space, and the norms defined on V and W induce a norm on this space of linear maps. If $T : V \rightarrow W$ is a linear map between normed spaces V and W , then the **operator norm** is defined as

$$\|T\|_{\text{op}} = \max_{\substack{v \in V \\ v \neq 0}} \frac{\|Tv\|_W}{\|v\|_V}$$

An important class of this general definition is when the domain and codomain are \mathbb{R}^n and \mathbb{R}^m , and the p -norm is used in both cases. Then for a matrix $A \in \mathbb{R}^{m \times n}$, we can define the matrix p -norm

$$\|A\|_p = \max_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p}$$

In the special cases $p = 1, 2, \infty$ we have

$$\begin{aligned}
\|A\|_1 &= \max_{1 \leq j \leq n} \sum_{i=1}^m |A_{ij}| \\
\|A\|_\infty &= \max_{1 \leq i \leq m} \sum_{j=1}^n |A_{ij}| \\
\|A\|_2 &= \sigma_1(A)
\end{aligned}$$

where σ_1 denotes the largest singular value. Note that the induced 1- and ∞ -norms are simply the maximum absolute column and row sums, respectively. The induced 2-norm (often called the **spectral norm**) simplifies to σ_1 by the properties of Rayleigh quotients proved earlier; clearly

$$\arg \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \arg \max_{x \neq 0} \frac{\|Ax\|_2^2}{\|x\|_2^2} = \arg \max_{x \neq 0} \frac{x^\top A^\top A x}{x^\top x}$$

and we have seen that the rightmost expression is maximized by an eigenvector of $A^\top A$ corresponding to its largest eigenvalue, $\lambda_{\max}(A^\top A) = \sigma_1^2(A)$.

By definition, these induced matrix norms have the important property that

$$\|Ax\|_p \leq \|A\|_p \|x\|_p$$

for any x . They are also **submultiplicative** in the following sense.

Proposition 8. $\|AB\|_p \leq \|A\|_p \|B\|_p$

Proof. For any x ,

$$\|ABx\|_p \leq \|A\|_p \|Bx\|_p \leq \|A\|_p \|B\|_p \|x\|_p$$

so

$$\|AB\|_p = \max_{x \neq 0} \frac{\|ABx\|_p}{\|x\|_p} \leq \max_{x \neq 0} \frac{\|A\|_p \|B\|_p \|x\|_p}{\|x\|_p} = \|A\|_p \|B\|_p$$

□

These are not the only matrix norms, however. Another frequently used is the **Frobenius norm**

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2} = \sqrt{\text{tr}(A^\top A)} = \sqrt{\sum_{i=1}^{\min(m,n)} \sigma_i^2(A)}$$

The first equivalence follows straightforwardly by expanding the definitions of matrix multiplication and trace. For the second, observe that (writing $A = U\Sigma V^\top$ as before)

$$\text{tr}(A^\top A) = \text{tr}(V\Sigma^\top \Sigma V^\top) = \text{tr}(V^\top V \Sigma^\top \Sigma) = \text{tr}(\Sigma^\top \Sigma) = \sum_{i=1}^{\min(m,n)} \sigma_i^2(A)$$

using the cyclic property of trace and orthogonality of V .

A matrix norm $\|\cdot\|$ is said to be **unitary invariant** if

$$\|UAV\| = \|A\|$$

for all orthogonal U and V of appropriate size. Unitary invariant norms essentially depend only on the singular values of a matrix, since for such norms,

$$\|A\| = \|U\Sigma V^\top\| = \|\Sigma\|$$

Two particular norms we have seen, the spectral norm and the Frobenius norm, can be expressed solely in terms of a matrix's singular values.

Proposition 9. *The spectral norm and the Frobenius norm are unitary invariant.*

Proof. For the Frobenius norm, the claim follows from

$$\text{tr}((UAV)^\top UAV) = \text{tr}(V^\top A^\top U^\top UAV) = \text{tr}(VV^\top A^\top A) = \text{tr}(A^\top A)$$

For the spectral norm, recall that $\|Ux\|_2 = \|x\|_2$ for any orthogonal U . Thus

$$\|UAV\|_2 = \max_{x \neq 0} \frac{\|UAVx\|_2}{\|x\|_2} = \max_{x \neq 0} \frac{\|AVx\|_2}{\|x\|_2} = \max_{y \neq 0} \frac{\|Ay\|_2}{\|y\|_2} = \|A\|_2$$

where we have used the change of variable $y = V^\top x$, which satisfies $\|y\|_2 = \|x\|_2$. Since V^\top is invertible, x and y are in one-to-one correspondence, and in particular $y = 0$ if and only if $x = 0$. Hence maximizing over $y \neq 0$ is equivalent to maximizing over $x \neq 0$. \square

1.14 Low-rank approximation

An important practical application of the SVD is to compute **low-rank approximations** to matrices. That is, given some matrix, we want to find another matrix of the same dimensions but lower rank such that the two matrices are close as measured by some norm. Such an approximation can be used to reduce the amount of data needed to store a matrix, while retaining most of its information.

A remarkable result known as the **Eckart-Young-Mirsky theorem** tells us that the optimal matrix can be computed easily from the SVD, as long as the norm in question is unitary invariant (e.g., the spectral norm or Frobenius norm).

Theorem 5. (*Eckart-Young-Mirsky*) Let $\|\cdot\|$ be a unitary invariant matrix norm. Suppose $A \in \mathbb{R}^{m \times n}$, where $m \geq n$, has singular value decomposition $A = \sum_{i=1}^n \sigma_i u_i v_i^\top$. Then the best rank- k approximation to A , where $k \leq \text{rank}(A)$, is given by

$$A_k = \sum_{i=1}^k \sigma_i u_i v_i^\top$$

in the sense that

$$\|A - A_k\| \leq \|A - \tilde{A}\|$$

for any $\tilde{A} \in \mathbb{R}^{m \times n}$ with $\text{rank}(\tilde{A}) \leq k$.

The proof of the general case requires a fair amount of work, so we prove only the special case where $\|\cdot\|$ is the spectral norm.

Proof. First we compute

$$\|A - A_k\|_2 = \left\| \sum_{i=1}^n \sigma_i u_i v_i^\top - \sum_{i=1}^k \sigma_i u_i v_i^\top \right\|_2 = \left\| \sum_{i=k+1}^n \sigma_i u_i v_i^\top \right\|_2 = \sigma_{k+1}$$

Let $\tilde{A} \in \mathbb{R}^{m \times n}$ have $\text{rank}(\tilde{A}) \leq k$. Then by the Fundamental Theorem of Linear Algebra,

$$\dim \text{null}(\tilde{A}) = n - \text{rank}(\tilde{A}) \geq n - k$$

It follows that

$$\text{null}(\tilde{A}) \cap \text{span}\{v_1, \dots, v_{k+1}\}$$

is non-trivial (has a nonzero element), because otherwise there would be at least $(n - k) + (k + 1) = n + 1$ linearly independent vectors in \mathbb{R}^n , which is impossible. Therefore let z be some element of

the intersection, and assume without loss of generality that it has unit norm: $\|z\|_2 = 1$. Expand $z = \alpha_1 v_1 + \cdots + \alpha_{k+1} v_{k+1}$, noting that

$$1 = \|z\|_2^2 = \|\alpha_1 v_1 + \cdots + \alpha_{k+1} v_{k+1}\|_2^2 = \alpha_1^2 + \cdots + \alpha_{k+1}^2$$

by the Pythagorean theorem. Thus

$$\begin{aligned} \|A - \tilde{A}\|_2 &\geq \|(A - \tilde{A})z\|_2 && \text{by def., and } \|z\|_2 = 1 \\ &= \|Az\|_2 && z \in \text{null}(\tilde{A}) \\ &= \left\| \sum_{i=1}^n \sigma_i u_i v_i^\top z \right\|_2 \\ &= \left\| \sum_{i=1}^{k+1} \sigma_i \alpha_i u_i \right\|_2 \\ &= \sqrt{(\sigma_1 \alpha_1)^2 + \cdots + (\sigma_{k+1} \alpha_{k+1})^2} && \text{Pythagorean theorem again} \\ &\geq \sigma_{k+1} \sqrt{\alpha_1^2 + \cdots + \alpha_{k+1}^2} && \sigma_{k+1} \leq \sigma_i \text{ for } i \leq k \\ &= \|A - A_k\|_2 && \text{using our earlier results} \end{aligned}$$

as was to be shown. \square

A measure of the quality of the approximation is given by

$$\frac{\|A_k\|_F^2}{\|A\|_F^2} = \frac{\sigma_1^2 + \cdots + \sigma_k^2}{\sigma_1^2 + \cdots + \sigma_r^2}$$

Ideally, this ratio will be close to 1, indicating that most of the information was retained.

1.15 Pseudoinverses

Let $A \in \mathbb{R}^{m \times n}$. If $m \neq n$, then A cannot possibly be invertible. However, there is a generalization of the inverse known as the **Moore-Penrose pseudoinverse**, denoted $A^\dagger \in \mathbb{R}^{n \times m}$, which always exists and is defined uniquely by the following properties:

- (i) $AA^\dagger A = A$
- (ii) $A^\dagger AA^\dagger = A^\dagger$
- (iii) AA^\dagger is symmetric
- (iv) $A^\dagger A$ is symmetric

If A is invertible, then $A^\dagger = A^{-1}$. More generally, we can compute the pseudoinverse of a matrix from its singular value decomposition: if $A = U\Sigma V^\top$, then

$$A^\dagger = V\Sigma^\dagger U^\top$$

where Σ^\dagger can be computed from Σ by taking the transpose and inverting the nonzero singular values on the diagonal. Verifying that this matrix satisfies the properties of the pseudoinverse is straightforward and left as an exercise to the reader.

1.16 Some useful matrix identities

1.16.1 Matrix-vector product as linear combination of matrix columns

Proposition 10. *Let $x \in \mathbb{R}^n$ be a vector and $A \in \mathbb{R}^{m \times n}$ a matrix with columns a_1, \dots, a_n . Then*

$$Ax = \sum_{i=1}^n x_i a_i$$

This identity is extremely useful in understanding linear operators in terms of their matrices' columns. The proof is very simple (consider each element of Ax individually and expand by definitions) but it is a good exercise to convince yourself.

1.16.2 Sum of outer products as matrix-matrix product

An **outer product** is an expression of the form ab^\top , where $a \in \mathbb{R}^m$ and $b \in \mathbb{R}^n$. By inspection it is not hard to see that such an expression yields an $m \times n$ matrix such that

$$(ab^\top)_{ij} = a_i b_j$$

It is not immediately obvious, but the sum of outer products is actually equivalent to an appropriate matrix-matrix product! We formalize this statement as

Proposition 11. *Let $a_1, \dots, a_k \in \mathbb{R}^m$ and $b_1, \dots, b_k \in \mathbb{R}^n$. Then*

$$\sum_{\ell=1}^k a_\ell b_\ell^\top = AB^\top$$

where

$$A = [a_1 \quad \dots \quad a_k], \quad B = [b_1 \quad \dots \quad b_k]$$

Proof. For each (i, j) , we have

$$\left(\sum_{\ell=1}^k a_\ell b_\ell^\top \right)_{ij} = \sum_{\ell=1}^k (a_\ell b_\ell^\top)_{ij} = \sum_{\ell=1}^k (a_\ell)_i (b_\ell)_j = \sum_{\ell=1}^k A_{i\ell} B_{j\ell}$$

This last expression should be recognized as an inner product between the i th row of A and the j th row of B , or equivalently the j th column of B^\top . Hence by the definition of matrix multiplication, it is equal to $(AB^\top)_{ij}$. \square